

InnoDB Plugin 1.0 for MySQL 5.1 User's Guide

InnoDB Plugin 1.0 for MySQL 5.1 User's Guide

Abstract

This is the User's Guide for InnoDB Plugin 1.0.8 for MySQL 5.1.

Starting with version 5.1, MySQL AB has promoted the idea of a “pluggable” [storage engine architecture](#), which permits multiple storage engines to be added to MySQL. Beginning with MySQL version 5.1, it is possible for users to swap out one version of InnoDB and use another. The pluggable storage engine architecture also permits Innobase Oy to release new versions of InnoDB containing bug fixes and new features independently of the release cycle for MySQL.

This User's Guide documents the installation and removal procedures and the additional features of the InnoDB Plugin 1.0.8 for MySQL 5.1.

WARNING: Because the InnoDB Plugin introduces a new file format, restrictions apply to the use of a database created with the InnoDB Plugin with earlier versions of InnoDB, when using `mysqldump` or MySQL replication and if you use the InnoDB Hot Backup utility. See [Section 1.5, “Operational Restrictions”](#).

For legal information, see the [Legal Notices](#).

Document generated on: 2014-01-30 (revision: 37573)

Table of Contents

Preface and Legal Notices	vii
1 Introduction to the InnoDB Plugin	1
1.1 Overview	1
1.2 Features of the InnoDB Plugin	1
1.3 Obtaining and Installing the InnoDB Plugin	3
1.4 Viewing the InnoDB Plugin Version Number	3
1.5 Operational Restrictions	4
2 Fast Index Creation in the InnoDB Storage Engine	5
2.1 Overview of Fast Index Creation	5
2.2 Examples	5
2.3 Implementation	6
2.4 Concurrency Considerations	7
2.5 Crash Recovery	7
2.6 Limitations	8
3 InnoDB Data Compression	9
3.1 Overview of Table Compression	9
3.2 Enabling Compression for a Table	9
3.2.1 Configuration Parameters for Compression	10
3.2.2 SQL Compression Syntax Warnings and Errors	11
3.3 Tuning InnoDB Compression	12
3.3.1 When to Use Compression	13
3.3.2 Monitoring Compression at Runtime	15
3.4 How Compression Works in InnoDB	16
3.4.1 Compression Algorithms	16
3.4.2 InnoDB Data Storage and Compression	16
3.4.3 Compression and the InnoDB Buffer Pool	18
3.4.4 Compression and the InnoDB Log Files	18
4 InnoDB File-Format Management	19
4.1 Overview of InnoDB File Formats	19
4.2 Named File Formats	19
4.3 Enabling File Formats	20
4.4 File Format Compatibility	20
4.4.1 Startup File Format Compatibility Checking	21
4.4.2 Table-Access File Format Compatibility Checking	22
4.5 Identifying the File Format in Use	23
4.6 Downgrading the File Format	24
4.7 Future InnoDB File Formats	24
5 InnoDB Row Storage and Row Formats	25
5.1 Storage of Variable-Length Columns	25
5.2 COMPACT and REDUNDANT Row Formats	25
5.3 DYNAMIC Row Format	25
5.4 Specifying a Table's Row Format	26
6 InnoDB INFORMATION_SCHEMA Tables	27
6.1 Overview of InnoDB Support in INFORMATION_SCHEMA	27
6.2 Information Schema Tables about Compression	27
6.2.1 INNODB_CMP and INNODB_CMP_RESET	27
6.2.2 INNODB_CMPMEM and INNODB_CMPMEM_RESET	28
6.2.3 Using the Compression Information Schema Tables	29
6.3 Information Schema Tables about Transactions	29
6.3.1 INNODB_TRX	29
6.3.2 INNODB_LOCKS	30

6.3.3	<code>INNODB_LOCK_WAITS</code>	31
6.3.4	Using the Transaction Information Schema Tables	31
6.4	Notes on Locking in InnoDB	36
6.4.1	Understanding InnoDB Locking	36
6.4.2	Rapidly Changing Internal Data	36
6.4.3	Possible Inconsistency with <code>PROCESSLIST</code>	37
7	Performance and Scalability Enhancements	39
7.1	Overview	39
7.2	Faster Locking for Improved Scalability	40
7.3	Using Operating System Memory Allocators	40
7.4	Controlling InnoDB Insert Buffering	41
7.5	Controlling Adaptive Hash Indexing	42
7.6	Changes Regarding Thread Concurrency	42
7.7	Changes in the Read Ahead Algorithm	43
7.8	Multiple Background I/O Threads	44
7.9	Group Commit	44
7.10	Controlling the Master Thread I/O Rate	45
7.11	Controlling the Flushing Rate of Dirty Pages	45
7.12	Using the <code>PAUSE</code> instruction in InnoDB spin loops	46
7.13	Control of Spin Lock Polling	46
7.14	Making Buffer Cache Scan Resistant	46
7.14.1	Guidelines for <code>innodb_old_blocks_pct</code> and <code>innodb_old_blocks_time</code>	48
7.15	Improvements to Crash Recovery Performance	48
8	Changes for Flexibility, Ease of Use and Reliability	49
8.1	Overview	49
8.2	Enabling New File Formats	49
8.3	Dynamic Control of System Configuration Parameters	50
8.3.1	Dynamically Changing <code>innodb_file_per_table</code>	50
8.3.2	Dynamically Changing <code>innodb_stats_on_metadata</code>	50
8.3.3	Dynamically Changing <code>innodb_lock_wait_timeout</code>	51
8.3.4	Dynamically Changing <code>innodb_adaptive_hash_index</code>	51
8.4	<code>TRUNCATE TABLE</code> Reclaims Space	51
8.5	InnoDB Strict Mode	52
8.6	Controlling Optimizer Statistics Estimation	52
8.7	Better Error Handling when Dropping Indexes	53
8.8	More Compact Output of <code>SHOW ENGINE INNODB MUTEX</code>	54
8.9	More Read Ahead Statistics	54
9	Installing the InnoDB Plugin	57
9.1	Overview of Installing the InnoDB Plugin	57
9.2	Checking for Compatible Version Levels	58
9.3	Installing the Precompiled InnoDB Plugin as a Shared Library	58
9.3.1	Installing the InnoDB Plugin as a Shared Library on Unix or Linux	59
9.3.2	Installing the Binary InnoDB Plugin as a Shared Library on Microsoft Windows	62
9.3.3	Errors When Installing the InnoDB Plugin on Microsoft Windows	64
9.4	Building the InnoDB Plugin from Source Code	65
9.4.1	Building the InnoDB Plugin on Linux or Unix	66
9.4.2	Building the InnoDB Plugin on Microsoft Windows	67
9.5	Configuring the InnoDB Plugin	68
9.6	Frequently Asked Questions about Plugin Installation	69
9.6.1	Should I use the InnoDB-supplied plugin or the one that is included with MySQL 5.1.38 or higher?	69
9.6.2	Why doesn't the MySQL service on Windows start after the replacement?	69
9.6.3	The Plugin is installed... now what?	69
9.6.4	Once the Plugin is installed, is it permanent?	69

10 Upgrading the InnoDB Plugin	71
10.1 Upgrading the Dynamic InnoDB Plugin	71
10.2 Upgrading a Statically Built InnoDB Plugin	71
10.3 Converting Compressed Tables Created Before Version 1.0.2	72
11 Downgrading from the InnoDB Plugin	73
11.1 Overview	73
11.2 The Built-in InnoDB, the Plugin and File Formats	73
11.3 How to Downgrade	74
11.3.1 Converting Tables	74
11.3.2 Adjusting the Configuration	74
11.3.3 Uninstalling a Dynamic Library	74
11.3.4 Uninstalling a Statically Built InnoDB Plugin	75
11.4 Possible Problems	75
11.4.1 Accessing <code>COMPRESSED</code> or <code>DYNAMIC</code> Tables	75
11.4.2 Issues with UNDO and REDO	76
11.4.3 Issues with the Doublewrite Buffer	76
11.4.4 Issues with the Insert Buffer	77
12 InnoDB Plugin Change History	79
12.1 Changes in InnoDB Plugin 1.0.9 and Higher	79
12.2 Changes in InnoDB Plugin 1.0.8 (May, 2010)	79
12.3 Changes in InnoDB Plugin 1.0.7 (April, 2010)	79
12.4 Changes in InnoDB Plugin 1.0.6 (November 27, 2009)	80
12.5 Changes in InnoDB Plugin 1.0.5 (November 18, 2009)	80
12.6 Changes in InnoDB Plugin 1.0.4 (August 11, 2009)	81
12.7 Changes in InnoDB Plugin 1.0.3 (March 11, 2009)	82
12.8 Changes in InnoDB Plugin 1.0.2 (December 1, 2008)	83
12.9 Changes in InnoDB Plugin 1.0.1 (May 8, 2008)	83
12.10 Changes in InnoDB Plugin 1.0.0 (April 15, 2008)	84
A Third-Party Software	85
A.1 Performance Patches from Google	85
A.2 Multiple Background I/O Threads Patch from Percona	86
A.3 Performance Patches from Sun Microsystems	86
B Using the InnoDB Plugin with MySQL 5.1.30 or Earlier	89
C List of Parameters Changed in the InnoDB Plugin 1.0	91
C.1 New Parameters	91
C.2 Deprecated Parameters	93
C.3 Parameters with New Defaults	93
InnoDB Glossary	95
Index	151

Preface and Legal Notices

This is the User's Guide for the InnoDB Plugin 1.0.8 for MySQL 5.1.

Legal Notices

Copyright © 1997, 2014, Oracle and/or its affiliates. All rights reserved.

This software and related documentation are provided under a license agreement containing restrictions on use and disclosure and are protected by intellectual property laws. Except as expressly permitted in your license agreement or allowed by law, you may not use, copy, reproduce, translate, broadcast, modify, license, transmit, distribute, exhibit, perform, publish, or display any part, in any form, or by any means. Reverse engineering, disassembly, or decompilation of this software, unless required by law for interoperability, is prohibited.

The information contained herein is subject to change without notice and is not warranted to be error-free. If you find any errors, please report them to us in writing.

If this software or related documentation is delivered to the U.S. Government or anyone licensing it on behalf of the U.S. Government, the following notice is applicable:

U.S. GOVERNMENT RIGHTS Programs, software, databases, and related documentation and technical data delivered to U.S. Government customers are "commercial computer software" or "commercial technical data" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, the use, duplication, disclosure, modification, and adaptation shall be subject to the restrictions and license terms set forth in the applicable Government contract, and, to the extent applicable by the terms of the Government contract, the additional rights set forth in FAR 52.227-19, Commercial Computer Software License (December 2007). Oracle USA, Inc., 500 Oracle Parkway, Redwood City, CA 94065.

This software is developed for general use in a variety of information management applications. It is not developed or intended for use in any inherently dangerous applications, including applications which may create a risk of personal injury. If you use this software in dangerous applications, then you shall be responsible to take all appropriate fail-safe, backup, redundancy, and other measures to ensure the safe use of this software. Oracle Corporation and its affiliates disclaim any liability for any damages caused by use of this software in dangerous applications.

Oracle is a registered trademark of Oracle Corporation and/or its affiliates. MySQL is a trademark of Oracle Corporation and/or its affiliates, and shall not be used without Oracle's express written authorization. Other names may be trademarks of their respective owners.

This software and documentation may provide access to or information on content, products, and services from third parties. Oracle Corporation and its affiliates are not responsible for and expressly disclaim all warranties of any kind with respect to third-party content, products, and services. Oracle Corporation and its affiliates will not be responsible for any loss, costs, or damages incurred due to your access to or use of third-party content, products, or services.

This document in any form, software or printed matter, contains proprietary information that is the exclusive property of Oracle. Your access to and use of this material is subject to the terms and conditions of your Oracle Software License and Service Agreement, which has been executed and with which you agree to comply. This document and information contained herein may not be disclosed, copied, reproduced, or distributed to anyone outside Oracle without prior written consent of Oracle or as specifically provided below. This document is not part of your license agreement nor can it be incorporated into any contractual agreement with Oracle or its subsidiaries or affiliates.

This documentation is NOT distributed under a GPL license. Use of this documentation is subject to the following terms:

You may create a printed copy of this documentation solely for your own personal use. Conversion to other formats is allowed as long as the actual content is not altered or edited in any way. You shall not publish or distribute this documentation in any form or on any media, except if you distribute the documentation in a manner similar to how Oracle disseminates it (that is, electronically for download on a Web site with the software) or on a CD-ROM or similar medium, provided however that the documentation is disseminated together with the software on the same medium. Any other use, such as any dissemination of printed copies or use of this documentation, in whole or in part, in another publication, requires the prior written consent from an authorized representative of Oracle. Oracle and/or its affiliates reserve any and all rights to this documentation not expressly granted above.

For more information on the terms of this license, or for details on how the MySQL documentation is built and produced, please visit [MySQL Contact & Questions](#).

For additional licensing information, including licenses for third-party libraries used by MySQL products, see [Preface and Legal Notices](#).

For help with using MySQL, please visit either the [MySQL Forums](#) or [MySQL Mailing Lists](#) where you can discuss your issues with other MySQL users.

For additional documentation on MySQL products, including translations of the documentation into other languages, and downloadable versions in variety of formats, including HTML and PDF formats, see the [MySQL Documentation Library](#).

Chapter 1 Introduction to the InnoDB Plugin

Table of Contents

1.1 Overview	1
1.2 Features of the InnoDB Plugin	1
1.3 Obtaining and Installing the InnoDB Plugin	3
1.4 Viewing the InnoDB Plugin Version Number	3
1.5 Operational Restrictions	4

1.1 Overview

The unique architecture of MySQL permits multiple storage engines with different capabilities to be accessed through the same SQL language and APIs. Starting with version 5.1, MySQL AB has promoted the idea of a “[pluggable storage engine architecture](#)”, which permits multiple storage engines to be added to MySQL. Currently, however, most users have accessed only those storage engines that are distributed by MySQL AB, and are linked into the binary (executable) releases.

Since 2001, MySQL AB has distributed the InnoDB transactional storage engine with its releases (both source and binary). Beginning with MySQL version 5.1, it is possible for users to swap out one version of InnoDB and use another. The pluggable storage engine architecture also permits Innobase Oy to release new versions of InnoDB containing bug fixes and new features independently of the release cycle for MySQL. Users can thus take advantage of these new versions of InnoDB in the context of their deployed MySQL installations.

1.2 Features of the InnoDB Plugin

The InnoDB Plugin for MySQL contains several important new features:

- [Viewing the InnoDB Plugin version number](#)
- [Fast index creation: add or drop indexes without copying the data](#)
- [Data compression: shrink tables, to significantly reduce storage and I/O](#)
- [New row format: fully off-page storage of long `BLOB`, `TEXT`, and `VARCHAR` columns](#)
- [File format management: protects upward and downward compatibility](#)
- [`INFORMATION_SCHEMA` tables: information about compression and locking](#)
- [Performance and scalability enhancements:](#)
 - [Section 7.2, “Faster Locking for Improved Scalability”](#)
 - [Section 7.3, “Using Operating System Memory Allocators”](#)
 - [Section 7.4, “Controlling InnoDB Insert Buffering”](#)
 - [Section 7.5, “Controlling Adaptive Hash Indexing”](#)
 - [Section 7.6, “Changes Regarding Thread Concurrency”](#)
 - [Section 7.7, “Changes in the Read Ahead Algorithm”](#)

- [Section 7.8, “Multiple Background I/O Threads”](#)
- [Section 7.9, “Group Commit”](#)
- [Section 7.10, “Controlling the Master Thread I/O Rate”](#)
- [Section 7.11, “Controlling the Flushing Rate of Dirty Pages”](#)
- [Section 7.12, “Using the PAUSE instruction in InnoDB spin loops”](#)
- [Section 7.13, “Control of Spin Lock Polling”](#)
- [Section C.3, “Parameters with New Defaults”](#)
- [Section 7.14, “Making Buffer Cache Scan Resistant”](#)
- [Section 7.15, “Improvements to Crash Recovery Performance”](#)
- Other changes for flexibility, ease of use and reliability:
 - [Section 8.2, “Enabling New File Formats”](#)
 - [Section 8.3, “Dynamic Control of System Configuration Parameters”](#)
 - [Section 8.4, “TRUNCATE TABLE Reclaims Space”](#)
 - [Section 8.5, “InnoDB Strict Mode”](#)
 - [Section 8.6, “Controlling Optimizer Statistics Estimation”](#)
 - [Section 8.7, “Better Error Handling when Dropping Indexes”](#)
 - [Section 8.8, “More Compact Output of SHOW ENGINE INNODB MUTEX”](#)
 - [Section 8.9, “More Read Ahead Statistics”](#)

Note that the ability to use data compression and the new row format require the use of a new InnoDB file format called “Barracuda”. The previous file format, used by the built-in InnoDB in MySQL versions 5.0 and 5.1 is now called “Antelope” and does not support these features, but does support the other features introduced with the InnoDB Plugin.

The InnoDB Plugin is upward compatible from standard InnoDB as built in to, and distributed with, MySQL. Existing databases can be used with the InnoDB Plugin for MySQL. As described in [Section 9.5, “Configuring the InnoDB Plugin”](#), the new parameter `innodb_file_format` can help protect upward and downward compatibility between InnoDB versions and database files, allowing users to enable or disable use of new features that can only be used with certain versions of InnoDB.

The built-in InnoDB in MySQL since version 5.0.21 has a safety feature that prevents it from opening tables that are in an unknown format. However, as noted in [Section 11.2, “The Built-in InnoDB, the Plugin and File Formats”](#), the system tablespace may contain references to new-format tables that will confuse the built-in InnoDB in MySQL. These references will be cleared in a “slow” shutdown of the InnoDB Plugin.

With previous versions of InnoDB, no error would be returned until you try to access a table that is in a format “too new” for the software. Beginning with version 1.0.1 of the InnoDB Plugin, however, to provide early feedback, InnoDB will check the system tablespace to ensure that the file format used in the database is enabled for use before it will start. See [Section 4.4.1, “Startup File Format Compatibility Checking”](#) for the details.

1.3 Obtaining and Installing the InnoDB Plugin

From the [MySQL download site](#), you can download the MySQL Server 5.1.38 and up, containing an additional, upgraded version of InnoDB with additional features beyond the built-in InnoDB provided by MySQL. You can configure your server to use the InnoDB Plugin rather than the built-in InnoDB, to take advantage of the extra features and performance enhancements in the Plugin. The InnoDB Plugin for MySQL is licensed under the same license that MySQL uses (GPLv2). It is available at no charge for use and can be freely redistributed, subject to the same requirements of the GPL as is MySQL itself (see the [GNU General Public License, version 2](#)). Support is available for the InnoDB Plugin for MySQL through the normal MySQL support channel, the [MySQL bug database](#).

In many environments, the InnoDB plugin can dynamically be added to a MySQL instance without relinking the MySQL server. The plugin version of InnoDB then “takes over” from the statically linked InnoDB that is part of the `mysqld` binary. In other environments, it is currently necessary to build the entire MySQL server, including the InnoDB plugin, from source code.

On Linux, Unix and Windows, it is a simple matter of installing the InnoDB Plugin for MySQL using the `INSTALL PLUGIN` statement. When the InnoDB Plugin for MySQL is installed, it replaces the statically-linked version of InnoDB that is incorporated in the MySQL binary as distributed by MySQL AB.

On platforms where the dynamic plugin is not available, users must download the source code for the InnoDB Plugin for MySQL and build MySQL from source. Building from source also facilitates the distribution of a MySQL server where the InnoDB Plugin is already “installed”, so all users of an organization can use the new capabilities without the `INSTALL` step. The procedure for building from source code is documented in [Section 9.4, “Building the InnoDB Plugin from Source Code”](#).

Full instructions are provided in [Chapter 9, *Installing the InnoDB Plugin*](#).

1.4 Viewing the InnoDB Plugin Version Number

InnoDB Plugin releases are numbered with version numbers independent of MySQL release numbers. The initial release of the InnoDB Plugin is version 1.0, and it is designed to work with MySQL 5.1.

- The first component of the InnoDB Plugin version number designates a major release level.
- The second component corresponds to the MySQL release. The digit 0 corresponds to MySQL 5.1.
- The third component indicates the specific release of the InnoDB Plugin (at a given major release level and for a specific MySQL release); only bug fixes and minor functional changes are introduced at this level.

Once you have installed the InnoDB Plugin, you can check its version number in three ways:

- In the error log, it is printed during startup.
- `SELECT * FROM information_schema.plugins;`
- `SELECT @@innodb_version;`

The InnoDB Plugin writes its version number to the error log, which can be helpful in diagnosis of errors:

```
091105 12:28:06 InnoDB Plugin 1.0.5 started; log sequence number 46509
```

Note that the `PLUGIN_VERSION` column in the table `INFORMATION_SCHEMA.PLUGINS` does not display the third component of the version number, only the first and second components, as in 1.0.

1.5 Operational Restrictions

Because the InnoDB Plugin introduces a new file format, with new on-disk data structures within both the database and log files, there are important restrictions on the use of the plugin in typical user environments. Specifically, you should pay special attention to the information presented here about file format compatibility with respect to the following scenarios:

- Downgrading from the InnoDB Plugin to the built-in InnoDB, or otherwise using different versions of InnoDB with database files created by the InnoDB Plugin
- Using `mysqldump`.
- Using MySQL replication.
- Using InnoDB Hot Backup.

WARNING: Once you use the InnoDB Plugin on a set of database files, take care to avoid crashes and corruptions when using those files with an earlier version of InnoDB, as might happen by opening the database with MySQL when the plugin is not installed. It is **strongly** recommended that you use a “slow shutdown” (`SET GLOBAL innodb_fast_shutdown=0`) when stopping the MySQL server when the InnoDB Plugin is enabled. This will ensure log files and other system information written by the plugin will not cause problems when using a prior version of InnoDB. See [Section 11.3, “How to Downgrade”](#).

Because of these considerations, and although it may be useful in certain circumstances to use the plugin in a temporary way as just described, many users will find it preferable to test their application with the plugin and use it on an on-going basis, without reverting back to the standard, built-in InnoDB.

WARNING: If you dump a database containing compressed tables with `mysqldump`, the dump file may contain `CREATE TABLE` commands that attempt to create compressed tables, or those using `ROW_FORMAT=DYNAMIC` in the new database. Therefore, you should be sure the new database is running the InnoDB Plugin, with the proper settings for `innodb_file_format` and `innodb_file_per_table`, if you want to have the tables re-created as they exist in the original database. Typically, however, when the `mysqldump` file is loaded, MySQL and InnoDB will ignore `CREATE TABLE` options they do not recognize, and the table(s) will be created in a format used by the running server.

WARNING: If you use MySQL replication, you should be careful to ensure all slaves are configured with the InnoDB Plugin, with the same settings for `innodb_file_format` and `innodb_file_per_table`. If you do not do so, and you create tables that require the new “Barracuda” file format, replication errors may occur. If a slave MySQL server is running the built-in InnoDB, it will ignore the `CREATE TABLE` options to create a compressed table or one with `ROW_FORMAT=DYNAMIC`, and create the table uncompressed, with `ROW_FORMAT=COMPACT`.

WARNING: The current version of InnoDB Hot Backup does not support the new “Barracuda” file format. Using InnoDB Hot Backup Version 3 to backup databases in this format will cause unpredictable behavior. A future version of InnoDB Hot Backup will support databases used with the InnoDB Plugin. As an alternative, you may back up such databases with `mysqldump`.

Chapter 2 Fast Index Creation in the InnoDB Storage Engine

Table of Contents

2.1 Overview of Fast Index Creation	5
2.2 Examples	5
2.3 Implementation	6
2.4 Concurrency Considerations	7
2.5 Crash Recovery	7
2.6 Limitations	8

2.1 Overview of Fast Index Creation

In MySQL versions up to 5.0, adding or dropping an index on a table with existing data can be very slow if the table has many rows. The `CREATE INDEX` and `DROP INDEX` commands work by creating a new, empty table defined with the requested set of indexes. It then copies the existing rows to the new table one-by-one, updating the indexes as it goes. Inserting entries into the indexes in this fashion, where the key values are not sorted, requires random access to the index nodes, and is far from optimal. After all rows from the original table are copied, the old table is dropped and the copy is renamed with the name of the original table.

Beginning with version 5.1, MySQL allows a storage engine to create or drop indexes without copying the contents of the entire table. The standard built-in InnoDB in MySQL version 5.1, however, does not take advantage of this capability. With the InnoDB Plugin, however, users can in most cases add and drop indexes much more efficiently than with prior releases.

In InnoDB, the rows of a table are stored in a clustered (or primary key) index, forming what some database systems call an “index-organized table”. Changing the clustered index requires copying the data, even with the InnoDB Plugin. However, adding or dropping a secondary index with the InnoDB Plugin is much faster, since it does not involve copying the data.

This new mechanism also means that you can generally speed the overall process of creating and loading an indexed table by creating the table with only the clustered index, and adding the secondary indexes after the data is loaded.

Although no syntax changes are required in the `CREATE INDEX` or `DROP INDEX` commands, some factors affect the performance, space usage, and semantics of this operation (see [Section 2.6, “Limitations”](#)).

Because the ability to create and drop indexes does not require use of a new on-disk file format, it is possible to temporarily use the InnoDB Plugin to create or drop an index, and then fall back to using the standard built-in InnoDB in MySQL for normal operations if you wish. See [Chapter 11, *Downgrading from the InnoDB Plugin*](#) for more information.

2.2 Examples

It is possible to create multiple indexes on a table with one `ALTER TABLE` command. This is relatively efficient, because the clustered index of the table needs to be scanned only once (although the data is sorted separately for each new index). For example:

```
CREATE TABLE T1(A INT PRIMARY KEY,  
  B INT, C CHAR(1)) ENGINE=InnoDB;  
INSERT INTO T1 VALUES
```

```
(1,2,'a'), (2,3,'b'), (3,2,'c'), (4,3,'d'), (5,2,'e');  
COMMIT;  
ALTER TABLE T1 ADD INDEX (B), ADD UNIQUE INDEX (C);
```

The above commands will create table `T1` with the clustered index (primary key) on column `A`, insert several rows, and then build two new indexes on columns `B` and `C`. If there were many rows inserted into `T1` before the `ALTER TABLE` command, this approach would be much more efficient than creating the table with all its indexes before loading the data.

You may also create the indexes one at a time, but then the clustered index of the table is scanned (as well as sorted) once for each `CREATE INDEX` command. Thus, the following commands are not as efficient as the `ALTER TABLE` command above, even though neither requires recreating the clustered index for table `T1`.

```
CREATE INDEX B ON T1 (B);  
CREATE UNIQUE INDEX C ON T1 (C);
```

Dropping indexes in the InnoDB Plugin does not require any copying of table data. Thus, you can equally quickly drop multiple indexes with a single `ALTER TABLE` command or multiple `DROP INDEX` commands:

```
ALTER TABLE T1 DROP INDEX B, DROP INDEX C;
```

or

```
DROP INDEX B ON T1;  
DROP INDEX C ON T1;
```

Restructuring the clustered index in InnoDB always requires copying the data in the table. For example, if you create a table without a primary key, InnoDB chooses one for you, which may be the first `UNIQUE` key defined on `NOT NULL` columns, or a system-generated key. Defining a `PRIMARY KEY` later causes the data to be copied, as in the following example:

```
CREATE TABLE T2 (A INT, B INT) ENGINE=InnoDB;  
INSERT INTO T2 VALUES (NULL, 1);  
ALTER TABLE T2 ADD PRIMARY KEY (B);
```

Note that when you create a `UNIQUE` or `PRIMARY KEY` index, InnoDB must do some extra work. For `UNIQUE` indexes, InnoDB checks that the table contains no duplicate values for the key. For a `PRIMARY KEY` index, InnoDB also checks that none of the `PRIMARY KEY` columns contains a `NULL`. It is best to define the primary key when you create a table, so you need not rebuild the table later.

2.3 Implementation

InnoDB has two types of indexes: the clustered index and secondary indexes. Since the clustered index contains the data values in its B-tree nodes, adding or dropping a clustered index does involve copying the data, and creating a new copy of the table. A secondary index, however, contains only the index key and the value of the primary key. This type of index may be created or dropped without copying the data in the clustered index. Furthermore, because the secondary index contains the values of the primary key (used to access the clustered index when needed), when you change the definition of the primary key, thus recreating the clustered index, all secondary indexes are recreated as well.

Dropping a secondary index is simple. Only the internal InnoDB system tables and the MySQL data dictionary tables need to be updated to reflect the fact that the index no longer exists. InnoDB returns the storage used for the index to the tablespace that contained it, so that new indexes or additional table rows may use the space.

To add a secondary index to an existing table, InnoDB scans the table, and sorts the rows using memory buffers and temporary files in order by the value(s) of the secondary index key column(s). The B-tree is then built in key-value order, which is more efficient than inserting rows into an index in random order with respect to the key values. Because the B-tree nodes are split when they fill, building the index in this way results in a higher fill-factor for the index, making it more efficient for subsequent access.

2.4 Concurrency Considerations

While a secondary index is being created or dropped, the table is locked in shared mode. That is, any writes to the table are blocked, but the data in the table may be read. When you alter the clustered index of a table, however, the table is locked in exclusive mode, because the data must be copied. Thus, during the creation of a new clustered index, all operations on the table are blocked.

Before it can start executing, a `CREATE INDEX` or `ALTER TABLE` command must always wait for currently executing transactions that are accessing the table to commit or rollback before it can proceed. In addition, `ALTER TABLE` commands that create a new clustered index must wait for all `SELECT` statements that access the table to complete (or their containing transactions to commit). Even though the original index exists throughout the creation of the new clustered index, no transactions whose execution spans the creation of the index can be accessing the table, because the original table must be dropped when clustered index is restructured.

Once a `CREATE INDEX` or `ALTER TABLE` command that creates a secondary index begins executing, queries may access the table for read access, but may not update the table. If an `ALTER TABLE` command is changing the clustered index, all queries must wait until the operation completes.

A newly-created secondary index contains only data that is current in the table as of the time the `CREATE INDEX` or `ALTER TABLE` command begins to execute. Specifically, a newly-created index contains only the versions of data as of the most-recently committed transactions prior to the creation of the index. The index thus does not contain any rows that were deleted (and therefore marked for deletion) by transactions that completed before the `CREATE INDEX` or `ALTER TABLE` began. Similarly, the index contains only current versions of every row, and none of the old versions of rows that were updated by transactions that ran before the index was created.

Because a newly-created index contains only information about data current at the time the index was created, queries that need to see data that was deleted or changed before the index was created cannot use the index. The only queries that could be affected by this limitation are those executing in transactions that began before the creation of the index was begun. For such queries, unpredictable results could occur. Newer queries can use the index.

2.5 Crash Recovery

No data is lost if the server crashes while an `ALTER TABLE` command is executing. Recovery, however, is different for clustered indexes and secondary indexes.

If the server crashes while creating a secondary index, upon recovery, InnoDB drops any partially created indexes. All you need to do to create the index is to re-run the `ALTER TABLE` or `CREATE INDEX` command.

However, when a crash occurs during the creation of a clustered index, recovery is somewhat more complicated, because the data in the table must be copied to an entirely new clustered index. Remember that all InnoDB tables are stored as clustered indexes. In the following discussion, we use the word table and clustered index interchangeably.

The InnoDB Plugin creates the new clustered index by copying the existing data from the original table to a temporary table that has the desired index structure. Once the data is completely copied to this temporary

table, the original table is renamed with a different temporary table name. The temporary table comprising the new clustered index is then renamed with the name of the original table, and the original table is then dropped from the database.

If a system crash occurs while creating a new clustered index, no data is lost, but users must complete the recovery process using the temporary tables that exist during the process.

Users rarely re-create a clustered index or re-define primary keys on large tables. Because system crashes are uncommon and the situation described here is rare, this manual does not provide information on recovering from this scenario. Instead, contact MySQL support.

2.6 Limitations

Take the following considerations into account when creating or dropping indexes using the InnoDB Plugin:

- During index creation, files are written to the temporary directory (`$TMPDIR` on Unix, `%TEMP%` on Windows, or the value of `--tmpdir` configuration variable). Each temporary file is large enough to hold one column that makes up the new index, and each one is removed as soon as it is merged into the final index.
- Due to a limitation of MySQL, the table is copied, rather than using “Fast Index Creation” when you create an index on a `TEMPORARY TABLE`. This has been reported as MySQL Bug #39833.
- To avoid consistency issues between the InnoDB data dictionary and the MySQL data dictionary, the table is copied, rather than using Fast Index Creation when you use the `ALTER TABLE ... RENAME COLUMN` syntax.
- The command `ALTER IGNORE TABLE t ADD UNIQUE INDEX` does not delete duplicate rows. This has been reported as MySQL Bug #40344. The `IGNORE` keyword is ignored. If any duplicate rows exist, the operation fails with the following error message:

```
ERROR 23000: Duplicate entry '347' for key 'pl'
```

- As noted above, a newly-created index contains only information about data current at the time the index was created. Therefore, you should not run queries in a transaction that might use a secondary index that did not exist at the beginning of the transaction. There is no way for InnoDB to access “old” data that is consistent with the rest of the data read by the transaction. See the discussion of locking in [Section 2.4, “Concurrency Considerations”](#).

Prior to InnoDB Plugin 1.0.4, unexpected results could occur if a query attempts to use an index created after the start of the transaction containing the query. If an old transaction attempts to access a “too new” index, InnoDB Plugin 1.0.4 and later reports an error:

```
ERROR HY000: Table definition has changed, please retry transaction
```

As the error message suggests, committing (or rolling back) the transaction, and restarting it, cures the problem.

- InnoDB Plugin 1.0.2 introduces some improvements in error handling when users attempt to drop indexes. See section [Section 8.7, “Better Error Handling when Dropping Indexes”](#) for details.
- MySQL 5.1 does not support efficient creation or dropping of `FOREIGN KEY` constraints. Therefore, if you use `ALTER TABLE` to add or remove a `REFERENCES` constraint, the child table will be copied, rather than using “Fast Index Creation”.

Chapter 3 InnoDB Data Compression

Table of Contents

3.1 Overview of Table Compression	9
3.2 Enabling Compression for a Table	9
3.2.1 Configuration Parameters for Compression	10
3.2.2 SQL Compression Syntax Warnings and Errors	11
3.3 Tuning InnoDB Compression	12
3.3.1 When to Use Compression	13
3.3.2 Monitoring Compression at Runtime	15
3.4 How Compression Works in InnoDB	16
3.4.1 Compression Algorithms	16
3.4.2 InnoDB Data Storage and Compression	16
3.4.3 Compression and the InnoDB Buffer Pool	18
3.4.4 Compression and the InnoDB Log Files	18

3.1 Overview of Table Compression

Over the years, processors and cache memories have become much faster, but mass storage based on rotating magnetic disks has not kept pace. While the storage capacity of disks has grown by about a factor of 1,000 in the past decade, random seek times and data transfer rates are still severely limited by mechanical constraints. Therefore, many workloads are I/O-bound. The idea of data compression is to pay a small cost in increased CPU utilization for the benefit of smaller databases and reduced I/O to improve throughput, potentially significantly.

The ability to compress user data is an important new capability of the InnoDB Plugin. Compressed tables reduce the size of the database on disk, resulting in fewer reads and writes needed to access the user data. For many InnoDB workloads and many typical user tables (especially with read-intensive applications where sufficient memory is available to keep frequently-used data in memory), compression not only significantly reduces the storage required for the database, but also improves throughput by reducing the I/O workload, at a modest cost in processing overhead. The storage cost savings can be important, but the reduction in I/O costs can be even more valuable. Compression can be especially important for [SSD](#) storage devices, because they tend to have lower capacity than [HDD](#) devices.

3.2 Enabling Compression for a Table

The usual (uncompressed) size of InnoDB data pages is 16KB. Beginning with the InnoDB Plugin, you can use the attributes `ROW_FORMAT=COMPRESSED`, `KEY_BLOCK_SIZE`, or both in the `CREATE TABLE` and `ALTER TABLE` statements to enable table compression. Depending on the combination of option values, InnoDB attempts to compress each page to 1KB, 2KB, 4KB, 8KB, or 16KB.



Note

The term `KEY_BLOCK_SIZE` does not refer to a “key”, but simply specifies the size of compressed pages to use for the table. Likewise, in the InnoDB Plugin, compression is applicable to tables, not to individual rows, despite the option name `ROW_FORMAT`. Because the InnoDB storage engine cannot add syntax to SQL statements, the InnoDB Plugin re-uses the clauses originally defined for [MyISAM](#).

To create a compressed table, you might use a statement like this:

```
CREATE TABLE name
  (column1 INT PRIMARY KEY)
ENGINE=InnoDB
ROW_FORMAT=COMPRESSED
KEY_BLOCK_SIZE=4;
```

If you specify `ROW_FORMAT=COMPRESSED` but not `KEY_BLOCK_SIZE`, the default compressed page size of 8KB is used. If `KEY_BLOCK_SIZE` is specified, you can omit the attribute `ROW_FORMAT=COMPRESSED`.

Setting `KEY_BLOCK_SIZE=16` most often does not result in much compression, since the normal InnoDB page size is 16KB. However, this setting may be useful for tables with many long `BLOB`, `VARCHAR` or `TEXT` columns, because such values often do compress well, and might therefore require fewer “overflow” pages as described later in this section.

Note that compression is specified on a table-by-table basis. All indexes of a table (including the clustered index) are compressed using the same page size, as specified on the `CREATE TABLE` or `ALTER TABLE` statement. Table attributes such as `ROW_FORMAT` and `KEY_BLOCK_SIZE` are not part of the `CREATE INDEX` syntax, and are ignored if they are specified (although you see them in the output of the `SHOW CREATE TABLE` statement).

3.2.1 Configuration Parameters for Compression

Compressed tables are stored in a format that previous versions of InnoDB cannot process. To preserve downward compatibility of database files, compression can be specified only when the “Barracuda” data file format is enabled using the configuration parameter `innodb_file_format`.

Table compression is also not available for the InnoDB system tablespace. The system tablespace (space 0, the `ibdata*` files) may contain user data, but it also contains internal InnoDB system information, and therefore is never compressed. Thus, compression applies only to tables (and indexes) stored in their own tablespaces.

To use compression, enable the “file per table” mode using the configuration parameter `innodb_file_per_table` and enable the “Barracuda” disk file format using the parameter `innodb_file_format`. You can set these parameters in the MySQL option file `my.cnf` or `my.ini`, but both are dynamic parameters that you can change with the `SET` statement without shutting down the MySQL server, as noted in [Section 9.5, “Configuring the InnoDB Plugin”](#).

Specifying `ROW_FORMAT=COMPRESSED` or a `KEY_BLOCK_SIZE` in the `CREATE TABLE` or `ALTER TABLE` statements if the “Barracuda” file format has not been enabled produces these warnings that you can view with the `SHOW WARNINGS` statement:

Level	Code	Message
Warning	1478	InnoDB: KEY_BLOCK_SIZE requires innodb_file_per_table.
Warning	1478	InnoDB: KEY_BLOCK_SIZE requires innodb_file_format=1.
Warning	1478	InnoDB: ignoring KEY_BLOCK_SIZE=4.
Warning	1478	InnoDB: ROW_FORMAT=COMPRESSED requires innodb_file_per_table.
Warning	1478	InnoDB: assuming ROW_FORMAT=COMPACT.



Note

These messages are only warnings, not errors, and the table is created as if the options were not specified. Enabling InnoDB “strict mode” (see [Section 8.5, “InnoDB Strict Mode”](#)) causes InnoDB to generate an error, not a warning, for these cases. In strict mode, the table is not created if the current configuration does not permit using compressed tables.

The “non-strict” behavior is intended to permit you to import a `mysqldump` file into a database that does not support compressed tables, even if the source database contained compressed tables. In that case, the InnoDB Plugin creates the table in `ROW_FORMAT=COMPACT` instead of preventing the operation.

When you import the dump file into a new database, if you want to have the tables re-created as they exist in the original database, ensure the server is running the InnoDB Plugin with the proper settings for the configuration parameters `innodb_file_format` and `innodb_file_per_table`,

3.2.2 SQL Compression Syntax Warnings and Errors

The attribute `KEY_BLOCK_SIZE` is permitted only when `ROW_FORMAT` is specified as `COMPRESSED` or is omitted. Specifying a `KEY_BLOCK_SIZE` with any other `ROW_FORMAT` generates a warning that you can view with `SHOW WARNINGS`. However, the table is non-compressed; the specified `KEY_BLOCK_SIZE` is ignored).

Level	Code	Message
Warning	1478	InnoDB: ignoring KEY_BLOCK_SIZE=n unless ROW_FORMAT=COMPRESSED.

If you are running in InnoDB strict mode, the combination of a `KEY_BLOCK_SIZE` with any `ROW_FORMAT` other than `COMPRESSED` generates an error, not a warning, and the table is not created.

Table 3.1, “Meaning of `CREATE TABLE` and `ALTER TABLE` Options” summarizes how the various options on `CREATE TABLE` and `ALTER TABLE` are handled.

Table 3.1 Meaning of `CREATE TABLE` and `ALTER TABLE` Options

Option	Usage	Description
<code>ROW_FORMAT=REDUNDANT</code>	Storage format used prior to MySQL 5.0.3	Less efficient than <code>ROW_FORMAT=COMPACT</code> ; for backward compatibility
<code>ROW_FORMAT=COMPACT</code>	Default storage format since MySQL 5.0.3	Stores a prefix of 768 bytes of long column values in the clustered index page, with the remaining bytes stored in an overflow page
<code>ROW_FORMAT=DYNAMIC</code>	Available only with <code>innodb_file_format=Barracuda</code>	Store values within the clustered index page if they fit; if not, stores only a 20-byte pointer to an overflow page (no prefix)
<code>ROW_FORMAT=COMPRESSED</code>	Available only with <code>innodb_file_format=Barracuda</code>	Compresses the table and indexes using zlib to default compressed page size of 8K bytes; implies <code>ROW_FORMAT=DYNAMIC</code>
<code>KEY_BLOCK_SIZE=n</code>	Available only with <code>innodb_file_format=Barracuda</code>	Specifies compressed page size of 1, 2, 4, 8 or 16K bytes; implies <code>ROW_FORMAT=DYNAMIC</code> and <code>ROW_FORMAT=COMPRESSED</code>

Table 3.2, “`CREATE/ALTER TABLE Warnings and Errors when InnoDB Strict Mode is OFF`” summarizes error conditions that occur with certain combinations of configuration parameters and options on the `CREATE TABLE` or `ALTER TABLE` statements, and how the options appear in the output of `SHOW TABLE STATUS`.

When InnoDB strict mode is `OFF`, InnoDB creates or alters the table, but may ignore certain settings, as shown below. You can see the warning messages in the MySQL error log. When InnoDB strict mode is `ON`, these specified combinations of options generate errors, and the table is not created or altered. You can see the full description of the error condition with `SHOW ERRORS`. For example:

```
mysql> CREATE TABLE x (id INT PRIMARY KEY, c INT)
```

```
-> ENGINE=INNODB KEY_BLOCK_SIZE=33333;

ERROR 1005 (HY000): Can't create table 'test.x' (errno: 1478)

mysql> SHOW ERRORS;
+-----+-----+-----+
| Level | Code | Message                                     |
+-----+-----+-----+
| Error | 1478 | InnoDB: invalid KEY_BLOCK_SIZE=33333.     |
| Error | 1005 | Can't create table 'test.x' (errno: 1478) |
+-----+-----+-----+

2 rows in set (0.00 sec)
```

Table 3.2 CREATE/ALTER TABLE Warnings and Errors when InnoDB Strict Mode is OFF

Syntax	Warning or Error Condition	Resulting ROW_FORMAT, as shown in SHOW TABLE STATUS
ROW_FORMAT=REDUNDANT	None	REDUNDANT
ROW_FORMAT=COMPACT	None	COMPACT
ROW_FORMAT=COMPRESSED or ROW_FORMAT=DYNAMIC or KEY_BLOCK_SIZE is specified	Ignored unless you override the default settings for innodb_file_format and innodb_file_per_table	COMPACT
Invalid KEY_BLOCK_SIZE is specified (not 1, 2, 4, 8 or 16)	KEY_BLOCK_SIZE is ignored	the requested one, or COMPACT by default
ROW_FORMAT=COMPRESSED and valid KEY_BLOCK_SIZE are specified	None; KEY_BLOCK_SIZE specified is used, not the 8K default	COMPRESSED
KEY_BLOCK_SIZE is specified with REDUNDANT, COMPACT or DYNAMIC row format	KEY_BLOCK_SIZE is ignored	REDUNDANT, COMPACT or DYNAMIC
ROW_FORMAT is not one of REDUNDANT, COMPACT, DYNAMIC or COMPRESSED	Ignored if recognized by the MySQL parser. Otherwise, an error is issued.	COMPACT or N/A

When InnoDB strict mode is ON (`innodb_strict_mode=1`), the InnoDB Plugin rejects invalid `ROW_FORMAT` or `KEY_BLOCK_SIZE` parameters. For compatibility with the built-in InnoDB in MySQL, InnoDB strict mode is not enabled by default, and in this default non-strict mode, the InnoDB Plugin issues warnings (not errors) for ignored invalid parameters.

Note that it is not possible to see the chosen `KEY_BLOCK_SIZE` using `SHOW TABLE STATUS`. The statement `SHOW CREATE TABLE` displays the `KEY_BLOCK_SIZE` (even if it was ignored by InnoDB). The real compressed page size inside InnoDB cannot be displayed by MySQL.

3.3 Tuning InnoDB Compression

Most often, the internal optimizations in InnoDB described in a later portion of this section, ensure that the system runs well with compressed data. However, because the efficiency of compression depends on the nature of your data, there are some factors you should consider to get best performance. You need to

choose which tables to compress, and what compressed page size to use. You may also want to adjust the size of the buffer pool based on run-time performance characteristics such as the amount of time the system spends compressing and uncompressing data.

3.3.1 When to Use Compression

In general, compression works best on tables that include a reasonable number of character string columns and where the data is read far more often than it is written. Because there are no guaranteed ways to predict whether or not compression benefits a particular situation, always test with a specific workload and data set running on a representative configuration. Consider the following factors when deciding which tables to compress.

3.3.1.1 Data Characteristics and Compression

A key determinant of the efficiency of compression in reducing the size of data files is the nature of the data itself. Recall that compression works by identifying repeated strings of bytes in a block of data. Completely randomized data is the worst case. Typical data often has repeated values, and so compresses effectively. Character strings often compress well, whether defined in [CHAR](#), [VARCHAR](#), [TEXT](#) or [BLOB](#) columns. On the other hand, tables containing mostly binary data (integers or floating point numbers) or data that is previously compressed (for example JPEG or PNG images) may not generally compress well, significantly or at all.

Compression is chosen on a table by table basis with the InnoDB Plugin, and a table and all of its indexes use the same (compressed) page size. It might be that the primary key (clustered) index, which contains the data for all columns of a table, compresses more effectively than the secondary indexes. For those cases where there are long rows, the use of compression may result in long column values being stored “off-page”, as discussed in [Section 5.3, “DYNAMIC Row Format”](#). Those overflow pages may compress well. Given these considerations, for many applications, some tables compress more effectively than others, and you may find that your workload performs best only with a subset of tables compressed.

Experimenting is the only way to determine whether or not to compress a particular table. InnoDB compresses data in 16K chunks corresponding to the uncompressed page size, and in addition to user data, the page format includes some internal system data that is not compressed. Compression utilities compress an entire stream of data, and so may find more repeated strings across the entire input stream than InnoDB would find in a table compressed in 16K chunks. But you can get a sense of how compression efficiency by using a utility that implements LZ77 compression (such as [gzip](#) or WinZip) on your data file.

Another way to test compression on a specific table is to copy some data from your uncompressed table to a similar, compressed table (having all the same indexes) and look at the size of the resulting file. When you do so (if nothing else using compression is running), you can examine the ratio of successful compression operations to overall compression operations. (In the [INNODB_CMP](#) table, compare [COMPRESS_OPS](#) to [COMPRESS_OPS_OK](#). See [INNODB_CMP](#) for more information.) If a high percentage of compression operations complete successfully, the table might be a good candidate for compression.

3.3.1.2 Compression and Application and Schema Design

Decide whether to compress data in your application or in the InnoDB table. It is usually not sensible to store data that is compressed by an application in an InnoDB compressed table. Further compression is extremely unlikely, and the attempt to compress just wastes CPU cycles.

Compressing in the Database

The InnoDB table compression is automatic and applies to all columns and index values. The columns can still be tested with operators such as [LIKE](#), and sort operations can still use indexes even when the index values are compressed. Because indexes are often a significant fraction of the total size of a database,

compression could result in significant savings in storage, I/O or processor time. The compression and decompression operations happen on the database server, which likely is a powerful system that is sized to handle the expected load.

Compressing in the Application

If you compress data such as text in your application, before it is inserted into the database, You might save overhead for data that does not compress well by compressing some columns and not others. This approach uses CPU cycles for compression and uncompression on the client machine rather than the database server, which might be appropriate for a distributed application with many clients, or where the client machine has spare CPU cycles.

Hybrid Approach

Of course, it is possible to combine these approaches. For some applications, it may be appropriate to use some compressed tables and some uncompressed tables. It may be best to externally compress some data (and store it in uncompressed InnoDB tables) and allow InnoDB to compress (some of) the other tables in the application. As always, up-front design and real-life testing are valuable in reaching the right decision.

3.3.1.3 Workload Characteristics and Compression

In addition to choosing which tables to compress (and the page size), the workload is another key determinant of performance. If the application is dominated by reads, rather than updates, fewer pages need to be reorganized and recompressed after the index page runs out of room for the per-page “modification log” that InnoDB maintains for compressed data. If the updates predominantly change non-indexed columns or those containing `BLOBs` or large strings that happen to be stored “off-page”, the overhead of compression may be acceptable. If the only changes to a table are `INSERTs` that use a monotonically increasing primary key, and there are few secondary indexes, there is little need to reorganize and recompress index pages. Since InnoDB can “delete-mark” and delete rows on compressed pages “in place” by modifying uncompressed data, `DELETE` operations on a table are relatively efficient.

For some environments, the time it takes to load data can be as important as run-time retrieval. Especially in data warehouse environments, many tables may be read-only or read-mostly. In those cases, it might or might not be acceptable to pay the price of compression in terms of increased load time, unless the resulting savings in fewer disk reads or in storage cost is significant.

Fundamentally, compression works best when the CPU time is available for compressing and uncompressing data. Thus, if your workload is I/O bound, rather than CPU-bound, you may find that compression can improve overall performance. Therefore when you test your application performance with different compression configurations, it is important to test on a platform similar to the planned configuration of the production system.

3.3.1.4 Configuration Characteristics and Compression

Reading and writing database pages from and to disk is the slowest aspect of system performance. Therefore, compression attempts to reduce I/O by using CPU time to compress and uncompress data, and thus is most effective when I/O is a relatively scarce resource compared to processor cycles.

This is often especially the case when running in a multi-user environment with fast, multi-core CPUs. When a page of a compressed table is in memory, InnoDB often uses an additional 16K in the buffer pool for an uncompressed copy of the page. The adaptive LRU algorithm in the InnoDB Plugin attempts to balance the use of memory between compressed and uncompressed pages to take into account whether the workload is running in an I/O-bound or CPU-bound manner. Nevertheless, a configuration with more memory dedicated to the InnoDB buffer pool tends to run better when using compressed tables than a configuration where memory is highly constrained.

3.3.1.5 Choosing the Compressed Page Size

The optimal setting of the compressed page size depends on the type and distribution of data that the table and its indexes contain. The compressed page size should always be bigger than the maximum record size, or operations may fail as noted in the discussion of compression internals later in this section.

Setting the compressed page size too large wastes some space, but the pages do not have to be compressed as often. If the compressed page size is set too small, inserts or updates may require time-consuming recompression, and the B-tree nodes may have to be split more frequently, leading to bigger data files and less efficient indexing.

Typically, one would set the compressed page size to 8K or 4K bytes. Given that the maximum InnoDB record size is around 8K, `KEY_BLOCK_SIZE=8` is usually a safe choice.

3.3.2 Monitoring Compression at Runtime

The current version of the InnoDB Plugin provides only a limited means to monitor the performance of compression at runtime. Overall application performance, CPU and I/O utilization and the size of disk files are the best indicators of how effective compression is for your application.

The InnoDB Plugin does include some Information Schema tables (see [Example 6.1, “Using the Compression Information Schema Tables”](#)) that reflect the internal use of memory and the rates of compression used overall. The `INNODB_CMP` tables report information about compression activity for each compressed page size (`KEY_BLOCK_SIZE`) in use. The information in these tables is system-wide, and includes summary data across all compressed tables in your database. You can use this data to help decide whether or not to compress a table by examining these tables when no other compressed tables are being accessed.

The key statistics to consider are the number of, and amount of time spent performing, compression and uncompression operations. Since InnoDB must split B-tree nodes when they are too full to contain the compressed data following a modification, you should also compare the number of “successful” compression operations with the number of such operations overall. Based on the information in the `INNODB_CMP` tables and overall application performance and hardware resource utilization, you may decide to make changes in your hardware configuration, adjust the size of the InnoDB buffer pool, choose a different page size, or select a different set of tables to compress.

If the amount of CPU time required for compressing and uncompressing is high, changing to faster CPUs, or those with more cores, can help improve performance with the same data, application workload and set of compressed tables. You may also benefit by increasing the size of the InnoDB buffer pool, so that more uncompressed pages can stay in memory, reducing the need to uncompress pages which exist in memory only in compressed form.

A large number of compression operations overall (compared to the number of `INSERT`, `UPDATE` and `DELETE` operations in your application and the size of the database) could indicate that some of your compressed tables are being updated too heavily for effective compression. You may want to choose a larger page size, or be more selective about which tables you compress.

If the number of “successful” compression operations (`COMPRESS_OPS_OK`) is a high percentage of the total number of compression operations (`COMPRESS_OPS`), then the system is likely performing well. However, if the ratio is low, then InnoDB is being caused to reorganize, recompress and split B-tree nodes more often than is desirable. In this case, you may want to avoid compressing some tables or choose a larger `KEY_BLOCK_SIZE` for some of the tables for which you are using compression. You may not want to compress tables which cause the number of “compression failures” in your application to be more than 1% or 2% of the total (although this may be acceptable during a data load, for example, if your application does not encounter such a ratio during normal operations).

3.4 How Compression Works in InnoDB

This section describes some internal implementation details about compression in InnoDB. The information presented here may be helpful in tuning for performance, but is not necessary to know for basic use of compression.

3.4.1 Compression Algorithms

Some operating systems implement compression at the file system level. Files are typically divided into fixed-size blocks that are compressed into variable-size blocks, which easily leads into fragmentation. Every time something inside a block is modified, the whole block is recompressed before it is written to disk. These properties make this compression technique unsuitable for use in an update-intensive database system.

The InnoDB Plugin implements a novel type of compression with the help of the well-known [zlib library](#), which implements the LZ77 compression algorithm. This compression algorithm is mature, robust, and efficient in both CPU utilization and in reduction of data size. The algorithm is “lossless”, so that the original uncompressed data can always be reconstructed from the compressed form. LZ77 compression works by finding sequences of data that are repeated within the data to be compressed. The patterns of values in your data determine how well it compresses, but typical user data often compresses by 50% or more.

Unlike compression performed by an application, or compression features of some other database management systems, InnoDB compression applies both to user data and to indexes. In many cases, indexes can constitute 40-50% or more of the total database size, so this difference is significant. When compression is working well for a data set, the size of the InnoDB data files (the `.ibd` files) is 25% to 50% of the uncompressed size or possibly smaller. Depending on the workload, this smaller database can in turn lead to a reduction in I/O, and an increase in throughput, at a modest cost in terms of increased CPU utilization.

3.4.2 InnoDB Data Storage and Compression

All user data in InnoDB is stored in pages comprising a B-tree index (the so-called clustered index). In some other database systems, this type of index is called an “index-organized table”. Each row in the index node contains the values of the (user-specified or system-generated) primary key and all the other columns of the table.

Secondary indexes in InnoDB are also B-trees, containing pairs of values: the index key and a pointer to a row in the clustered index. The pointer is in fact the value of the primary key of the table, which is used to access the clustered index if columns other than the index key and primary key are required. Secondary index records must always fit on a single B-tree page.

The compression of B-tree nodes (of both clustered and secondary indexes) is handled differently from compression of overflow pages used to store long [VARCHAR](#), [BLOB](#), or [TEXT](#) columns, as explained in the following sections.

3.4.2.1 Compression of B-Tree Pages

Because they are frequently updated, B-tree pages require special treatment. It is important to minimize the number of times B-tree nodes are split, as well as to minimize the need to uncompress and recompress their content.

One technique InnoDB uses is to maintain some system information in the B-tree node in uncompressed form, thus facilitating certain in-place updates. For example, this allows rows to be delete-marked and deleted without any compression operation.

In addition, InnoDB attempts to avoid unnecessary uncompression and recompression of index pages when they are changed. Within each B-tree page, the system keeps an uncompressed “modification log” to record changes made to the page. Updates and inserts of small records may be written to this modification log without requiring the entire page to be completely reconstructed.

When the space for the modification log runs out, InnoDB uncompresses the page, applies the changes and recompresses the page. If recompression fails, the B-tree nodes are split and the process is repeated until the update or insert succeeds.

Generally, InnoDB requires that each B-tree page can accommodate at least two records. For compressed tables, this requirement has been relaxed. Leaf pages of B-tree nodes (whether of the primary key or secondary indexes) only need to accommodate one record, but that record must fit in uncompressed form, in the per-page modification log. Starting with InnoDB Plugin version 1.0.2, and if InnoDB strict mode is `ON`, the InnoDB Plugin checks the maximum row size during `CREATE TABLE` or `CREATE INDEX`. If the row does not fit, the following error message is issued: `ERROR HY000: Too big row.`

If you create a table when InnoDB strict mode is `OFF`, and a subsequent `INSERT` or `UPDATE` statement attempts to create an index entry that does not fit in the size of the compressed page, the operation fails with `ERROR 42000: Row size too large.` (This error message does not name the index for which the record is too large, or mention the length of the index record or the maximum record size on that particular index page.) To solve this problem, rebuild the table with `ALTER TABLE` and select a larger compressed page size (`KEY_BLOCK_SIZE`), shorten any column prefix indexes, or disable compression entirely with `ROW_FORMAT=DYNAMIC` or `ROW_FORMAT=COMPACT`.

3.4.2.2 Compressing BLOB, VARCHAR and TEXT Columns

In a clustered index, `BLOB`, `VARCHAR` and `TEXT` columns that are not part of the primary key may be stored on separately allocated (“overflow”) pages. We call these “off-page columns” whose values are stored on singly-linked lists of overflow pages.

For tables created in `ROW_FORMAT=DYNAMIC` or `ROW_FORMAT=COMPRESSED`, the values of `BLOB`, `TEXT` or `VARCHAR` columns may be stored fully off-page, depending on their length and the length of the entire row. For columns that are stored off-page, the clustered index record only contains 20-byte pointers to the overflow pages, one per column. Whether any columns are stored off-page depends on the page size and the total size of the row. When the row is too long to fit entirely within the page of the clustered index, InnoDB chooses the longest columns for off-page storage until the row fits on the clustered index page. As noted above, if a row does not fit by itself on a compressed page, an error occurs.

Tables created in previous versions of InnoDB use the “Antelope” file format, which supports only `ROW_FORMAT=REDUNDANT` and `ROW_FORMAT=COMPACT`. In these formats, InnoDB stores the first 768 bytes of `BLOB`, `VARCHAR` and `TEXT` columns in the clustered index record along with the primary key. The 768-byte prefix is followed by a 20-byte pointer to the overflow pages that contain the rest of the column value.

When a table is in `COMPRESSED` format, all data written to overflow pages is compressed “as is”; that is, InnoDB applies the zlib compression algorithm to the entire data item. Other than the data, compressed overflow pages contain an uncompressed header and trailer comprising a page checksum and a link to the next overflow page, among other things. Therefore, very significant storage savings can be obtained for longer `BLOB`, `TEXT` or `VARCHAR` columns if the data is highly compressible, as is often the case with text data (but not previously compressed images).

The overflow pages are of the same size as other pages. A row containing ten columns stored off-page occupies ten overflow pages, even if the total length of the columns is only 8K bytes. In an uncompressed table, ten uncompressed overflow pages occupy 160K bytes. In a compressed table with an 8K page size, they occupy only 80K bytes. Thus, it is often more efficient to use compressed table format for tables with long column values.

Using a 16K compressed page size can reduce storage and I/O costs for [BLOB](#), [VARCHAR](#) or [TEXT](#) columns, because such data often compress well, and might therefore require fewer “overflow” pages, even though the B-tree nodes themselves take as many pages as in the uncompressed form.

3.4.3 Compression and the InnoDB Buffer Pool

In a compressed InnoDB table, every compressed page (whether 1K, 2K, 4K or 8K) corresponds to an uncompressed page of 16K bytes. To access the data in a page, InnoDB must read the compressed page from disk (unless it is already in memory), and then uncompress the page to its original 16K byte form. This section describes how InnoDB manages the buffer pool with respect to pages of compressed tables.

To minimize I/O and to reduce the need to uncompress a page, at times the buffer pool contains both the compressed and uncompressed form of a database page. However, to make room for other required database pages, InnoDB may “evict” from the buffer pool an uncompressed page, while leaving the compressed page in memory. Or, if a page has not been accessed in a while, the compressed form of the page may be written to disk, to free space for other data. Thus, at any given time, the buffer pool may (a) not contain any copy of a given database page, (b) contain only the compressed form of the page, or (c) contain both the compressed and uncompressed forms of the page.

InnoDB keeps track of which pages to retain in memory and which to evict using a least-recently-used (LRU) list, so that “hot” or frequently accessed data tends to stay in memory. When compressed tables are accessed, InnoDB uses an adaptive LRU algorithm to achieve an appropriate balance of compressed and uncompressed pages in memory. This adaptive algorithm is sensitive to whether the system is running in an I/O-bound or CPU-bound manner.

The essential idea is to avoid spending too much processing time uncompressing pages when the CPU is busy, and to avoid doing excess I/O when the CPU has spare cycles that can be used for uncompressing compressed pages (that may already be in memory). When the system is I/O-bound, the algorithm prefers to evict the uncompressed copy of a page rather than both copies, to make more room for other disk pages to become memory resident. When the system is CPU-bound, InnoDB prefers to evict both the compressed and uncompressed page, so that more memory can be used for “hot” pages and reducing the need to uncompress data in memory only in compressed form.

3.4.4 Compression and the InnoDB Log Files

Before (but not necessarily at the same time as) a compressed page is written to a database file, InnoDB writes a copy of the page to the redo log (if it has been recompressed since the last time it was written to the database). This is done to ensure that redo logs will always be usable, even if a future version of InnoDB uses a slightly different compression algorithm. Therefore, some increase in the size of log files, or a need for more frequent checkpoints, can be expected when using compression. The amount of increase in the log file size or checkpoint frequency depends on the number of times compressed pages are modified in a way that requires reorganization and recompression.

Note that the redo log file format (and the database file format) are different from previous releases when using compression. The current release of InnoDB Hot Backup (version 3) therefore does not support databases that use compression. Only databases using the file format “Antelope” can be backed up online by InnoDB Hot Backup.

Chapter 4 InnoDB File-Format Management

Table of Contents

4.1 Overview of InnoDB File Formats	19
4.2 Named File Formats	19
4.3 Enabling File Formats	20
4.4 File Format Compatibility	20
4.4.1 Startup File Format Compatibility Checking	21
4.4.2 Table-Access File Format Compatibility Checking	22
4.5 Identifying the File Format in Use	23
4.6 Downgrading the File Format	24
4.7 Future InnoDB File Formats	24

4.1 Overview of InnoDB File Formats

As InnoDB evolves, new on-disk data structures are sometimes required to support new features. This release of InnoDB introduces two such new data structures: compressed tables (see [Chapter 3, InnoDB Data Compression](#)), and long variable-length columns stored off-page (see [Chapter 5, InnoDB Row Storage and Row Formats](#)). These features both require use of the new [Barracuda](#) file format.



Note

These new data structures are not compatible with prior versions of InnoDB. The other new features of the InnoDB Plugin are compatible with the original [Antelope](#) file format and do not require the Barracuda file format.

In general, a newer version of InnoDB may create a table or index that cannot safely be read or written with a prior version of InnoDB without risk of crashes, hangs, wrong results or corruptions. The InnoDB Plugin introduces a new mechanism to guard against these conditions, and to help preserve compatibility among database files and versions of InnoDB. This mechanism lets you take advantage of some new features of an InnoDB release (e.g., performance improvements and bug fixes), and still preserve the option of using your database with a prior version of InnoDB, by precluding the use of new features that create downward incompatible on-disk data structures.

4.2 Named File Formats

The InnoDB Plugin introduces the idea of a named file format and a configuration parameter to enable the use of features that require use of that format. The new file format is the [Barracuda](#) format, and the file format supported by prior releases of InnoDB is known as [Antelope](#). Compressed tables and the new row format that stores long columns “off-page” require the use of the Barracuda file format or newer. Future versions of InnoDB may introduce a series of file formats, identified with the names of animals, in ascending alphabetic order.

Beginning with this release, every InnoDB per-table tablespace file is labeled with a file format identifier. This does not apply to the system tablespace (the `ibdata` files) but only the files of separate tablespaces (the `*.ibd` files where tables and indexes are stored in their own tablespace). As noted below, however, the system tablespace is tagged with the “highest” file format in use in a group of InnoDB database files, and this tag is checked when the files are opened.

In this release, when you create a compressed table, or a table with `ROW_FORMAT=DYNAMIC`, the file header for the corresponding `.ibd` file and the table type in the InnoDB data dictionary are updated

with the identifier for the “Barracuda” file format. From that point forward, the table cannot be used with a version of InnoDB that does not support this new file format. To protect against anomalous behavior, InnoDB version 5.0.21 and later performs a compatibility check when the table is opened, as described below. (Note that the `ALTER TABLE` command in many cases, causes a table to be recreated and thereby change its properties. The special case of adding or dropping indexes without rebuilding the table is described in [Chapter 2, Fast Index Creation in the InnoDB Storage Engine](#).)

If a version of InnoDB supports a particular file format (whether or not it is enabled), you can access and even update any table that requires that format or an earlier format. Only the creation of new tables using new features is limited based on the particular file format enabled. Conversely, if a tablespace contains a table or index that uses a file format that is not supported by the currently running software, it cannot be accessed at all, even for read access.

The only way to “downgrade” an InnoDB tablespace to an earlier file format is to copy the data to a new table, in a tablespace that uses the earlier format. This can be done with the `ALTER TABLE` command, as described in [Section 4.6, “Downgrading the File Format”](#).

The easiest way to determine the file format of an existing InnoDB tablespace is to examine the properties of the table it contains, using the `SHOW TABLE STATUS` command or querying the table `INFORMATION_SCHEMA.TABLES`. If the `Row_format` of the table is reported as `'Compressed'` or `'Dynamic'`, the tablespace containing the table uses the “Barracuda” format. Otherwise, it uses the prior InnoDB file format, “Antelope”.

4.3 Enabling File Formats

The new configuration parameter `innodb_file_format` controls whether such commands as `CREATE TABLE` and `ALTER TABLE` can be used to create tables that depend on support for the “Barracuda” file format.

The file format is a dynamic, global parameter that can be specified in the MySQL option file (`my.cnf` or `my.ini`) or changed with the `SET GLOBAL` command, as described in [Section 9.5, “Configuring the InnoDB Plugin”](#).

4.4 File Format Compatibility

To avoid confusion, for the purposes of this discussion we define the term “ib-file set” to mean the set of operating system files that InnoDB manages as a unit. The ib-file set includes the following files:

- The system tablespace (one or more `ibdata` files) that contain internal system information (including internal catalogs and undo information) and may include user data and indexes.
- Zero or more single-table tablespaces (also called “file per table” files, named `*.ibd` files).
- (Usually two) InnoDB log files (`ib_logfile0` and `ib_logfile1`), used for crash recovery and in backups.

This collection of files is transactionally consistent, and recoverable as a unit. An “ib-file set” specifically does not include the related MySQL `.frm` files that contain metadata about InnoDB tables. The `.frm` files are created and managed exclusively by MySQL, and can sometimes get out of sync with the internal metadata in InnoDB.

Instead of “ib-file set”, we might call such a collection a “database”. However, MySQL uses the word “database” to mean a logical collection of tables, what other systems term a “schema” or “catalog”. Given MySQL terminology, multiple tables (even from more than one database) can be stored in a single “ib-file set”.

The InnoDB Plugin incorporates several checks to guard against the possible crashes and data corruptions that might occur if you use an ib-file set in a file format that is not supported by the software release in use. These checks take place when the server is started, and when you first access a table. This section describes these checks, how you can control them, and error and warning conditions that may arise.

4.4.1 Startup File Format Compatibility Checking

To prevent possible crashes or data corruptions when InnoDB Plugin opens an ib-file set, it checks that it can fully support the file formats in use within the ib-file set. If the system is restarted following a crash, or a “fast shutdown” (i.e., `innodb_fast_shutdown` is greater than zero), there may be on-disk data structures (such as redo or undo entries, or doublewrite pages) that are in a “too-new” format for the current software. During the recovery process, serious damage can be done to your data files if these data structures are accessed. The startup check of the file format occurs before any recovery process begins, thereby preventing the problems described in [Section 11.4, “Possible Problems”](#).

Beginning with version 1.0.1 of the InnoDB Plugin, the system tablespace records an identifier or tag for the “highest” file format used by any table in any of the tablespaces that is part of the ib-file set. Checks against this file format tag are controlled by the new configuration parameter `innodb_file_format_check`, which is `ON` by default.

If the file format tag in the system tablespace is newer or higher than the highest version supported by the particular currently executing software and if `innodb_file_format_check` is `ON`, the following error is issued when the server is started:

```
InnoDB: Error: the system tablespace is in a file format that this version doesn't support
```

You can also set `innodb_file_format` to a file format name. Doing so prevents the InnoDB Plugin from starting if the current software does not support the file format specified. It also sets the “high water mark” to the value you specify. The ability to set `innodb_file_format_check` will be useful (with future releases of InnoDB) if you manually “downgrade” all of the tables in an ib-file set (as described in [Chapter 11, *Downgrading from the InnoDB Plugin*](#)). You can then rely on the file format check at startup if you subsequently use an older version of InnoDB to access the ib-file set.

In some limited circumstances, you might want to start the server and use an ib-file set that is in a “too new” format (one that is not supported by the software you are using). If you set the configuration parameter `innodb_file_format_check` to `OFF`, the InnoDB Plugin opens the database, but issues this warning message in the error log:

```
InnoDB: Warning: the system tablespace is in a
file format that this version doesn't support
```



Note

This is a very dangerous setting, as it permits the recovery process to run, possibly corrupting your database if the previous shutdown was a crash or “fast shutdown”. You should only set `innodb_file_format_check` to `OFF` if you are sure that the previous shutdown was done with `innodb_fast_shutdown=0`, so that essentially no recovery process occurs. In a future release, this parameter setting may be renamed from `OFF` to `UNSAFE`. (However, until there are newer releases of InnoDB that support additional file formats, even disabling the startup checking is in fact “safe”.)

Note that the parameter `innodb_file_format_check` affects only what happens when a database is opened, not subsequently. Conversely, the parameter `innodb_file_format` (which enables a specific

format) only determines whether or not a new table can be created in the enabled format and has no effect on whether or not a database can be opened.

The file format tag is a “high water mark”, and as such it is increased after the server is started, if a table in a “higher” format is created or an existing table is accessed for read or write (assuming its format is supported). If you access an existing table in a format higher than the format the running software supports, the system tablespace tag is not updated, but table-level compatibility checking applies (and an error is issued), as described in [Section 4.4.2, “Table-Access File Format Compatibility Checking”](#). Any time the high water mark is updated, the value of `innodb_file_format_check` is updated as well, so the command `SELECT @@innodb_file_format_check;` displays the name of the newest file format known to be used by tables in the currently open ib-file set and supported by the currently executing software.

To best illustrate this behavior, consider the scenario described in [Table 4.1, “InnoDB Data File Compatibility and Related InnoDB Parameters”](#). Imagine that some future version of InnoDB supports the “Cheetah” format and that an ib-file set has been used with that version.

Table 4.1 InnoDB Data File Compatibility and Related InnoDB Parameters

innodb file format check	innodb file format	Highest file format used in ib-file set	Highest file format supported by InnoDB	Result
OFF	Antelope or Barracuda	Barracuda	Barracuda	Database can be opened; tables can be created which require “Antelope” or “Barracuda” file format
OFF	Antelope or Barracuda	Cheetah	Barracuda	Database can be opened with a warning, since the database contains files in a “too new” format; tables can be created which require “Antelope” or “Barracuda” file format; tables in “Cheetah” format cannot be accessed
OFF	Cheetah	Barracuda	Barracuda	Database cannot be opened; <code>innodb_file_format</code> cannot be set to “Cheetah”
ON	Antelope or Barracuda	Barracuda	Barracuda	Database can be opened; tables can be created which require “Antelope” or “Barracuda” file format
ON	Antelope or Barracuda	Cheetah	Barracuda	Database cannot be opened, since the database contains files in a “too new” format (“Cheetah”)
ON	Cheetah	Barracuda	Barracuda	Database cannot be opened; <code>innodb_file_format</code> cannot be set to “Cheetah”

4.4.2 Table-Access File Format Compatibility Checking

When a table is first accessed, InnoDB (including some releases prior to InnoDB Plugin 1.0) check that the file format of the tablespace in which the table is stored is fully supported. This check prevents crashes or corruptions that would otherwise occur when tables using a “too new” data structure are encountered.

Note that all tables using any file format supported by a release can be read or written (assuming the user has sufficient privileges). The setting of the system configuration parameter `innodb_file_format` can prevent creating a new table that uses specific file formats, even if they are supported by a given release. Such a setting might be used to preserve backward compatibility, but it does not prevent accessing any table that uses any supported format.

As noted in [Section 4.2, “Named File Formats”](#), versions of InnoDB older than 5.0.21 cannot reliably use database files created by newer versions if a new file format was used when a table was created. To

prevent various error conditions or corruptions, InnoDB checks file format compatibility when it opens a file (e.g., upon first access to a table). If the currently running version of InnoDB does not support the file format identified by the table type in the InnoDB data dictionary, MySQL reports the following error:

```
ERROR 1146 (42S02): Table 'test.t1' doesn't exist
```

Furthermore, InnoDB writes a message to the error log:

```
InnoDB: table test/t1: unknown table type 33
```

The table type should be equal to the tablespace flags, which contains the file format version as discussed in [Section 4.5, “Identifying the File Format in Use”](#).

Versions of InnoDB prior to 4.1 did not include table format identifiers in the database files, and versions prior to 5.0.21 did not include a table format compatibility check. Therefore, there is no way to ensure proper operations if a table in a “too new” format is used with versions of InnoDB prior to 5.0.21.

The new file format management capability introduced with the InnoDB Plugin (comprising tablespace tagging and run-time checks) allows InnoDB to verify as soon as possible that the running version of software can properly process the tables existing in the database.

If you permit InnoDB to open a database containing files in a format it does not support (by setting the parameter `innodb_file_format_check` to `OFF`), the table-level checking described in this section still applies.

Users are *strongly* urged not to use database files that contain “Barracuda” file format tables with releases of InnoDB older than the InnoDB Plugin. It is possible to “downgrade” such tables to the “Antelope” format (that the built-in InnoDB in MySQL up to version 5.1 supports) with the procedure described in [Section 4.6, “Downgrading the File Format”](#).

4.5 Identifying the File Format in Use

Although you may have enabled a given `innodb_file_format` at a particular time, unless you create a new table, the database file format is unchanged. If you do create a new table, the tablespace containing the table is tagged with the “earliest” or “simplest” file format that is required for the table’s features. For example, if you enable file format “Barracuda”, and create a new table that is not compressed and does not use `ROW_FORMAT=DYNAMIC`, the new tablespace that contains the table is tagged as using file format “Antelope”.

It is easy to identify the file format used by a given tablespace or table. The table uses the “Barracuda” format if the `Row_format` reported by `SHOW CREATE TABLE` or `INFORMATION_SCHEMA.TABLES` is one of `'Compressed'` or `'Dynamic'`. (Please note that the `Row_format` is a separate column, and ignore the contents of the `Create_options` column, which may contain the string `ROW_FORMAT`.) If the table in a tablespace uses neither of those features, the file uses the format supported by prior releases of InnoDB, now called file format “Antelope”. Then, the `Row_format` is one of `'Redundant'` or `'Compact'`.

The file format identifier is written as part of the tablespace flags (a 32-bit number) in the `*.ibd` file in the 4 bytes starting at position 54 of the file, most significant byte first. (The first byte of the file is byte zero.) On some systems, you can display these bytes in hexadecimal with the command `od -t x1 -j 54 -N 4 tablename.ibd`. If all bytes are zero, the tablespace uses the “Antelope” file format (which is the format used by the standard built-in InnoDB in MySQL up to version 5.1). Otherwise, the least significant bit should be set in the tablespace flags, and the file format identifier is written in the bits 5 through 11. (Divide the tablespace flags by 32 and take the remainder after dividing the integer part of the result by 128.)

4.6 Downgrading the File Format

Each InnoDB tablespace file (with a name matching `*.ibd`) is tagged with the file format used to create its table and indexes. The way to downgrade the tablespace is to re-create the table and its indexes. The easiest way to recreate a table and its indexes is to use the command:

```
ALTER TABLE t ROW_FORMAT=COMPACT;
```

on each table that you want to downgrade. The `COMPACT` row format uses the file format “Antelope”. It was introduced in MySQL 5.0.3.

4.7 Future InnoDB File Formats

The file format used by the standard built-in InnoDB in MySQL 5.1 is the “Antelope” format, and the new file format introduced with the InnoDB Plugin 1.0 is the “Barracuda” format. No definitive plans have been made to introduce new features that would require additional new file formats. However, the file format mechanism introduced with the InnoDB Plugin allows for further enhancements.

For the sake of completeness, these are the file format names that might be used for future file formats: Antelope, Barracuda, Cheetah, Dragon, Elk, Fox, Gazelle, Hornet, Impala, Jaguar, Kangaroo, Leopard, Moose, Nautilus, Ocelot, Porpoise, Quail, Rabbit, Shark, Tiger, Urchin, Viper, Whale, Xenops, Yak and Zebra. These file formats correspond to the internal identifiers 0..25.

Chapter 5 InnoDB Row Storage and Row Formats

Table of Contents

5.1 Storage of Variable-Length Columns	25
5.2 COMPACT and REDUNDANT Row Formats	25
5.3 DYNAMIC Row Format	25
5.4 Specifying a Table's Row Format	26

5.1 Storage of Variable-Length Columns

All data in InnoDB is stored in database pages comprising a B-tree index (the so-called clustered index or primary key index). The essential idea is that the nodes of the B-tree contain, for each primary key value (whether user-specified or generated or chosen by the system), the values of the remaining columns of the row as well as the key. In some other database systems, a clustered index is called an “index-organized table”. Secondary indexes in InnoDB are also B-trees, containing pairs of values of the index key and the value of the primary key, which acts as a pointer to the row in the clustered index.

There is an exception to this rule. Variable-length columns (such as **BLOB** and **VARCHAR**) that are too long to fit on a B-tree page are stored on separately allocated disk (“overflow”) pages. We call these “off-page columns”. The values of such columns are stored on singly-linked lists of overflow pages, and each such column has its own list of one or more overflow pages. In some cases, all or a prefix of the long column values is stored in the B-tree, to avoid wasting storage and eliminating the need to read a separate page.

The new “Barracuda” file format provides a new option (**KEY_BLOCK_SIZE**) to control how much column data is stored in the clustered index, and how much is placed on overflow pages.

5.2 **COMPACT** and **REDUNDANT** Row Formats

Previous versions of InnoDB used an unnamed file format (now called “Antelope”) for database files. With that format, tables were defined with **ROW_FORMAT=COMPACT** (or **ROW_FORMAT=REDUNDANT**) and InnoDB stored up to the first 768 bytes of variable-length columns (such as **BLOB** and **VARCHAR**) in the index record within the B-tree node, with the remainder stored on the overflow page(s).

To preserve compatibility with those prior versions, tables created with the InnoDB Plugin use the prefix format, unless one of **ROW_FORMAT=DYNAMIC** or **ROW_FORMAT=COMPRESSED** is specified (or implied) on the **CREATE TABLE** command.

With the “Antelope” file format, if the value of a column is not longer than 768 bytes, no overflow page is needed, and some savings in I/O may result, since the value is in the B-tree node. This works well for relatively short **BLOB**s, but may cause B-tree nodes to fill with data rather than key values, thereby reducing their efficiency. Tables with many **BLOB** columns could cause B-tree nodes to become too full of data, and contain too few rows, making the entire index less efficient than if the rows were shorter or if the column values were stored off-page.

5.3 **DYNAMIC** Row Format

When **innodb_file_format** is set to “Barracuda” and a table is created with **ROW_FORMAT=DYNAMIC** or **ROW_FORMAT=COMPRESSED**, long column values are stored fully off-page, and the clustered index record contains only a 20-byte pointer to the overflow page.

Whether any columns are stored off-page depends on the page size and the total size of the row. When the row is too long, InnoDB chooses the longest columns for off-page storage until the clustered index record fits on the B-tree page.

The `DYNAMIC` row format maintains the efficiency of storing the entire row in the index node if it fits (as do the `COMPACT` and `REDUNDANT` formats), but this new format avoids the problem of filling B-tree nodes with a large number of data bytes of long columns. The `DYNAMIC` format is predicated on the idea that if a portion of a long data value is stored off-page, it is usually most efficient to store all of the value off-page. With `DYNAMIC` format, shorter columns are likely to remain in the B-tree node, minimizing the number of overflow pages needed for any given row.

5.4 Specifying a Table's Row Format

The row format used for a table is specified with the `ROW_FORMAT` clause of the `CREATE TABLE` and `ALTER TABLE` commands. Note that `COMPRESSED` format implies `DYNAMIC` format. See [Section 3.2, “Enabling Compression for a Table”](#) for more details on the relationship between this clause and other clauses of these commands.

Chapter 6 InnoDB `INFORMATION_SCHEMA` Tables

Table of Contents

6.1 Overview of InnoDB Support in <code>INFORMATION_SCHEMA</code>	27
6.2 Information Schema Tables about Compression	27
6.2.1 <code>INNODB_CMP</code> and <code>INNODB_CMP_RESET</code>	27
6.2.2 <code>INNODB_CMPMEM</code> and <code>INNODB_CMPMEM_RESET</code>	28
6.2.3 Using the Compression Information Schema Tables	29
6.3 Information Schema Tables about Transactions	29
6.3.1 <code>INNODB_TRX</code>	29
6.3.2 <code>INNODB_LOCKS</code>	30
6.3.3 <code>INNODB_LOCK_WAITS</code>	31
6.3.4 Using the Transaction Information Schema Tables	31
6.4 Notes on Locking in InnoDB	36
6.4.1 Understanding InnoDB Locking	36
6.4.2 Rapidly Changing Internal Data	36
6.4.3 Possible Inconsistency with <code>PROCESSLIST</code>	37

6.1 Overview of InnoDB Support in `INFORMATION_SCHEMA`

The `INFORMATION_SCHEMA` tables `INNODB_BUFFER_PAGE`, `INNODB_BUFFER_PAGE_LRU`, `INNODB_BUFFER_POOL_STATS`, `INNODB_CMP`, `INNODB_CMP_RESET`, `INNODB_CMPMEM`, `INNODB_CMPMEM_RESET`, `INNODB_TRX`, `INNODB_LOCKS` and `INNODB_LOCK_WAITS` contain live information about the InnoDB buffer pool, compressed InnoDB tables, the compressed InnoDB buffer pool, all transactions currently executing inside InnoDB, the locks that transactions hold and those that are blocking transactions waiting for access to a resource (a table or row).

Note that the Information Schema tables are themselves plugins to the MySQL server. As such they need to be `INSTALLED` as described in [Chapter 9, *Installing the InnoDB Plugin*](#). If they are installed, but the InnoDB storage engine plugin is not installed, these tables appear to be empty.

Following is a description of the new Information Schema tables introduced in the InnoDB Plugin, and some examples of their use.

6.2 Information Schema Tables about Compression

Two new pairs of Information Schema tables provided by the InnoDB Plugin can give you some insight into how well compression is working overall. One pair of tables contains information about the number of compression operations and the amount of time spent performing compression. Another pair of tables contains information on the way memory is allocated for compression.

6.2.1 `INNODB_CMP` and `INNODB_CMP_RESET`

The tables `INNODB_CMP` and `INNODB_CMP_RESET` contain status information on the operations related to compressed tables, which are covered in [Chapter 3, *InnoDB Data Compression*](#). The compressed page size is in the column `PAGE_SIZE`.

These two tables have identical contents, but reading from `INNODB_CMP_RESET` resets the statistics on compression and uncompression operations. For example, if you archived the output of `INNODB_CMP_RESET` every 60 minutes, it would show the hourly statistics. If you never read

INNODB_CMP_RESET and monitored the output of INNODB_CMP instead, it would show the cumulated statistics since InnoDB was started.

Table 6.1 Columns of INNODB_CMP and INNODB_CMP_RESET

Column name	Description
PAGE_SIZE	Compressed page size in bytes.
COMPRESS_OPS	Number of times a B-tree page of the size PAGE_SIZE has been compressed. Pages are compressed whenever an empty page is created or the space for the uncompressed modification log runs out.
COMPRESS_OPS_OK	Number of times a B-tree page of the size PAGE_SIZE has been successfully compressed. This count should never exceed COMPRESS_OPS.
COMPRESS_TIME	Total time in seconds spent in attempts to compress B-tree pages of the size PAGE_SIZE.
UNCOMPRESS_OPS	Number of times a B-tree page of the size PAGE_SIZE has been uncompressed. B-tree pages are uncompressed whenever compression fails or at first access when the uncompressed page does not exist in the buffer pool.
UNCOMPRESS_TIME	Total time in seconds spent in uncompressing B-tree pages of the size PAGE_SIZE.

6.2.2 INNODB_CMPMEM and INNODB_CMPMEM_RESET

You may consider the tables INNODB_CMPMEM and INNODB_CMPMEM_RESET as the status information on the compressed pages that reside in the buffer pool. Please consult [Chapter 3, InnoDB Data Compression](#) for further information on compressed tables and the use of the buffer pool. The tables INNODB_CMP and INNODB_CMP_RESET should provide more useful statistics on compression.

The InnoDB Plugin uses a so-called “buddy allocator” system to manage memory allocated to pages of various sizes, from 1KB to 16KB. Each row of the two tables described here corresponds to a single page size.

These two tables have identical contents, but reading from INNODB_CMPMEM_RESET resets the statistics on relocation operations. For example, if every 60 minutes you archived the output of INNODB_CMPMEM_RESET, it would show the hourly statistics. If you never read INNODB_CMPMEM_RESET and monitored the output of INNODB_CMPMEM instead, it would show the cumulated statistics since InnoDB was started.

Table 6.2 Columns of INNODB_CMPMEM and INNODB_CMPMEM_RESET

Column name	Description
PAGE_SIZE	Block size in bytes. Each record of this table describes blocks of this size.
PAGES_USED	Number of blocks of the size PAGE_SIZE that are currently in use.
PAGES_FREE	Number of blocks of the size PAGE_SIZE that are currently available for allocation. This column shows the external fragmentation in the memory pool. Ideally, these numbers should be at most 1.
RELOCATION_OPS	Number of times a block of the size PAGE_SIZE has been relocated. The buddy system can relocate the allocated “buddy neighbor” of a freed block when it tries to form a bigger freed block. Reading from the table INNODB_CMPMEM_RESET resets this count.
RELOCATION_TIME	Total time in microseconds spent in relocating blocks of the size PAGE_SIZE. Reading from the table INNODB_CMPMEM_RESET resets this count.

6.2.3 Using the Compression Information Schema Tables

Example 6.1 Using the Compression Information Schema Tables

The following is sample output from a database that contains compressed tables (see [Chapter 3, InnoDB Data Compression](#), `INNODB_CMP`, and `INNODB_CMPMEM`).

The following table shows the contents of `INFORMATION_SCHEMA.INNODB_CMP` under a light workload. The only compressed page size that the buffer pool contains is 8K. Compressing or uncompressing pages has consumed less than a second since the time the statistics were reset, because the columns `COMPRESS_TIME` and `UNCOMPRESS_TIME` are zero.

page size	compress ops	compress ops ok	compress time	uncompress ops	uncompress time
1024	0	0	0	0	0
2048	0	0	0	0	0
4096	0	0	0	0	0
8192	1048	921	0	61	0
16384	0	0	0	0	0

According to `INNODB_CMPMEM`, there are 6169 compressed 8KB pages in the buffer pool.

The following table shows the contents of `INFORMATION_SCHEMA.INNODB_CMPMEM` under light load. We can see that some memory is unusable due to fragmentation of the InnoDB memory allocator for compressed pages: `SUM(PAGE_SIZE*PAGES_FREE)=6784`. This is because small memory allocation requests are fulfilled by splitting bigger blocks, starting from the 16K blocks that are allocated from the main buffer pool, using the buddy allocation system. The fragmentation is this low because some allocated blocks have been relocated (copied) to form bigger adjacent free blocks. This copying of `SUM(PAGE_SIZE*RELOCATION_OPS)` bytes has consumed less than a second (`SUM(RELOCATION_TIME)=0`).

page size	pages used	pages free	relocation ops	relocation time
1024	0	0	0	0
2048	0	1	0	0
4096	0	1	0	0
8192	6169	0	5	0
16384	0	0	0	0

6.3 Information Schema Tables about Transactions

Three new Information Schema tables introduced in the InnoDB Plugin make it much easier to monitor transactions and diagnose possible locking problems. The three tables are `INNODB_TRX`, `INNODB_LOCKS` and `INNODB_LOCK_WAITS`.

6.3.1 `INNODB_TRX`

Contains information about every transaction currently executing inside InnoDB, including whether the transaction is waiting for a lock, when the transaction started, and the particular SQL statement the transaction is executing.

Table 6.3 INNODB_TRX Columns

Column name	Description
TRX_ID	Unique transaction ID number, internal to InnoDB.
TRX_WEIGHT	The weight of a transaction, reflecting (but not necessarily the exact count of) the number of rows altered and the number of rows locked by the transaction. To resolve a deadlock, InnoDB selects the transaction with the smallest weight as the “victim” to rollback. Transactions that have changed non-transactional tables are considered heavier than others, regardless of the number of altered and locked rows.
TRX_STATE	Transaction execution state. One of 'RUNNING', 'LOCK WAIT', 'ROLLING BACK' or 'COMMITTING'.
TRX_STARTED	Transaction start time; the transaction is created by executing a transactional query.
TRX_REQUESTED_LOCK_ID	ID of the lock the transaction is currently waiting for (if TRX_STATE is 'LOCK WAIT', otherwise NULL). Details about the lock can be found by joining with INNODB_LOCKS on LOCK_ID.
TRX_WAIT_STARTED	Time when the transaction started waiting on the lock (if TRX_STATE is 'LOCK WAIT', otherwise NULL).
TRX_MYSQL_THREAD_ID	MySQL thread ID. Can be used for joining with PROCESSLIST on ID. See Section 6.4.3, “Possible Inconsistency with PROCESSLIST”.
TRX_QUERY	The SQL query that is being executed by the transaction.

6.3.2 INNODB_LOCKS

Each transaction in InnoDB that is waiting for another transaction to release a lock (`INNODB_TRX.TRX_STATE='LOCK WAIT'`) is blocked by exactly one “blocking lock request”. That blocking lock request is for a row or table lock held by another transaction in an incompatible mode. The waiting or blocked transaction cannot proceed until the other transaction commits or rolls back, thereby releasing the requested lock. For every blocked transaction, `INNODB_LOCKS` contains one row that describes each lock the transaction has requested, and for which it is waiting. `INNODB_LOCKS` also contains one row for each lock that is blocking another transaction, whatever the state of the transaction that holds the lock ('RUNNING', 'LOCK WAIT', 'ROLLING BACK' or 'COMMITTING'). The lock that is blocking a transaction is always held in a mode (read vs. write, shared vs. exclusive) incompatible with the mode of requested lock.

Table 6.4 INNODB_LOCKS Columns

Column name	Description
LOCK_ID	Unique lock ID number, internal to InnoDB. Should be treated as an opaque string. Although <code>LOCK_ID</code> currently contains <code>TRX_ID</code> , the format of the data in <code>LOCK_ID</code> is not guaranteed to remain the same in future releases. You should not write programs that parse the <code>LOCK_ID</code> value.
LOCK_TRX_ID	ID of the transaction holding this lock. Details about the transaction can be found by joining with <code>INNODB_TRX</code> on <code>TRX_ID</code> .
LOCK_MODE	Mode of the lock. One of 'S', 'X', 'IS', 'IX', 'S,GAP', 'X,GAP', 'IS,GAP', 'IX,GAP', or 'AUTO_INC' for shared, exclusive, intention shared, intention exclusive row locks, shared and exclusive gap locks, intention shared and intention exclusive gap locks, and auto-increment table level lock, respectively. Refer to the

Column name	Description
	sections InnoDB Lock Modes and The InnoDB Transaction Model and Locking of the MySQL Manual for information on InnoDB locking.
LOCK_TYPE	Type of the lock. One of 'RECORD' or 'TABLE' for record (row) level or table level locks, respectively.
LOCK_TABLE	Name of the table that has been locked or contains locked records.
LOCK_INDEX	Name of the index if LOCK_TYPE= 'RECORD' , otherwise NULL.
LOCK_SPACE	Tablespace ID of the locked record if LOCK_TYPE= 'RECORD' , otherwise NULL.
LOCK_PAGE	Page number of the locked record if LOCK_TYPE= 'RECORD' , otherwise NULL.
LOCK_REC	Heap number of the locked record within the page if LOCK_TYPE= 'RECORD' , otherwise NULL.
LOCK_DATA	Primary key of the locked record if LOCK_TYPE= 'RECORD' , otherwise NULL. This column contains the value(s) of the primary key column(s) in the locked row, formatted as a valid SQL string (ready to be copied to SQL commands). If there is no primary key then the InnoDB internal unique row ID number is used. When the page containing the locked record is not in the buffer pool (in the case that it was paged out to disk while the lock was held), InnoDB does not fetch the page from disk, to avoid unnecessary disk operations. Instead, LOCK_DATA is set to NULL.

6.3.3 INNODB_LOCK_WAITS

Using this table, you can tell which transactions are waiting for a given lock, or for which lock a given transaction is waiting. This table contains one or more rows for each *blocked* transaction, indicating the lock it has requested and any locks that are blocking that request. The [REQUESTED_LOCK_ID](#) refers to the lock that a transaction is requesting, and the [BLOCKING_LOCK_ID](#) refers to the lock (held by another transaction) that is preventing the first transaction from proceeding. For any given blocked transaction, all rows in [INNODB_LOCK_WAITS](#) have the same value for [REQUESTED_LOCK_ID](#) and different values for [BLOCKING_LOCK_ID](#).

Table 6.5 [INNODB_LOCK_WAITS](#) Columns

Column name	Description
REQUESTING_TRX_ID	ID of the requesting transaction.
REQUESTED_LOCK_ID	ID of the lock for which a transaction is waiting. Details about the lock can be found by joining with INNODB_LOCKS on LOCK_ID .
BLOCKING_TRX_ID	ID of the blocking transaction.
BLOCKING_LOCK_ID	ID of a lock held by a transaction blocking another transaction from proceeding. Details about the lock can be found by joining with INNODB_LOCKS on LOCK_ID .

6.3.4 Using the Transaction Information Schema Tables

Example 6.2 Identifying Blocking Transactions

It is sometimes helpful to be able to identify which transaction is blocking another. You can use the Information Schema tables to find out which transaction is waiting for another, and which resource is being requested.

Suppose you have the following scenario, with three users running concurrently. Each user (or session) corresponds to a MySQL thread, and executes one transaction after another. Consider the state of

the system when these users have issued the following commands, but none has yet committed its transaction:

- User A:

```
BEGIN;
SELECT a FROM t FOR UPDATE;
SELECT SLEEP(100);
```

- User B:

```
SELECT b FROM t FOR UPDATE;
```

- User C:

```
SELECT c FROM t FOR UPDATE;
```

In this scenario, you can use this query to see who is waiting for whom:

```
SELECT r.trx_id waiting_trx_id,
       r.trx_mysql_thread_id waiting_thread,
       r.trx_query waiting_query,
       b.trx_id blocking_trx_id,
       b.trx_mysql_thread_id blocking_thread,
       b.trx_query blocking_query
FROM   information_schema.innodb_lock_waits w
INNER JOIN information_schema.innodb_trx b ON
        b.trx_id = w.blocking_trx_id
INNER JOIN information_schema.innodb_trx r ON
        r.trx_id = w.requesting_trx_id;
```

waiting trx id	waiting thread	waiting query	blocking trx id	blocking thread	blocking query
A4	6	SELECT b FROM t FOR UPDATE	A3	5	SELECT SLEEP(100)
A5	7	SELECT c FROM t FOR UPDATE	A3	5	SELECT SLEEP(100)
A5	7	SELECT c FROM t FOR UPDATE	A4	6	SELECT b FROM t FOR UPDATE

In the above result, you can identify users by the “waiting query” or “blocking query”. As you can see:

- User B (trx id 'A4', thread 6) and User C (trx id 'A5', thread 7) are both waiting for User A (trx id 'A3', thread 5).
- User C is waiting for User B as well as User A.

You can see the underlying data in the tables [INNODB_TRX](#), [INNODB_LOCKS](#), and [INNODB_LOCK_WAITS](#).

The following table shows some sample contents of INFORMATION_SCHEMA.INNODB_TRX.

trx id	trx state	trx started	trx requested lock id	trx wait started	trx weight	trx mysql thread id	trx query
A3	RUNNING	2008-01-15 16:44:54	NULL	NULL	2	5	SELECT SLEEP(100)

trx id	trx state	trx started	trx requested lock id	trx wait started	trx weight	trx mysql thread id	trx query
A4	LOCK WAIT	2008-01-15 16:45:09	A4:1:3:2	2008-01-15 16:45:09	2	6	SELECT b FROM t FOR UPDATE
A5	LOCK WAIT	2008-01-15 16:45:14	A5:1:3:2	2008-01-15 16:45:14	2	7	SELECT c FROM t FOR UPDATE

The following table shows some sample contents of `INFORMATION_SCHEMA.INNODB_LOCKS`.

lock id	lock trx id	lock mode	lock type	lock table	lock index	lock space	lock page	lock rec	lock data
A3:1:3:2	A3	X	RECORD	`test`.`t`	`PRIMARY`	1	3	2	0x0200
A4:1:3:2	A4	X	RECORD	`test`.`t`	`PRIMARY`	1	3	2	0x0200
A5:1:3:2	A5	X	RECORD	`test`.`t`	`PRIMARY`	1	3	2	0x0200

The following table shows some sample contents of `INFORMATION_SCHEMA.INNODB_LOCK_WAITS`.

requesting trx id	requested lock id	blocking trx id	blocking lock id
A4	A4:1:3:2	A3	A3:1:3:2
A5	A5:1:3:2	A3	A3:1:3:2
A5	A5:1:3:2	A4	A4:1:3:2

Example 6.3 More Complex Example of Transaction Data in Information Schema Tables

Sometimes you would like to correlate the internal InnoDB locking information with session-level information maintained by MySQL. For example, you might like to know, for a given InnoDB transaction ID, the corresponding MySQL session ID and name of the user that may be holding a lock, and thus blocking another transaction.

The following output from the `INFORMATION_SCHEMA` tables is taken from a somewhat loaded system.

As can be seen in the following tables, there are several transactions running.

The following `INNODB_LOCKS` and `INNODB_LOCK_WAITS` tables shows that:

- Transaction `77F` (executing an `INSERT`) is waiting for transactions `77E`, `77D` and `77B` to commit.
- Transaction `77E` (executing an `INSERT`) is waiting for transactions `77D` and `77B` to commit.
- Transaction `77D` (executing an `INSERT`) is waiting for transaction `77B` to commit.
- Transaction `77B` (executing an `INSERT`) is waiting for transaction `77A` to commit.
- Transaction `77A` is running, currently executing `SELECT`.
- Transaction `E56` (executing an `INSERT`) is waiting for transaction `E55` to commit.
- Transaction `E55` (executing an `INSERT`) is waiting for transaction `19C` to commit.
- Transaction `19C` is running, currently executing an `INSERT`.

Note that there may be an inconsistency between queries shown in the two tables `INNODB_TRX.TRX_QUERY` and `PROCESSLIST.INFO`. The current transaction ID for a thread, and the

query being executed in that transaction, may be different in these two tables for any given thread. See [Section 6.4.3, “Possible Inconsistency with PROCESSLIST”](#) for an explanation.

The following table shows the contents of `INFORMATION_SCHEMA.PROCESSLIST` in a system running a heavy workload.

ID	USER	HOST	DB	COMMAND	TIME	STATE	INFO
384	root	localhost	test	Query	10	update	insert into t2 values ...
257	root	localhost	test	Query	3	update	insert into t2 values ...
130	root	localhost	test	Query	0	update	insert into t2 values ...
61	root	localhost	test	Query	1	update	insert into t2 values ...
8	root	localhost	test	Query	1	update	insert into t2 values ...
4	root	localhost	test	Query	0	preparing	SELECT * FROM processlist
2	root	localhost	test	Sleep	566		NULL

The following table shows the contents of `INFORMATION_SCHEMA.INNODB_TRX` in a system running a heavy workload.

trx id	trx state	trx started	trx requested lock id	trx wait started	trx weight	trx mysql thread id	trx query
77F	LOCK WAIT	2008-01-15 13:10:16	77F:806	2008-01-15 13:10:16	1	876	insert into t09 (D, B, C) values ...
77E	LOCK WAIT	2008-01-15 13:10:16	77E:806	2008-01-15 13:10:16	1	875	insert into t09 (D, B, C) values ...
77D	LOCK WAIT	2008-01-15 13:10:16	77D:806	2008-01-15 13:10:16	1	874	insert into t09 (D, B, C) values ...
77B	LOCK WAIT	2008-01-15 13:10:16	77B:733:12:1	2008-01-15 13:10:16	4	873	insert into t09 (D, B, C) values ...
77A	RUNNING	2008-01-15 13:10:16	NULL	NULL	4	872	select b, c from t09 where ...
E56	LOCK WAIT	2008-01-15 13:10:06	E56:743:6:2	2008-01-15 13:10:06	5	384	insert into t2 values ...
E55	LOCK WAIT	2008-01-15 13:10:06	E55:743:38:2	2008-01-15 13:10:13	965	257	insert into t2 values ...
19C	RUNNING	2008-01-15 13:09:10	NULL	NULL	2900	130	insert into t2 values ...

trx id	trx state	trx started	trx requested lock id	trx wait started	trx weight	trx mysql thread id	trx query
<i>E15</i>	RUNNING	2008-01-15 13:08:59	NULL	NULL	5395	61	insert into t2 values ...
<i>51D</i>	RUNNING	2008-01-15 13:08:47	NULL	NULL	9807	8	insert into t2 values ...

The following table shows the contents of `INFORMATION_SCHEMA.INNODB_LOCK_WAITS` in a loaded system

requesting trx id	requested lock id	blocking trx id	blocking lock id
<i>77F</i>	<i>77F</i> :806	<i>77E</i>	<i>77E</i> :806
<i>77F</i>	<i>77F</i> :806	<i>77D</i>	<i>77D</i> :806
<i>77F</i>	<i>77F</i> :806	<i>77B</i>	<i>77B</i> :806
<i>77E</i>	<i>77E</i> :806	<i>77D</i>	<i>77D</i> :806
<i>77E</i>	<i>77E</i> :806	<i>77B</i>	<i>77B</i> :806
<i>77D</i>	<i>77D</i> :806	<i>77B</i>	<i>77B</i> :806
<i>77B</i>	<i>77B</i> :733:12:1	<i>77A</i>	<i>77A</i> :733:12:1
<i>E56</i>	<i>E56</i> :743:6:2	<i>E55</i>	<i>E55</i> :743:6:2
<i>E55</i>	<i>E55</i> :743:38:2	<i>19C</i>	<i>19C</i> :743:38:2

The following table shows the contents of `INFORMATION_SCHEMA.INNODB_LOCKS` in a system running a heavy workload.

lock id	lock trx id	lock mode	lock type	lock table	lock index	lock space	lock page	lock rec	lock data
<i>77F</i> :806	<i>77F</i>	AUTO_INC	TABLE	<code>`test`.`t09`</code>	NULL	NULL	NULL	NULL	NULL
<i>77E</i> :806	<i>77E</i>	AUTO_INC	TABLE	<code>`test`.`t09`</code>	NULL	NULL	NULL	NULL	NULL
<i>77D</i> :806	<i>77D</i>	AUTO_INC	TABLE	<code>`test`.`t09`</code>	NULL	NULL	NULL	NULL	NULL
<i>77B</i> :806	<i>77B</i>	AUTO_INC	TABLE	<code>`test`.`t09`</code>	NULL	NULL	NULL	NULL	NULL
<i>77B</i> :733:12:1	<i>77B</i>	X	RECORD	<code>`test`.`t09`</code>	<code>`PRIMARY`</code>	733	12	1	supremum pseudo-record
<i>77A</i> :733:12:1	<i>77A</i>	X	RECORD	<code>`test`.`t09`</code>	<code>`PRIMARY`</code>	733	12	1	supremum pseudo-record
<i>E56</i> :743:6:2	<i>E56</i>	S	RECORD	<code>`test`.`t2`</code>	<code>`PRIMARY`</code>	743	6	2	0, 0
<i>E55</i> :743:6:2	<i>E55</i>	X	RECORD	<code>`test`.`t2`</code>	<code>`PRIMARY`</code>	743	6	2	0, 0

lock id	lock trx id	lock mode	lock type	lock table	lock index	lock space	lock page	lock rec	lock data
<i>E55</i> :743 :38:2	<i>E55</i>	S	RECORD	`test` .`t2`	`PRIMARY`	743	38	2	1922, 1922
<i>19C</i> :743 :38:2	<i>19C</i>	X	RECORD	`test` .`t2`	`PRIMARY`	743	38	2	1922, 1922

6.4 Notes on Locking in InnoDB

6.4.1 Understanding InnoDB Locking

When a transaction updates a row in a table, or locks it with `SELECT FOR UPDATE`, InnoDB establishes a list or queue of locks on that row. Similarly, InnoDB maintains a list of locks on a table for table-level locks transactions hold. If a second transaction wants to update a row or lock a table already locked by a prior transaction in an incompatible mode, InnoDB adds a lock request for the row to the corresponding queue. For a lock to be acquired by a transaction, all incompatible lock requests previously entered into the lock queue for that row or table must be removed (the transactions holding or requesting those locks either commit or rollback).

A transaction may have any number of lock requests for different rows or tables. At any given time, a transaction may be requesting a lock that is held by another transaction, in which case it is blocked by that other transaction. The requesting transaction must wait for the transaction that holds the blocking lock to commit or rollback. If a transaction is not waiting for a lock, it is in the `'RUNNING'` state. If a transaction is waiting for a lock, it is in the `'LOCK WAIT'` state.

The table `INNODB_LOCKS` holds one or more row for each `'LOCK WAIT'` transaction, indicating any lock requests that are preventing its progress. This table also contains one row describing each lock in a queue of locks pending for a given row or table. The table `INNODB_LOCK_WAITS` shows which locks already held by a transaction are blocking locks requested by other transactions.

6.4.2 Rapidly Changing Internal Data

The data exposed by the transaction and locking tables represent a glimpse into fast-changing data. This is not like other (user) tables, where the data only changes when application-initiated updates occur. The underlying data is internal system-managed data, and can change very quickly.

For performance reasons, and to minimize the chance of misleading `JOINS` between the `INFORMATION_SCHEMA` tables, InnoDB collects the required transaction and locking information into an intermediate buffer whenever a `SELECT` on any of the tables is issued. This buffer is refreshed only if more than 0.1 seconds has elapsed since the last time the buffer was read. The data needed to fill the three tables is fetched atomically and consistently and is saved in this global internal buffer, forming a point-in-time “snapshot”. If multiple table accesses occur within 0.1 seconds (as they almost certainly do when MySQL processes a join among these tables), then the same snapshot is used to satisfy the query.

A correct result is returned when you `JOIN` any of these tables together in a single query, because the data for the three tables comes from the same snapshot. Because the buffer is not refreshed with every query of any of these tables, if you issue separate queries against these tables within a tenth of a second, the results are the same from query to query. On the other hand, two separate queries of the same or different tables issued more than a tenth of a second apart may see different results, since the data come from different snapshots.

Because InnoDB must temporarily stall while the transaction and locking data is collected, too frequent queries of these tables can negatively impact performance as seen by other users.

As these tables contain sensitive information (at least `INNODB_LOCKS.LOCK_DATA` and `INNODB_TRX.TRX_QUERY`), for security reasons, only the users with the `PROCESS` privilege are allowed to `SELECT` from them.

6.4.3 Possible Inconsistency with `PROCESSLIST`

As just described, while the transaction and locking data is correct and consistent when these `INFORMATION_SCHEMA` tables are populated, the underlying data changes so fast that similar glimpses at other, similarly fast-changing data, may not be in sync. Thus, you should be careful in comparing the data in the InnoDB transaction and locking tables with that in the [The `INFORMATION_SCHEMA` `PROCESSLIST` Table](#). The data from the `PROCESSLIST` table does not come from the same snapshot as the data about locking and transactions. Even if you issue a single `SELECT` (joining `INNODB_TRX` and `PROCESSLIST`, for example), the content of those tables is generally not consistent. `INNODB_TRX` may reference rows that are not present in `PROCESSLIST` or the currently executing SQL query of a transaction, shown in `INNODB_TRX.TRX_QUERY` may be different from the one in `PROCESSLIST.INFO`. The query in `INNODB_TRX` is always consistent with the rest of `INNODB_TRX`, `INNODB_LOCKS` and `INNODB_LOCK_WAITS` when the data comes from the same snapshot.

Chapter 7 Performance and Scalability Enhancements

Table of Contents

7.1 Overview	39
7.2 Faster Locking for Improved Scalability	40
7.3 Using Operating System Memory Allocators	40
7.4 Controlling InnoDB Insert Buffering	41
7.5 Controlling Adaptive Hash Indexing	42
7.6 Changes Regarding Thread Concurrency	42
7.7 Changes in the Read Ahead Algorithm	43
7.8 Multiple Background I/O Threads	44
7.9 Group Commit	44
7.10 Controlling the Master Thread I/O Rate	45
7.11 Controlling the Flushing Rate of Dirty Pages	45
7.12 Using the PAUSE instruction in InnoDB spin loops	46
7.13 Control of Spin Lock Polling	46
7.14 Making Buffer Cache Scan Resistant	46
7.14.1 Guidelines for <code>innodb_old_blocks_pct</code> and <code>innodb_old_blocks_time</code>	48
7.15 Improvements to Crash Recovery Performance	48

7.1 Overview

InnoDB has always been highly efficient, and includes several unique architectural elements to assure high performance and scalability. InnoDB Plugin 1.0.8 includes several new features that take better advantage of recent advances in operating systems and hardware platforms, such as multi-core processors and improved memory allocation systems. In addition, this release permits you to better control the use of some InnoDB internal subsystems to achieve the best performance with your workload.

InnoDB Plugin 1.0.8 includes new capabilities in these areas:

- [Section 7.2, “Faster Locking for Improved Scalability”](#)
- [Section 7.3, “Using Operating System Memory Allocators”](#)
- [Section 7.4, “Controlling InnoDB Insert Buffering”](#)
- [Section 7.5, “Controlling Adaptive Hash Indexing”](#)
- [Section 7.6, “Changes Regarding Thread Concurrency”](#)
- [Section 7.7, “Changes in the Read Ahead Algorithm”](#)
- [Section 7.8, “Multiple Background I/O Threads”](#)
- [Section 7.9, “Group Commit”](#)
- [Section 7.10, “Controlling the Master Thread I/O Rate”](#)
- [Section 7.11, “Controlling the Flushing Rate of Dirty Pages”](#)
- [Section 7.12, “Using the PAUSE instruction in InnoDB spin loops”](#)
- [Section 7.13, “Control of Spin Lock Polling”](#)

- [Section 7.14, “Making Buffer Cache Scan Resistant”](#)
- [Section 7.15, “Improvements to Crash Recovery Performance”](#)

7.2 Faster Locking for Improved Scalability

In MySQL and InnoDB, multiple threads of execution access shared data structures. InnoDB synchronizes these accesses with its own implementation of mutexes and read/write locks. InnoDB has historically protected the internal state of a read/write lock with an InnoDB mutex. On Unix and Linux platforms, the internal state of an InnoDB mutex is protected by a Pthreads mutex, as in IEEE Std 1003.1c (POSIX.1c).

On many platforms, there is a more efficient way to implement mutexes and read/write locks. Atomic operations can often be used to synchronize the actions of multiple threads more efficiently than Pthreads. Each operation to acquire or release a lock can be done in fewer CPU instructions, and thus result in less wasted time when threads are contending for access to shared data structures. This in turn means greater scalability on multi-core platforms.

Beginning with InnoDB Plugin 1.0.3, InnoDB implements mutexes and read/write locks with the [built-in functions provided by the GNU Compiler Collection \(GCC\) for atomic memory access](#) instead of using the Pthreads approach previously used. More specifically, an InnoDB Plugin that is compiled with GCC version 4.1.2 or later will use the atomic builtins instead of a `pthread_mutex_t` to implement InnoDB mutexes and read/write locks.

On 32-bit Microsoft Windows, InnoDB has implemented mutexes (but not read/write locks) with hand-written assembler instructions. Beginning with Microsoft Windows 2000, it is possible to use functions for [Interlocked Variable Access](#) that are similar to the built-in functions provided by GCC. Beginning with InnoDB Plugin 1.0.4, InnoDB makes use of the Interlocked functions on Windows. Unlike the old hand-written assembler code, the new implementation supports read/write locks and 64-bit platforms.

Solaris 10 introduced library functions for atomic operations. Beginning with InnoDB Plugin 1.0.4, when InnoDB is compiled on Solaris 10 with a compiler that does not support the [built-in functions provided by the GNU Compiler Collection \(GCC\) for atomic memory access](#), the library functions will be used.

This change improves the scalability of InnoDB on multi-core systems. Note that the user does not have to set any particular parameter or option to take advantage of this new feature. This feature is enabled out-of-the-box on the platforms where it is supported. On platforms where the GCC, Windows, or Solaris functions for atomic memory access are not available, InnoDB will use the traditional Pthreads method of implementing mutexes and read/write locks.

When MySQL starts, InnoDB will write a message to the log file indicating whether atomic memory access will be used for mutexes, for mutexes and read/write locks, or neither. If suitable tools are used to build the InnoDB Plugin and the target CPU supports the atomic operations required, InnoDB will use the built-in functions for mutexing. If, in addition, the compare-and-swap operation can be used on thread identifiers (`pthread_t`), then InnoDB will use the instructions for read-write locks as well.

Note: If you are building from source, see [Section 9.4.1, “Building the InnoDB Plugin on Linux or Unix” \[66\]](#) to ensure that your build process properly takes advantage of your platform capabilities.

7.3 Using Operating System Memory Allocators

When InnoDB was developed, the memory allocators supplied with operating systems and run-time libraries were often lacking in performance and scalability. At that time, there were no memory allocator libraries tuned for multi-core CPUs. Therefore, InnoDB implemented its own memory allocator in the `mem` subsystem. This allocator is guarded by a single mutex, which may become a bottleneck. InnoDB also

implements a wrapper interface around the system allocator (`malloc` and `free`) that is likewise guarded by a single mutex.

Today, as multi-core systems have become more widely available, and as operating systems have matured, significant improvements have been made in the memory allocators provided with operating systems. New memory allocators perform better and are more scalable than they were in the past. The leading high-performance memory allocators include `Hoard`, `libumem`, `mtmalloc`, `ptmalloc`, `tbbmalloc`, and `TCMalloc`. Most workloads, especially those where memory is frequently allocated and released (such as multi-table joins) will benefit from using a more highly tuned memory allocator as opposed to the internal, InnoDB-specific memory allocator.

Beginning with InnoDB Plugin 1.0.3, you control whether InnoDB uses its own memory allocator or an allocator of the operating system, by setting the value of the new system configuration parameter `innodb_use_sys_malloc` in the MySQL option file (`my.cnf` or `my.ini`). If set to `ON` or `1` (the default), InnoDB will use the `malloc` and `free` functions of the underlying system rather than manage memory pools itself. This parameter is not dynamic, and takes effect only when the system is started. To continue to use the InnoDB memory allocator in InnoDB Plugin, you will have to set `innodb_use_sys_malloc` to `0`.

Note that when the InnoDB memory allocator is disabled, InnoDB will ignore the value of the parameter `innodb_additional_mem_pool_size`. The InnoDB memory allocator uses an additional memory pool for satisfying allocation requests without having to fall back to the system memory allocator. When the InnoDB memory allocator is disabled, all such allocation requests will be fulfilled by the system memory allocator.

Furthermore, since InnoDB cannot track all memory use when the system memory allocator is used (`innodb_use_sys_malloc` is `ON`), the section “BUFFER POOL AND MEMORY” in the output of the `SHOW ENGINE INNODB STATUS` command will only include the buffer pool statistics in the “Total memory allocated”. Any memory allocated using the `mem` subsystem or using `ut_malloc` will be excluded.

On Unix-like systems that use dynamic linking, replacing the memory allocator may be as easy as making the environment variable `LD_PRELOAD` or `LD_LIBRARY_PATH` point to the dynamic library that implements the allocator. On other systems, some relinking may be necessary. Please refer to the documentation of the memory allocator library of your choice.

7.4 Controlling InnoDB Insert Buffering

When `INSERTs` are done to a table, often the values of indexed columns (particularly the values of secondary keys) are not in sorted order. This means that the inserts of such values into secondary B-tree indexes is “random”, and this can cause excessive I/O if the entire index does not fit in memory. InnoDB has an insert buffer that caches changes to secondary index entries when the relevant page is not in the buffer pool, thus avoiding I/O operations by not reading in the page from the disk. The buffered changes are written into a special insert buffer tree and are subsequently merged when the page is loaded to the buffer pool. The InnoDB main thread merges buffered changes when the server is nearly idle.

Usually, this process will result in fewer disk reads and writes, especially during bulk inserts. However, the insert buffer tree will occupy a part of the buffer pool. If the working set almost fits in the buffer pool, it may be useful to disable insert buffering. If the working set entirely fits in the buffer pool, insert buffering will not be used anyway, because the index would exist in memory.

Beginning with InnoDB Plugin 1.0.3, you can control whether InnoDB performs insert buffering with the system configuration parameter `innodb_change_buffering`. The allowed values of `innodb_change_buffering` are `none` (do not buffer any operations) and `inserts` (buffer insert operations, the default). You can set the value of this parameter in the MySQL option file (`my.cnf` or `my.ini`) or change it dynamically with the `SET GLOBAL` command, which requires the `SUPER` privilege.

Changing the setting affects the buffering of new operations; the merging of already buffered entries is not affected.

7.5 Controlling Adaptive Hash Indexing

If a table fits almost entirely in main memory, the fastest way to perform queries on it is to use hash indexes rather than B-tree lookups. InnoDB monitors searches on each index defined for a table. If it notices that certain index values are being accessed frequently, it automatically builds an in-memory hash table for that index. Based on the pattern of searches that InnoDB observes, it will build a hash index using a prefix of the index key. The prefix of the key can be any length, and it may be that only a subset of the values in the B-tree will appear in the hash index. InnoDB builds hash indexes on demand for those pages of the index that are often accessed.

The adaptive hash index mechanism allows InnoDB to take advantage of large amounts of memory, something typically done only by database systems specifically designed for databases that reside entirely in memory. Normally, the automatic building and use of adaptive hash indexes will improve performance. However, sometimes, the read/write lock that guards access to the adaptive hash index may become a source of contention under heavy workloads, such as multiple concurrent joins.

You can monitor the use of the adaptive hash index and the contention for its use in the “SEMAPHORES” section of the output of the `SHOW ENGINE INNODB STATUS` command. If you see many threads waiting on an RW-latch created in `btr0sea.c`, then it might be useful to disable adaptive hash indexing.

The configuration parameter `innodb_adaptive_hash_index` can be set to disable or enable the adaptive hash index. See [Section 8.3.4, “Dynamically Changing `innodb_adaptive_hash_index`”](#) for details.

7.6 Changes Regarding Thread Concurrency

InnoDB uses operating system threads to process requests from user transactions. (Transactions may issue many requests to InnoDB before they commit or roll back.) On today's modern operating systems and servers with multi-core processors, where context switching is efficient, most workloads will run well without any limit on the number of concurrent threads. Thanks to several scalability improvements in InnoDB Plugin 1.0.3, and further changes in release 1.0.4, there should be less need to artificially limit the number of concurrently executing threads inside InnoDB.

However, for some situations, it may be helpful to minimize context switching between threads. InnoDB can use a number of techniques to limit the number of concurrently executing operating system threads (and thus the number of requests that are processed at any one time). When InnoDB receives a new request from a user session, if the number of threads concurrently executing is at a pre-defined limit, the new request will sleep for a short time before it tries again. A request that cannot be rescheduled after the sleep is put in a first-in/first-out queue and eventually will be processed. Threads waiting for locks are not counted in the number of concurrently executing threads.

The limit on the number of concurrent threads is given by the settable global variable `innodb_thread_concurrency`. Once the number of executing threads reaches this limit, additional threads will sleep for a number of microseconds, set by the system configuration parameter `innodb_thread_sleep_delay`, before being placed into the queue.

The default value for `innodb_thread_concurrency` and the implied default limit on the number of concurrent threads has been changed in various releases of MySQL and the InnoDB Plugin. Starting with InnoDB Plugin 1.0.3, the default value of `innodb_thread_concurrency` is 0, so that by default there is no limit on the number of concurrently executing threads, as shown in [Table 7.1, “Changes to `innodb_thread_concurrency`”](#).

Table 7.1 Changes to `innodb_thread_concurrency`

InnoDB Version	MySQL Version	Default value	Default limit of concurrent threads	Value to allow unlimited threads
Built-in	Earlier than 5.1.11	20	No limit	20 or higher
Built-in	5.1.11 and newer	8	8	0
InnoDB Plugin before 1.0.3	(corresponding to Plugin)	8	8	0
InnoDB Plugin 1.0.3 and newer	(corresponding to Plugin)	0	No limit	0

Note that InnoDB will cause threads to sleep only when the number of concurrent threads is limited. When there is no limit on the number of threads, all will contend equally to be scheduled. That is, if `innodb_thread_concurrency` is 0, the value of `innodb_thread_sleep_delay` is ignored.

When there is a limit on the number of threads, InnoDB reduces context switching overhead by permitting multiple requests made during the execution of a single SQL statement to enter InnoDB without observing the limit set by `innodb_thread_concurrency`. Since an SQL statement (such as a join) may comprise multiple row operations within InnoDB, InnoDB assigns “tickets” that allow a thread to be scheduled repeatedly with minimal overhead.

When starting to execute a new SQL statement, a thread will have no tickets, and it must observe `innodb_thread_concurrency`. Once the thread is entitled to enter InnoDB, it will be assigned a number of tickets that it can use for subsequently entering InnoDB. If the tickets run out, `innodb_thread_concurrency` will be observed again and further tickets will be assigned. The number of tickets to assign is specified by the global option `innodb_concurrency_tickets`, which is 500 by default. A thread that is waiting for a lock will be given one ticket once the lock becomes available.

The correct values of these variables are dependent on your environment and workload. You will need to try a range of different values to determine what value works for your applications. Before limiting the number of concurrently executing threads, you should review configuration options that may improve the performance of InnoDB on multi-core and multi-processor computers, such as `innodb_use_sys_malloc` and `innodb_adaptive_hash_index`.

7.7 Changes in the Read Ahead Algorithm

A read ahead request is an I/O request to prefetch multiple pages in the buffer cache asynchronously in anticipation that these pages will be needed in the near future. InnoDB has historically used two read ahead algorithms to improve I/O performance.

Random read-ahead is a technique that predicts when pages might be needed soon based on pages already in the buffer pool, regardless of the order in which those pages were read. If 13 consecutive pages from the same extent are found in the buffer pool, InnoDB asynchronously issues a request to prefetch the remaining pages of the extent. This feature was initially removed from InnoDB starting with InnoDB Plugin 1.0.4 and turned off with MySQL 5.5. It is available once again starting in MySQL 5.1.59 and 5.5.16 and higher, turned off by default. To enable this feature, set the configuration variable `innodb_random_read_ahead`.

Linear read ahead is based on the access pattern of the pages in the buffer cache, not just their number. In releases before 1.0.4, if most pages belonging to some extent are accessed sequentially, InnoDB will issue an asynchronous prefetch request for the entire next extent when it reads in the last page of the current extent. Beginning with InnoDB Plugin 1.0.4, users can control when InnoDB performs a read

ahead, by adjusting the number of sequential page accesses required to trigger an asynchronous read request using the new configuration parameter `innodb_read_ahead_threshold`.

If the number of pages read from an extent of 64 pages is greater or equal to `innodb_read_ahead_threshold`, InnoDB will initiate an asynchronous read ahead of the entire following extent. Thus, this parameter controls how sensitive InnoDB is to the pattern of page accesses within an extent in deciding whether to read the following extent asynchronously. The higher the value, the more strict will be the access pattern check. For example, if you set the value to 48, InnoDB will trigger a linear read ahead request only when 48 pages in the current extent have been accessed sequentially. If the value is 8, InnoDB would trigger an asynchronous read ahead even if as few as 8 pages in the extent were accessed sequentially.

The new configuration parameter `innodb_read_ahead_threshold` may be set to any value from 0-64. The default value is 56, meaning that an asynchronous read ahead is performed only when 56 of the 64 pages in the extent are accessed sequentially. You can set the value of this parameter in the MySQL option file (`my.cnf` or `my.ini`), or change it dynamically with the `SET GLOBAL` command, which requires the `SUPER` privilege.

Starting with InnoDB Plugin 1.0.5 more statistics are provided through `SHOW ENGINE INNODB STATUS` command to measure the effectiveness of the read ahead algorithm. See [Section 8.9, “More Read Ahead Statistics”](#) for more information.

7.8 Multiple Background I/O Threads

InnoDB uses background threads to service various types of I/O requests. Starting from InnoDB Plugin 1.0.4, the number of background threads tasked with servicing read and write I/O on data pages is configurable. In previous versions of InnoDB, there was only one thread each for read and write on non-Windows platforms. On Windows, the number of background threads was controlled by `innodb_file_io_threads`. The configuration parameter `innodb_file_io_threads` has been removed in InnoDB Plugin 1.0.4. If you try to set a value for this parameter, a warning will be written to the log file and the value will be ignored.

In place of `innodb_file_io_threads`, two new configuration parameters are introduced in the InnoDB Plugin 1.0.4, which are effective on all supported platforms. The two parameters `innodb_read_io_threads` and `innodb_write_io_threads` signify the number of background threads used for read and write requests respectively. You can set the value of these parameters in the MySQL option file (`my.cnf` or `my.ini`). These parameters cannot be changed dynamically. The default value for these parameters is 4 and the permissible values range from 1-64.

The purpose of this change is to make InnoDB more scalable on high end systems. Each background thread can handle up to 256 pending I/O requests. A major source of background I/O is the read ahead requests. InnoDB tries to balance the load of incoming requests in such way that most of the background threads share work equally. InnoDB also attempts to allocate read requests from the same extent to the same thread to increase the chances of coalescing the requests together. If you have a high end I/O subsystem and you see more than 64 times `innodb_read_io_threads` pending read requests in `SHOW ENGINE INNODB STATUS`, then you may gain by increasing the value of `innodb_read_io_threads`.

7.9 Group Commit

InnoDB, like any other ACID compliant database engine, is required to flush the redo log of a transaction before it is committed. Historically InnoDB used group commit functionality to group multiple such flush requests together to avoid one flush for each commit. With group commit, InnoDB can issue a single write to the log file to effectuate the commit action for multiple user transactions that commit at about the same time, significantly improving throughput.

Group commit in InnoDB worked until MySQL 4.x. With the introduction of support for the distributed transactions and Two Phase Commit (2PC) in MySQL 5.0, group commit functionality inside InnoDB was broken.

Beginning with InnoDB Plugin 1.0.4, the group commit functionality inside InnoDB works with the Two Phase Commit protocol in MySQL. Re-enabling of the group commit functionality fully ensures that the ordering of commit in the MySQL binlog and the InnoDB logfile is the same as it was before. It means it is **totally safe to use InnoDB Hot Backup with InnoDB Plugin 1.0.4**.

Group commit is transparent to the user and nothing needs to be done by the user to take advantage of this significant performance improvement.

7.10 Controlling the Master Thread I/O Rate

The master thread in InnoDB performs various tasks in the background. Most of these tasks are I/O related like flushing of the dirty pages from the buffer cache or writing the buffered inserts to the appropriate secondary indexes. The master thread attempts to perform these tasks in a way that does not adversely affect the normal working of the server. It tries to estimate the free I/O bandwidth available and tune its activities to take advantage of this free capacity. Historically, InnoDB has used a hard coded value of 100 IOPs (input/output operations per second) as the total I/O capacity of the server.

Beginning with InnoDB Plugin 1.0.4, a new configuration parameter is introduced to indicate the overall I/O capacity available to InnoDB. The new parameter `innodb_io_capacity` should be set to approximately the number of I/O operations that the system can perform per second. The value will of course depend on your system configuration. When `innodb_io_capacity` is set, the master threads estimates the I/O bandwidth available for background tasks based on the set value. Setting the value to `100` reverts to the old behavior.

You can set the value of `innodb_io_capacity` to any number 100 or greater, and the default value is `200`. Typically, values around the previous default of 100 are appropriate for consumer-level storage devices, such as hard drives up to 7200 RPMs. Faster hard drives, RAID configurations, and SSDs benefit from higher values. You can set the value of this parameter in the MySQL option file (`my.cnf` or `my.ini`) or change it dynamically with the `SET GLOBAL` command, which requires the `SUPER` privilege.

7.11 Controlling the Flushing Rate of Dirty Pages

InnoDB performs certain tasks in the background, including flushing of dirty pages (those pages that have been changed but are not yet written to the database files) from the buffer cache, a task performed by the “master thread”. Currently, the master thread aggressively flushes buffer pool pages if the percentage of dirty pages in the buffer pool exceeds `innodb_max_dirty_pages_pct`.

This behavior can cause temporary reductions in throughput when excessive buffer pool flushing takes place, limiting the I/O capacity available for ordinary read and write activity. Beginning with release 1.0.4, InnoDB Plugin uses a new algorithm to estimate the required rate of flushing based on the speed of redo log generation and the current rate of flushing. The intent of this change is to smooth overall performance, eliminating steep dips in throughput, by ensuring that buffer flush activity keeps up with the need to keep the buffer pool “clean”.

Remember that InnoDB uses its log files in a circular fashion. To make a log file (or a portion of it) reusable, InnoDB must flush to disk all dirty buffer pool pages whose redo entries are contained in that portion of the log file. When required, InnoDB performs a so-called “sharp checkpoint” by flushing the appropriate dirty pages to make space available in the log file. If a workload is write intensive, it will generate a lot of redo information (writes to the log file). In this case, it is possible that available space in the log files will be used up, even though `innodb_max_dirty_pages_pct` is not reached. This will cause a sharp checkpoint, causing a temporary reduction in throughput.

Beginning with release 1.0.4, InnoDB Plugin uses a new heuristic-based algorithm to avoid such a scenario. The heuristic is a function of the number of dirty pages in the buffer cache and the rate at which redo is being generated. Based on this heuristic, the master thread will decide how many dirty pages to flush from the buffer cache each second. This self adapting heuristic is able to deal with sudden changes in the workload.

The primary aim of this feature is to smooth out I/O activity, avoiding sudden dips in throughput when flushing activity becomes high. Internal benchmarking has also shown that this algorithm not only maintains throughput over time, but can also improve overall throughput significantly.

Because adaptive flushing is a new feature that can significantly affect the I/O pattern of a workload, the InnoDB Plugin introduces a new configuration parameter that can be used to disable this feature. The default value of the new boolean parameter `innodb_adaptive_flushing` is `TRUE`, enabling the new algorithm. You can set the value of this parameter in the MySQL option file (`my.cnf` or `my.ini`) or change it dynamically with the `SET GLOBAL` command, which requires the `SUPER` privilege.

7.12 Using the PAUSE instruction in InnoDB spin loops

Synchronization inside InnoDB frequently involves the use of spin loops (where, while waiting, InnoDB executes a tight loop of instructions repeatedly to avoid having the InnoDB process and threads be rescheduled by the operating system). If the spin loops are executed too quickly, system resources are wasted, imposing a relatively severe penalty on transaction throughput. Most modern processors implement the PAUSE instruction for use in spin loops, so the processor can be more efficient.

Beginning with 1.0.4, the InnoDB Plugin uses a PAUSE instruction in its spin loops on all platforms where such an instruction is available. This technique increases overall performance with CPU-bound workloads, and has the added benefit of minimizing power consumption during the execution of the spin loops.

Using the PAUSE instruction in InnoDB spin loops is transparent to the user. User does not have to do anything to take advantage of this performance improvement.

7.13 Control of Spin Lock Polling

Many InnoDB mutexes and rw-locks are reserved for a short amount of time. On a multi-core system, it is often more efficient for a thread to actively poll a mutex or rw-lock for a while before sleeping. If the mutex or rw-lock becomes available during this polling period, the thread may continue immediately, in the same time slice. Alas, if a shared object is being polled too frequently by multiple threads, it may result in “cache ping-pong”, the shipping of cache lines between processors. InnoDB tries to avoid this by making threads busy, waiting a random time between subsequent polls. The delay is implemented as a busy loop.

Starting with InnoDB Plugin 1.0.4, it is possible to control the maximum delay between sampling a mutex or rw-lock using the new parameter `innodb_spin_wait_delay`. In the 100 MHz Pentium era, the unit of delay used to be one microsecond. The duration of the delay loop depends on the C compiler and the target processor. On a system where all processor cores share a fast cache memory, it might be useful to reduce the maximum delay or disable the busy loop altogether by setting `innodb_spin_wait_delay=0`. On a system that consists of multiple processor chips, the shipping of cache lines can be slower and it may be useful to increase the maximum delay.

The default value of `innodb_spin_wait_delay` is 6. The spin wait delay is a dynamic, global parameter that can be specified in the MySQL option file (`my.cnf` or `my.ini`) or changed at runtime with the command `SET GLOBAL innodb_spin_wait_delay=delay`, where `delay` is the desired maximum delay. Changing the setting requires the `SUPER` privilege.

7.14 Making Buffer Cache Scan Resistant

Historically, InnoDB has inserted newly read blocks into the middle of the list representing the buffer cache, to avoid pollution of the cache due to excessive read-ahead. The idea is that the read-ahead algorithm should not pollute the buffer cache by forcing the frequently accessed (“hot”) pages out of the LRU list. To achieve this, InnoDB internally maintains a pointer at 3/8 from the tail of the LRU list, and all newly read pages are inserted at this location in the LRU list. The pages are moved to the front of the list (the most-recently used end) when they are accessed from the buffer cache for the first time. Thus pages that are never accessed never make it to the front 5/8 of the LRU list.

The above arrangement logically divides the LRU list into two segments where the 3/8 pages downstream of the insertion point are considered “old” and are desirable victims for LRU eviction. Starting with InnoDB Plugin 1.0.5, this mechanism has been extended in two ways.

You can control the insertion point in the LRU list. A new configuration parameter `innodb_old_blocks_pct` now controls the percentage of “old” blocks in the LRU list. The default value of `innodb_old_blocks_pct` is 37, corresponding to the original fixed ratio of 3/8. The permissible value range is 5 to 95.

The optimization that keeps the buffer cache from being churned too much by read-ahead, is extended to avoid similar problems resulting from table or index scans. During an index scan, a data page is typically accessed a few times in quick succession and is then never touched again. InnoDB Plugin 1.0.5 introduces a new configuration parameter `innodb_old_blocks_time` which specifies the time window (in milliseconds) after the first access to a page during which it can be accessed without being moved to the front (most-recently used end) of the LRU list. The default value of `innodb_old_blocks_time` is 0, corresponding to the original behavior of moving a page to the MRU end of the LRU list on first access in the buffer pool.

Both the new parameters `innodb_old_blocks_pct` and `innodb_old_blocks_time` are dynamic, global and can be specified in the MySQL option file (`my.cnf` or `my.ini`) or changed at runtime with the `SET GLOBAL` command. Changing the setting requires the `SUPER` privilege.

To help you gauge the effect of setting these parameters, some additional statistics are reported by `SHOW ENGINE INNODB STATUS` command. The `BUFFER POOL AND MEMORY` section now looks like:

```
Total memory allocated 1107296256; in additional pool allocated 0
Dictionary memory allocated 80360
Buffer pool size 65535
Free buffers 0
Database pages 63920
Old database pages 23600
Modified db pages 34969
Pending reads 32
Pending writes: LRU 0, flush list 0, single page 0
Pages made young 414946, not young 2930673
1274.75 youngs/s, 16521.90 non-youngs/s
Pages read 486005, created 3178, written 160585
2132.37 reads/s, 3.40 creates/s, 323.74 writes/s
Buffer pool hit rate 950 / 1000, young-making rate 30 / 1000 not 392 / 1000
Pages read ahead 1510.10/s, evicted without access 0.00/s
LRU len: 63920, unzip_LRU len: 0
I/O sum[43690]:cur[221], unzip sum[0]:cur[0]
```

- `Old database pages` is the number of pages in the “old” segment of the LRU list.
- `Pages made young` and `not young` is the total number of “old” pages that have been made young or not respectively.
- `youngs/s` and `non-young/s` is the rate at which page accesses to the “old” pages have resulted in making such pages young or otherwise respectively since the last invocation of the command.

- `young-making rate` and `not` provides the same rate but in terms of overall buffer cache accesses instead of accesses just to the “old” pages.



Note

Per second averages provided in InnoDB Monitor output are based on the elapsed time between the current time and the last time InnoDB Monitor output was printed.

7.14.1 Guidelines for `innodb_old_blocks_pct` and `innodb_old_blocks_time`

The default values of both parameters leave the original behavior as of InnoDB Plugin 1.0.4 intact. To take advantage of this feature, you must set different values.

Because the effects of these parameters can vary widely based on your hardware configuration, your data, and the details of your workload, always benchmark to verify the effectiveness before changing these settings in any performance-critical or production environment.

In mixed workloads where most of the activity is OLTP type with periodic batch reporting queries which result in large scans, setting the value of `innodb_old_blocks_time` during the batch runs can help keep the working set of the normal workload in the buffer cache.

When scanning large tables that cannot fit entirely in the buffer pool, setting `innodb_old_blocks_pct` to a small value keeps the data that is only read once from consuming a significant portion of the buffer pool. For example, setting `innodb_old_blocks_pct=5` restricts this data that is only read once to 5% of the buffer pool.

When scanning small tables that do fit into memory, there is less overhead for moving pages around within the buffer pool, so you can leave `innodb_old_blocks_pct` at its default value, or even higher, such as `innodb_old_blocks_pct=50`.

The effect of the `innodb_old_blocks_time` parameter is harder to predict than the `innodb_old_blocks_pct` parameter, is relatively small, and varies more with the workload. To arrive at an optimal value, conduct your own benchmarks if the performance improvement from adjusting `innodb_old_blocks_pct` is not sufficient.

7.15 Improvements to Crash Recovery Performance

Starting with InnoDB Plugin 1.0.7, a number of optimizations speed up certain steps of the recovery that happens on the next startup after a crash. In particular, scanning the redo log and applying the redo log are faster. You do not need to take any actions to take advantage of this performance enhancement. If you kept the size of your logfiles artificially low because recovery took a long time, you can consider increasing the logfile size.

Chapter 8 Changes for Flexibility, Ease of Use and Reliability

Table of Contents

8.1 Overview	49
8.2 Enabling New File Formats	49
8.3 Dynamic Control of System Configuration Parameters	50
8.3.1 Dynamically Changing <code>innodb_file_per_table</code>	50
8.3.2 Dynamically Changing <code>innodb_stats_on_metadata</code>	50
8.3.3 Dynamically Changing <code>innodb_lock_wait_timeout</code>	51
8.3.4 Dynamically Changing <code>innodb_adaptive_hash_index</code>	51
8.4 <code>TRUNCATE TABLE</code> Reclaims Space	51
8.5 InnoDB Strict Mode	52
8.6 Controlling Optimizer Statistics Estimation	52
8.7 Better Error Handling when Dropping Indexes	53
8.8 More Compact Output of <code>SHOW ENGINE INNODB MUTEX</code>	54
8.9 More Read Ahead Statistics	54

8.1 Overview

This chapter describes several changes in the InnoDB Plugin that offer new flexibility and improve ease of use, reliability and performance:

- [Section 8.2, “Enabling New File Formats”](#)
- [Section 8.3, “Dynamic Control of System Configuration Parameters”](#)
- [Section 8.4, “`TRUNCATE TABLE` Reclaims Space”](#)
- [Section 8.5, “InnoDB Strict Mode”](#)
- [Section 8.6, “Controlling Optimizer Statistics Estimation”](#)
- [Section 8.7, “Better Error Handling when Dropping Indexes”](#)
- [Section 8.8, “More Compact Output of `SHOW ENGINE INNODB MUTEX`”](#)
- [Section 8.9, “More Read Ahead Statistics”](#)

8.2 Enabling New File Formats

The InnoDB Plugin introduces named file formats to improve compatibility between database file formats and various InnoDB versions.

To create new tables that require a new file format, you must enable the new “Barracuda” file format, using the configuration parameter `innodb_file_format`. The value of this parameter will determine whether it will be possible to create a table or index using compression or the new `DYNAMIC` row format. If you omit `innodb_file_format` or set it to “Antelope”, you preclude the use of new features that would make your database inaccessible to the built-in InnoDB in MySQL 5.1 and prior releases.

You can set the value of `innodb_file_format` on the command line when you start `mysqld`, or in the option file `my.cnf` (Unix operating systems) or `my.ini` (Windows). You can also change it dynamically with the `SET GLOBAL` command.

Further information about managing file formats is presented in [Chapter 4, InnoDB File-Format Management](#).

8.3 Dynamic Control of System Configuration Parameters

With the InnoDB Plugin it now is possible to change certain system configuration parameters dynamically, without shutting down and restarting the server as was previously necessary. This increases up-time and facilitates testing of various options. You can now set these parameters dynamically:

- `innodb_file_per_table`
- `innodb_stats_on_metadata`
- `innodb_lock_wait_timeout`
- `innodb_adaptive_hash_index`

8.3.1 Dynamically Changing `innodb_file_per_table`

Since MySQL version 4.1, InnoDB has provided two options for how tables are stored on disk. You can choose to create a new table and its indexes in the shared system tablespace (corresponding to the set of files named `ibdata` files), along with other internal InnoDB system information. Or, you can choose to use a separate file (an `.ibd` file) to store a new table and its indexes.

The tablespace style used for new tables is determined by the setting of the configuration parameter `innodb_file_per_table` at the time a table is created. Previously, the only way to set this parameter was in the MySQL option file (`my.cnf` or `my.ini`), and changing it required shutting down and restarting the server. Beginning with the InnoDB Plugin, the configuration parameter `innodb_file_per_table` is dynamic, and can be set `ON` or `OFF` using the `SET GLOBAL` command. The default setting is `OFF`, so new tables and indexes are created in the system tablespace. Dynamically changing the value of this parameter requires the `SUPER` privilege and immediately affects the operation of all connections.

Tables created when `innodb_file_per_table` is disabled cannot use the new compression capability, or use the new row format `DYNAMIC`. Tables created when `innodb_file_per_table` is enabled can use those new features, and each table and its indexes will be stored in a new `.ibd` file.

The ability to change the setting of `innodb_file_per_table` dynamically is useful for testing. As noted above, the parameter `innodb_file_format` is also dynamic, and must be set to “Barracuda” to create new compressed tables, or tables that use the new row format `DYNAMIC`. Since both parameters are dynamic, it is easy to experiment with these table formats and the downgrade procedure described in [Chapter 11, Downgrading from the InnoDB Plugin](#) without a system shutdown and restart.

Note that the InnoDB Plugin can add and drop a table’s secondary indexes without re-creating the table, but must recreate the table when you change the clustered (primary key) index (see [Chapter 2, Fast Index Creation in the InnoDB Storage Engine](#)). When a table is recreated as a result of creating or dropping an index, the table and its indexes will be stored in the shared system tablespace or in its own `.ibd` file just as if it were created using a `CREATE TABLE` command (and depending on the setting of `innodb_file_per_table`). When an index is created without rebuilding the table, the index is stored in the same file as the clustered index, regardless of the setting of `innodb_file_per_table`.

8.3.2 Dynamically Changing `innodb_stats_on_metadata`

As noted in [Section 8.6, “Controlling Optimizer Statistics Estimation”](#), the InnoDB Plugin allows you to control the way in which InnoDB gathers information about the number of distinct values in an index key. A related parameter, `innodb_stats_on_metadata`, has existed since MySQL release 5.1.17 to control

whether or not InnoDB performs statistics gathering when metadata statements are executed. See the MySQL manual on [InnoDB Startup Options and System Variables](#) for details.

Beginning with release 1.0.2 of the InnoDB Plugin, it is possible to change the setting of `innodb_stats_on_metadata` dynamically at runtime with the command `SET GLOBAL innodb_stats_on_metadata=mode`, where *mode* is either `ON` or `OFF` (or `1` or `0`). Changing this setting requires the `SUPER` privilege and immediately affects the operation of all connections.

8.3.3 Dynamically Changing `innodb_lock_wait_timeout`

When a transaction is waiting for a resource, it will wait for the resource to become free, or stop waiting and return with the error

```
ERROR HY000: Lock wait timeout exceeded; try restarting transaction
```

The length of time a transaction will wait for a resource before “giving up” is determined by the value of the configuration parameter `innodb_lock_wait_timeout`. The default setting for this parameter is 50 seconds. The minimum setting is 1 second, and values above 100,000,000 disable the timeout, so a transaction will wait “forever”. Following a timeout, the SQL statement that was executing will be rolled back. (In MySQL 5.0.12 and earlier, the transaction rolled back.) The user application may try the statement again (usually after waiting for a while), or rollback the entire transaction and restart.

Before InnoDB Plugin 1.0.2, the only way to set this parameter was in the MySQL option file (`my.cnf` or `my.ini`), and changing it required shutting down and restarting the server. Beginning with the InnoDB Plugin 1.0.2, the configuration parameter `innodb_lock_wait_timeout` can be set at runtime with the `SET GLOBAL` or `SET SESSION` commands. Changing the `GLOBAL` setting requires the `SUPER` privilege and affects the operation of all clients that subsequently connect. Any client can change the `SESSION` setting for `innodb_lock_wait_timeout`, which affects only that client.

8.3.4 Dynamically Changing `innodb_adaptive_hash_index`

As described in [Section 7.5, “Controlling Adaptive Hash Indexing”](#), it may be desirable, depending on your workload, to dynamically enable or disable the adaptive hash indexing scheme InnoDB uses to improve query performance.

Version 5.1.24 of MySQL introduced the start-up option `innodb_adaptive_hash_index` that allows the adaptive hash index to be disabled. It is enabled by default. Starting with InnoDB Plugin 1.0.3, the parameter can be modified by the `SET GLOBAL` command, without restarting the server. Changing the setting requires the `SUPER` privilege.

Disabling the adaptive hash index will empty the hash table immediately. Normal operations can continue while the hash table is emptied, and executing queries that have been using the hash table will access the index B-trees directly instead of attempting to utilize the hash index. When the adaptive hash index is enabled, the hash table will be populated during normal operation.

8.4 `TRUNCATE TABLE` Reclaims Space

Starting with the InnoDB Plugin, when the user requests to `TRUNCATE` a table that is stored in an `.ibd` file of its own (because `innodb_file_per_table` was enabled when the table was created), and if the table is not referenced in a `FOREIGN KEY` constraint, the InnoDB Plugin will drop and re-create the table in a new `.ibd` file. This operation is much faster than deleting the rows one by one, and will return disk space to the operating system and reduce the size of page-level backups.

Previous versions of InnoDB would re-use the existing `.ibd` file, thus releasing the space only to InnoDB for storage management, but not to the operating system. Note that when the table is truncated, the count of rows affected by the `TRUNCATE` command is an arbitrary number.

Note: if there are referential constraints between the table being truncated and other tables, MySQL instead automatically converts the `TRUNCATE` command to a `DELETE` command that operates row-by-row, so that `ON DELETE` operations can occur on “child” tables.

8.5 InnoDB Strict Mode

To guard against ignored typos and syntax errors in SQL, or other unintended consequences of various combinations of operational modes and SQL commands, the InnoDB Plugin provides a “strict mode” of operations. In this mode, InnoDB will raise error conditions in certain cases, rather than issue a warning and process the specified command (perhaps with some unintended defaults). This is analogous to MySQL's `sql_mode`, which controls what SQL syntax MySQL will accept, and determines whether it will silently ignore errors, or validate input syntax and data values. Note that there is no strict mode with the built-in InnoDB, so some commands that execute without errors with the built-in InnoDB will generate errors with the InnoDB Plugin, unless you disable strict mode.

In the InnoDB Plugin, the setting of InnoDB strict mode affects the handling of syntax errors on the `CREATE TABLE`, `ALTER TABLE` and `CREATE INDEX` commands. Starting with InnoDB Plugin version 1.0.2, the strict mode also enables a record size check, so that an `INSERT` or `UPDATE` will never fail due to the record being too large for the selected page size.

Using the new clauses and settings for `ROW_FORMAT` and `KEY_BLOCK_SIZE` on `CREATE TABLE` and `ALTER TABLE` commands and the `CREATE INDEX` can be confusing when not running in strict mode. Unless you run in strict mode, InnoDB will ignore certain syntax errors and will create the table or index, with only a warning in the message log. However if InnoDB strict mode is on, such errors will generate an immediate error and the table or index will not be created, thus saving time by catching the error at the time the command is issued.

The default for strict mode is off, but in the future, the default may be changed. It is best to start using strict mode with the InnoDB Plugin, and make sure your SQL scripts use commands that do not generate warnings or unintended effects.

InnoDB strict mode is set with the configuration parameter `innodb_strict_mode`, which can be specified as `on` or `off`. You can set the value on the command line when you start `mysqld`, or in the configuration file `my.cnf` (Unix operating systems) or `my.ini` (Windows). You can also enable or disable InnoDB strict mode at runtime with the command `SET [GLOBAL|SESSION] innodb_strict_mode=mode`, where `mode` is either `ON` or `OFF`. Changing the `GLOBAL` setting requires the `SUPER` privilege and affects the operation of all clients that subsequently connect. Any client can change the `SESSION` setting for `innodb_strict_mode`, which affects only that client.

8.6 Controlling Optimizer Statistics Estimation

The MySQL query optimizer uses estimated statistics about key distributions to select or avoid using an index in an execution plan, based on the relative selectivity of the index. Previously, InnoDB sampled 8 random pages from an index to get an estimate of the cardinality of (i.e., the number of distinct values in) the index. (This page sampling technique is frequently described as “index dives”.) This small number of page samples frequently was insufficient, and could give inaccurate estimates of an index's selectivity and thus lead to poor choices by the query optimizer.

To give users control over the quality of the statistics estimate (and thus better information for the query optimizer), the number of sampled pages now can be changed using the parameter `innodb_stats_sample_pages`.

This feature addresses user requests such as that as expressed in [MySQL Bug #25640: InnoDB Analyze Table Should Allow User Selection of Index Dives](#).

You can change the number of sampled pages using the global parameter `innodb_stats_sample_pages`, which can be set at runtime (i.e., it is a dynamic parameter). The default value for this parameter is 8, preserving the same behavior as in past releases.

Note that the value of `innodb_stats_sample_pages` affects the index sampling for *all* tables and indexes. You should also be aware that there are the following potentially significant impacts when you change the index sample size:

- small values like 1 or 2 can result in very inaccurate estimates of cardinality
- values much larger than 8 (say, 100), can cause a big slowdown in the time it takes to open a table or execute `SHOW TABLE STATUS`.
- the optimizer may choose very different query plans based on different estimates of index selectivity

Note that the cardinality estimation can be disabled for metadata commands such as `SHOW TABLE STATUS` by executing the command `SET GLOBAL innodb_stats_on_metadata=OFF` (or 0). Before InnoDB Plugin 1.0.2, this variable could only be set in the MySQL option file (`my.cnf` or `my.ini`), and changing it required shutting down and restarting the server.

The cardinality (the number of different key values) in every index of a table is calculated when a table is opened, at `SHOW TABLE STATUS` and `ANALYZE TABLE` and on other circumstances (like when the table has changed too much). Note that all tables are opened, and the statistics are re-estimated, when the `mysql` client starts if the auto-rehash setting is set on (the default). The auto-rehash feature enables automatic name completion of database, table, and column names for interactive users. You may prefer setting auto-rehash off to improve the start up time of the `mysql` client.

You should note that it does not make sense to increase the index sample size, then run `ANALYZE TABLE` and decrease sample size to attempt to obtain better statistics. This is because the statistics are not persistent. They are automatically recalculated at various times other than on execution of `ANALYZE TABLE`. Sooner or later the “better” statistics calculated by `ANALYZE` running with a high value of `innodb_stats_sample_pages` will be wiped away.

The estimated cardinality for an index will be more accurate with a larger number of samples, but each sample might require a disk read, so you do not want to make the sample size too large. You should choose a value for `innodb_stats_sample_pages` that results in reasonably accurate estimates for all tables in your database without requiring excessive I/O.

Although it is not possible to specify the sample size on a per-table basis, smaller tables generally would require fewer index samples than larger tables require. If your database has many large tables, you may want to consider using a higher value for `innodb_stats_sample_pages` than if you have mostly smaller tables.

8.7 Better Error Handling when Dropping Indexes

For efficiency, InnoDB requires an index to exist on foreign key columns so that `UPDATE` and `DELETE` operations on a “parent” table can easily check for the existence or non-existence of corresponding rows in the “child” table. To ensure that there is an appropriate index for such checks, MySQL will sometimes implicitly create or drop such indexes as a side-effect of `CREATE TABLE`, `CREATE INDEX`, and `ALTER TABLE` statements.

When you explicitly `DROP` an index, InnoDB will check that an index suitable for referential integrity checking will still exist following the `DROP` of the index. InnoDB will prevent you from dropping the last usable index for enforcing any given referential constraint. Users have been confused by this behavior, as reported in MySQL Bug #21395.

In releases prior to InnoDB Plugin 1.0.2, attempts to drop the only usable index would result in an error message such as

```
ERROR 1025 (HY000): Error on rename of './db2/#sql-18eb_3'  
to './db2/foo'(errno: 150)
```

Beginning with InnoDB Plugin 1.0.2, this error condition is reported with a more friendly message:

```
ERROR 1553 (HY000): Cannot drop index 'fooIdx':  
needed in a foreign key constraint
```

As a related matter, because all user data in InnoDB is maintained in the so-called “clustered index” (or primary key index), InnoDB ensures that there is such an index for every table, even if the user does not declare an explicit `PRIMARY KEY`. In such cases, InnoDB will create an implicit clustered index using the first columns of the table that have been declared `UNIQUE` and `NOT NULL`.

When the InnoDB Plugin is used with a MySQL version earlier than 5.1.29, an attempt to drop an implicit clustered index (the first `UNIQUE NOT NULL` index) will fail if the table does not contain a `PRIMARY KEY`. This has been reported as MySQL Bug #31233. Attempts to use the `DROP INDEX` or `ALTER TABLE` command to drop such an index will generate this error:

```
ERROR 42000: This table type requires a primary key
```

Beginning with MySQL 5.1.29 when using the InnoDB Plugin, attempts to drop such an index will copy the table, rebuilding the index using a different `UNIQUE NOT NULL` group of columns or a system-generated key. Note that all indexes will be re-created by copying the table, as described in [Section 2.3, “Implementation”](#).

In those versions of MySQL that are affected by this bug, one way to change an index of this type is to create a new table and copy the data into it using `INSERT INTO newtable SELECT * FROM oldtable`, and then `DROP` the old table and rename the new table.

However, if there are existing tables with references to the table whose index you are dropping, you will first need to use the `ALTER TABLE` command to remove foreign key references from or to other tables. Unfortunately, MySQL does not support dropping or creating `FOREIGN KEY` constraints, even though dropping a constraint would be trivial. Therefore, if you use `ALTER TABLE` to add or remove a `REFERENCES` constraint, the child table will be copied, rather than using “Fast Index Creation”.

8.8 More Compact Output of `SHOW ENGINE INNODB MUTEX`

The command `SHOW ENGINE INNODB MUTEX` displays information about InnoDB mutexes and rw-locks. It can be a useful tuning aid on multi-core systems. However, with a big buffer pool, the size of the output may be overwhelming. There is a mutex and rw-lock in each 16K buffer pool block. It is highly improbable that an individual block mutex or rw-lock could become a performance bottleneck, and there are 65,536 blocks per gigabyte.

Starting with InnoDB Plugin 1.0.4, `SHOW ENGINE INNODB MUTEX` will skip the mutexes and rw-locks of buffer pool blocks. Furthermore, it will not list any mutexes or rw-locks that have never been waited on (`os_waits=0`). Therefore, `SHOW ENGINE INNODB MUTEX` only displays information about such mutexes and rw-locks that does not belong to the buffer pool blocks and for whom there have been at least one OS level wait.

8.9 More Read Ahead Statistics

As described in [Section 7.7, “Changes in the Read Ahead Algorithm”](#) a read ahead request is an asynchronous IO request issued in anticipation that the page being read in will be used in near future.

It can be very useful if a DBA has the information about how many pages are read in as part of read ahead and how many of them are evicted from the buffer pool without ever being accessed. Based on this information a DBA can then fine tune the degree of aggressiveness of the read ahead using the parameter `innodb_read_ahead_threshold`.

Starting from InnoDB Plugin 1.0.5 two new status variables are added to the `SHOW STATUS` output. These global status variables `InnoDB_buffer_pool_read_ahead` and `InnoDB_buffer_pool_read_ahead_evicted` indicate the number of pages read in as part of read ahead and the number of such pages evicted without ever being accessed respectively. These counters provide global values since the start of the server. Please also note that the status variables `InnoDB_buffer_pool_read_ahead_rnd` and `InnoDB_buffer_pool_read_ahead_seq` have been removed from the `SHOW STATUS` output.

In addition to the two counters mentioned above `SHOW INNODB STATUS` will also show the rate at which the read ahead pages are being read in and the rate at which such pages are being evicted without being accessed. The per second averages are based on the statistics collected since the last invocation of `SHOW INNODB STATUS` and are displayed in the `BUFFER POOL AND MEMORY` section of the output.

Chapter 9 Installing the InnoDB Plugin

Table of Contents

9.1 Overview of Installing the InnoDB Plugin	57
9.2 Checking for Compatible Version Levels	58
9.3 Installing the Precompiled InnoDB Plugin as a Shared Library	58
9.3.1 Installing the InnoDB Plugin as a Shared Library on Unix or Linux	59
9.3.2 Installing the Binary InnoDB Plugin as a Shared Library on Microsoft Windows	62
9.3.3 Errors When Installing the InnoDB Plugin on Microsoft Windows	64
9.4 Building the InnoDB Plugin from Source Code	65
9.4.1 Building the InnoDB Plugin on Linux or Unix	66
9.4.2 Building the InnoDB Plugin on Microsoft Windows	67
9.5 Configuring the InnoDB Plugin	68
9.6 Frequently Asked Questions about Plugin Installation	69
9.6.1 Should I use the InnoDB-supplied plugin or the one that is included with MySQL 5.1.38 or higher?	69
9.6.2 Why doesn't the MySQL service on Windows start after the replacement?	69
9.6.3 The Plugin is installed... now what?	69
9.6.4 Once the Plugin is installed, is it permanent?	69

9.1 Overview of Installing the InnoDB Plugin

You can acquire the plugin in these formats:

- As a platform-specific executable binary file that is dynamically linked or “plugged in” to the MySQL server.
- In source code form, available under the GNU General Public License (GPL), version 2.



Notes:

While it is possible to use the source code to build a special version of MySQL containing the InnoDB Plugin, we recommend you install the binary shared library for the InnoDB Plugin instead, without building from source. Replacing the shared library is simpler and much less error prone than building from source.

The InnoDB Plugin is included in the MySQL distribution, starting from MySQL 5.1.38. From MySQL 5.1.46 and up, this is the only download location for the InnoDB Plugin; it is not available from the InnoDB web site. If you used any scripts or configuration files with the earlier InnoDB Plugin from the InnoDB web site, be aware that the filename of the shared library as supplied by MySQL is `ha_innodb_plugin.so` or `ha_innodb_plugin.dll`, as opposed to `ha_innodb.so` or `ha_innodb.dll` in the older Plugin downloaded from the InnoDB web site. You might need to change the applicable file names in your startup or configuration scripts.

This discussion pertains to using the InnoDB Plugin with the MySQL Community Edition, whether source or binary. Except for download locations for MySQL software, the procedures documented here should work without change when you use MySQL Enterprise.

Whether you dynamically install the binary InnoDB Plugin or build from source, configure MySQL by editing the configuration file to use InnoDB as the default

engine (if desired) and set appropriate configuration parameters to enable use of new InnoDB Plugin features, as described in [Section 9.5, “Configuring the InnoDB Plugin”](#).

At this time, the InnoDB Plugin has not been compiled or tested with the Intel C Compiler (`icc`), so you should use a version of MySQL compiled with the GNU Compiler Collection (`gcc`).

9.2 Checking for Compatible Version Levels

Use the following table to confirm that the version of the InnoDB Plugin (whether source or binary) is compatible with your platform, hardware type (including 32-bit vs 64-bit) and with your version of MySQL. In general, a specific release of the InnoDB Plugin is designed to be compatible with a range of MySQL versions, but this may vary, so check the information on the download page.

When building MySQL from source, you can generally use the source for the InnoDB Plugin in place of the source for the built-in InnoDB. However, due to limitations of MySQL, a given binary version of the InnoDB Plugin is compatible only with a specific version of MySQL, as follows.

Table 9.1 InnoDB Plugin Compatibility

InnoDB Plugin Release	MySQL Release (Binary Compatibility)	MySQL Release (Source Compatibility)
1.0.0	5.1.23	5.1.23 or newer
1.0.1	5.1.24	5.1.24 or newer
1.0.2	5.1.30	5.1.24 or newer
1.0.3	5.1.30 (not 5.1.31)	5.1.24 or newer
1.0.4	5.1.37	5.1.24 or newer
1.0.5	5.1.41	5.1.24 or newer
1.0.6	5.1.41	5.1.24 or newer
1.0.7 and higher	Incorporated into the applicable 5.1.x MySQL release; not separately downloadable.	N/A



Note

[MySQL Bug #42610](#) prevents using the binary InnoDB Plugin with MySQL 5.1.31 or 5.1.32. There is no binary InnoDB Plugin for MySQL 5.1.33. The only way to use InnoDB Plugin with MySQL 5.1.31 through 5.1.35 is by building from source. This issue was resolved in MySQL 5.1.37 and InnoDB Plugin 1.0.4.

9.3 Installing the Precompiled InnoDB Plugin as a Shared Library

The simplest way to install the InnoDB Plugin is to use a precompiled (binary) shared library file, when one is available. The procedures are similar for installing the InnoDB Plugin using the binary on different hardware and operating systems platforms, but the specific details differ between Unix or Linux and Microsoft Windows. See below for notes specific to your platform.

Note that due to [MySQL Bug #42610](#), the procedure of installing the binary InnoDB Plugin changed in MySQL 5.1.33. If your version of MySQL is older than 5.1.33, refer to [Appendix B, Using the InnoDB Plugin with MySQL 5.1.30 or Earlier](#).

The steps for installing the InnoDB Plugin as a shared library are as follows:

- Download, extract and install the suitable MySQL executable for your platform.
- Make sure the MySQL server is not running. If you have to shut down the database server, you use a special “slow” shutdown procedure, described later.
- On database startup, make MySQL ignore the builtin InnoDB, and load the InnoDB Plugin and all new InnoDB Information Schema tables implemented in the InnoDB Plugin, using one of these alternatives:
 - Edit the option file (my.cnf, or my.ini) to contain the necessary options.
 - Specify equivalent options on the MySQL command line.
 - Edit the option file to disable InnoDB, then use `INSTALL` statements on the MySQL command line after startup.

These procedures are described in detail in the following sections.

- Set appropriate configuration parameters to enable new InnoDB Plugin features.
- Start MySQL, and verify the installation of the plugins.

The following sections detail these steps for Unix or Linux systems, and for Microsoft Windows.

9.3.1 Installing the InnoDB Plugin as a Shared Library on Unix or Linux

For Unix and Linux systems, use the following procedure to install the InnoDB Plugin as a shared library:

1. Download, extract and install the suitable MySQL executable for your server platform and operating system from the [MySQL download section for MySQL Database Server 5.1](#). Be sure to use a 32-bit or 64-bit version as appropriate for your hardware and operating system.
2. Make sure the MySQL server is not running. If the server is running, do a “slow” shutdown by issuing the following command before performing the shutdown:

```
SET GLOBAL innodb_fast_shutdown=0;
```

Then finish the shutdown process, as described in [The Shutdown Process](#) in the MySQL documentation. This option setting performs a full purge and an insert buffer merge before the shutdown, which can typically take minutes, or even hours for very large and busy databases.

3. The InnoDB Plugin shared library is already installed in the directory `lib/plugin` as part of the MySQL installation.
4. Edit the option file (my.cnf) to ignore the builtin InnoDB, and load the InnoDB Plugin and all Information Schema tables implemented in the InnoDB Plugin when the server starts:

```
ignore_builtin_innodb
plugin-load=innodb=ha_innodb_plugin.so;innodb_trx=ha_innodb_plugin.so;
innodb_locks=ha_innodb_plugin.so;innodb_lock_waits=ha_innodb_plugin.so;
innodb_cmp=ha_innodb_plugin.so;innodb_cmp_reset=ha_innodb_plugin.so;
innodb_cmpmem=ha_innodb_plugin.so;innodb_cmpmem_reset=ha_innodb_plugin.so
```

Note that all plugins for `plugin-load` should be on the same line in the option file.

Alternatively, you can use the equivalent options on the MySQL command line:

```
mysqld --ignore-builtin-innodb --plugin-load=innodb=ha_innodb_plugin.so;
```

```
innodb_trx=ha_innodb_plugin.so;innodb_locks=ha_innodb_plugin.so;
innodb_lock_waits=ha_innodb_plugin.so;innodb_cmp=ha_innodb_plugin.so;
innodb_cmp_reset=ha_innodb_plugin.so;innodb_cmpmem=ha_innodb_plugin.so;
innodb_cmpmem_reset=ha_innodb_plugin.so
```

You can also install the InnoDB Plugin and the new InnoDB Information Schema tables implemented in `ha_innodb_plugin.so` with `INSTALL` commands:

```
INSTALL PLUGIN INNODB SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_TRX SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_LOCKS SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_LOCK_WAITS SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_CMP SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_CMP_RESET SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_CMPMEM SONAME 'ha_innodb_plugin.so';
INSTALL PLUGIN INNODB_CMPMEM_RESET SONAME 'ha_innodb_plugin.so';
```

If you use `INSTALL PLUGIN` statement to install the InnoDB Plugin and the Information Schema tables, ensure the following conditions are set up:

- In the `mysqld` command line or `my.cnf` option file, prepend each InnoDB option with `loose_`, so that MySQL will start even when InnoDB is unavailable. For example, write `loose_innodb_file_per_table` instead of `innodb_file_per_table`.
- Start the MySQL server while it is configured to skip loading the built-in InnoDB and to make MyISAM the default storage engine. This can be done by editing the option file `my.cnf` to contain these two lines:

```
ignore_builtin_innodb
default_storage_engine=MyISAM
```

Or, you can use the equivalent options on the MySQL command line:

```
mysqld --ignore-builtin-innodb --default-storage-engine=MyISAM ...
```

See the MySQL Manual section on [INSTALL PLUGIN Syntax](#) for information on how these commands work.

5. Edit the option file `my.cnf` to use InnoDB as the default engine (if desired) and set appropriate configuration parameters to enable use of new InnoDB Plugin features, as described in [Section 9.5, “Configuring the InnoDB Plugin”](#). In particular, we recommend that you set the following specific parameters as follows:

```
default-storage-engine=InnoDB
innodb_file_per_table=1
innodb_file_format=barracuda
innodb_strict_mode=1
```



IMPORTANT:

The MySQL server always must be started with the option `ignore_builtin_innodb`, as long as you want to use the InnoDB Plugin as a shared library. Also, remember that the startup option `skip_grant_tables` prevents MySQL from loading any plugins.

6. Verify the installation of the plugins with the MySQL command `SHOW PLUGINS`, which should produce the following output:

Name	Status	Type	Library	License
binlog	ACTIVE	STORAGE ENGINE	NULL	GPL
CSV	ACTIVE	STORAGE ENGINE	NULL	GPL
MEMORY	ACTIVE	STORAGE ENGINE	NULL	GPL
InnoDB	ACTIVE	STORAGE ENGINE	ha_innodb_plugin.so	GPL
INNODB_TRX	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
INNODB_LOCKS	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
INNODB_LOCK_WAITS	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
INNODB_CMP	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
INNODB_CMP_RESET	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
INNODB_CMPMEM	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
INNODB_CMPMEM_RESET	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.so	GPL
MRG_MYISAM	ACTIVE	STORAGE ENGINE	NULL	GPL
MyISAM	ACTIVE	STORAGE ENGINE	NULL	GPL

If the plugins fail to load properly, see [Section 9.3.1.1, “Errors When Installing the InnoDB Plugin on Unix or Linux”](#) for possible causes and corrections.

After verifying that the Plugin is recognized by MySQL, create an InnoDB table as another confirmation of success.

9.3.1.1 Errors When Installing the InnoDB Plugin on Unix or Linux

If MySQL or its associated daemon process cannot start, or a post-startup `INSTALL PLUGIN` statement fails, look at the MySQL error log (usually named `machine_name.err` and located in the MySQL `data` directory) for the detailed error message. The log is in chronological order, so look at the end of the file. Try to resolve the problem based on other information in the message.

The following table outlines installation-related InnoDB Plugin error conditions or messages and possible solutions.

Error Condition or Message	Possible Solution
<code>Can't open shared library library_name</code>	Diagnose the cause from the following message details.
<code>API version for STORAGE ENGINE plugin is too different</code>	The version of the Plugin is not compatible with the version of the MySQL server. Consult the compatibility chart .
<code>No such file or directory</code>	Check that the file <code>ha_innodb_plugin.so</code> or <code>.dll</code> was copied to the correct location. Confirm that you specified the right file name (<code>ha_innodb_plugin.so</code> or <code>.dll</code> for the library

Error Condition or Message	Possible Solution
	from the InnoDB web site; <code>ha_innodb_plugin.so</code> or <code>.dll</code> for the library supplied along with MySQL 5.1.38 and up.)
<code>Permission denied</code>	Check that the directory and file access permissions are set properly, or change them using <code>chmod</code> on Unix-like systems . The <code>mysqld</code> process must have permission to read (r) the file <code>ha_innodb_plugin.so</code> and to access files (x) in the plugin directory.
<code>wrong ELF class</code> or any other message	Ensure that <code>ha_innodb_plugin.so</code> is for the same system platform as <code>mysqld</code> . In particular, note that a 32-bit <code>mysqld</code> is unable to load a 64-bit plugin, and vice versa. Be sure to download an InnoDB Plugin that is compatible with your platform.

**Note**

The Information Schema tables are themselves plugins to the MySQL server, but they depend on having the InnoDB storage engine plugin installed as well. These tables will appear to be empty if the storage engine is not installed.

9.3.2 Installing the Binary InnoDB Plugin as a Shared Library on Microsoft Windows

The InnoDB Plugin is supported on any of the Windows operating system versions supported by MySQL. In particular, this includes Microsoft Windows 2008 Server, Windows Vista, Windows 2003 Server and Windows XP. Note that on Vista certain special procedures must be followed that are not documented here.

Use the following procedure to dynamically install the InnoDB Plugin on Microsoft Windows.

1. Download, extract and install the suitable MySQL executable for your server platform and operating system from the [MySQL download section for MySQL Database Server 5.1](#). Be sure to use a 32-bit or 64-bit version as appropriate for your hardware and Windows version.
2. Make sure the MySQL server is not running. You do a “slow” shutdown by issuing the following command before performing the shutdown:

```
SET GLOBAL innodb_fast_shutdown=0;
```

Then finish the shutdown process, as described in [The Shutdown Process](#) in the MySQL documentation. This option setting performs a full purge and an insert buffer merge before the shutdown, which can typically take minutes, or even hours for very large and busy databases.

3. The InnoDB Plugin shared library is already installed in the directory `lib\plugin` as part of the MySQL installation.
4. Edit the option file (`my.ini`) to ignore the builtin InnoDB, and load the InnoDB Plugin and all Information Schema tables implemented in the InnoDB Plugin when the server starts:

```
ignore_builtin_innodb
plugin-load=innodb=ha_innodb_plugin.dll;innodb_trx=ha_innodb_plugin.dll;
innodb_locks=ha_innodb_plugin.dll;innodb_lock_waits=ha_innodb_plugin.dll;
innodb_cmp=ha_innodb_plugin.dll;innodb_cmp_reset=ha_innodb_plugin.dll;
```

```
innodb_cmpmem=ha_innodb_plugin.dll;innodb_cmpmem_reset=ha_innodb_plugin.dll
```

**Note**

All plugins for `plugin-load` should be on the same line in the option file. Be careful when copying and pasting that the line does not split.

Alternatively, you can use the equivalent options on the MySQL command line:

```
mysqld --ignore-builtin-innodb --plugin-load=
innodb=ha_innodb_plugin.dll;
innodb_trx=ha_innodb_plugin.dll;innodb_locks=ha_innodb_plugin.dll;
innodb_lock_waits=ha_innodb_plugin.dll;innodb_cmp=ha_innodb_plugin.dll;
innodb_cmp_reset=ha_innodb_plugin.dll;innodb_cmpmem=ha_innodb_plugin.dll;
innodb_cmpmem_reset=ha_innodb_plugin.dll
```

You can also install the InnoDB Plugin and the new InnoDB Information Schema tables implemented in `ha_innodb_plugin.so` with `INSTALL` commands, as follows:

```
INSTALL PLUGIN INNODB SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_TRX SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_LOCKS SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_LOCK_WAITS SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_CMP SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_CMP_RESET SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_CMPMEM SONAME 'ha_innodb_plugin.dll';
INSTALL PLUGIN INNODB_CMPMEM_RESET SONAME 'ha_innodb_plugin.dll';
```

If you use `INSTALL PLUGIN` statements to install the InnoDB Plugin and the Information Schema tables, ensure the following conditions are set up:

- In the `mysqld` command line or `my.ini` option file, prepend each InnoDB option with `loose_`, so that MySQL will start even when InnoDB is unavailable. For example, write `loose_innodb_file_per_table` instead of `innodb_file_per_table`.
- Start the MySQL server while it is configured to skip loading the built-in InnoDB and to make MyISAM the default storage engine. This can be done by editing the option file `my.cnf` to contain these two lines:

```
ignore_builtin_innodb
default_storage_engine=MyISAM
```

Or, you can use the equivalent options on the MySQL command line:

```
mysqld --ignore-builtin-innodb --default-storage-engine=MyISAM ...
```

See the MySQL Manual section on `INSTALL PLUGIN Syntax` for information on how these commands work.

5. Edit the option file `my.ini` to use InnoDB as the default engine (if desired) and set appropriate configuration parameters to enable use of new InnoDB Plugin features, as described in [Section 9.5, “Configuring the InnoDB Plugin”](#). In particular, we recommend that you set the following specific parameters as follows:

```
default-storage-engine=InnoDB
innodb_file_per_table=1
```

```
innodb_file_format=barracuda
innodb_strict_mode=1
```

IMPORTANT: The MySQL server always must be started with the option `ignore_builtin_innodb`, as long as you want to use the dynamic InnoDB Plugin. Also, remember that the startup option `skip_grant_tables` prevents MySQL from loading any plugins.

- Verify the installation of the plugins with the MySQL command `SHOW PLUGINS`, which should produce the following output:

Name	Status	Type	Library	License
binlog	ACTIVE	STORAGE ENGINE	NULL	GPL
CSV	ACTIVE	STORAGE ENGINE	NULL	GPL
MEMORY	ACTIVE	STORAGE ENGINE	NULL	GPL
InnoDB	ACTIVE	STORAGE ENGINE	ha_innodb_plugin.dll	GPL
INNODB_TRX	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
INNODB_LOCKS	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
INNODB_LOCK_WAITS	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
INNODB_CMP	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
INNODB_CMP_RESET	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
INNODB_CMPMEM	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
INNODB_CMPMEM_RESET	ACTIVE	INFORMATION SCHEMA	ha_innodb_plugin.dll	GPL
MRG_MYISAM	ACTIVE	STORAGE ENGINE	NULL	GPL
MyISAM	ACTIVE	STORAGE ENGINE	NULL	GPL

If the plugins fail to load properly, see [Section 9.3.3, “Errors When Installing the InnoDB Plugin on Microsoft Windows”](#) for possible causes and corrections.

After verifying that the Plugin is recognized by MySQL, create an InnoDB table as another confirmation of success.

9.3.3 Errors When Installing the InnoDB Plugin on Microsoft Windows

If MySQL or the associated Windows service can not start, or a post-startup `INSTALL PLUGIN` statement fails, look at the MySQL error log (usually named `machine_name.err` and located in the MySQL `data` directory) for the detailed error message. The log is in chronological order, so look at the end of the file. Try to resolve the problem based on other information in the message.

The following table outlines installation-related InnoDB Plugin error conditions or messages and possible solutions.

Error Condition or Message	Possible Solution
Can't open shared library <i>library_name</i>	Diagnose the cause from the following message details.
API version for STORAGE ENGINE plugin is too different	The version of the Plugin is not compatible with the version of the MySQL server. Consult the compatibility chart .
No such file or directory	Check that the file <code>ha_innodb_plugin.so</code> or <code>.dll</code> was copied to the correct location. Confirm that you specified the right file name (<code>ha_innodb_plugin.so</code> or <code>.dll</code> for the library from the InnoDB web site; <code>ha_innodb_plugin.so</code> or <code>.dll</code> for the library supplied along with MySQL 5.1.38 and up.)
Permission denied	Check that the folder and file access permissions are set properly. The <code>mysqld</code> process must have permission to read the file <code>ha_innodb_plugin.dll</code> and to read files in the plugin folder. On Windows XP, file permissions can be seen or changed by right-clicking a file and pressing Properties, and then the Security Tab. To see the Security Tab, you may need to adjust the Folder Options on the Control Panel to turn off "Use Simple File Sharing".
Can't open shared library 'ha_innodb_plugin.dll' (errno: 0)	Ensure that <code>ha_innodb_plugin.dll</code> is for the same system platform as <code>mysqld</code> . In particular, note that a 32-bit <code>mysqld</code> is unable to load a 64-bit plugin, and vice versa.

Note: The Information Schema tables are themselves plugins to the MySQL server, but they depend on having the InnoDB storage engine plugin installed as well. These tables will appear to be empty if the storage engine is not installed.

9.4 Building the InnoDB Plugin from Source Code

Sometimes, you may wish to build the plugin from the source code using special compilation options, or there might be no binary plugin available for your server platform. With the resulting special version of MySQL containing the new InnoDB functionality, it is not necessary to `INSTALL` any plugins or be concerned about startup parameters that preclude loading plugins.

To build the InnoDB Plugin from the source code, you also need the MySQL source code and some software tools. You should become familiar with the MySQL manual section on [Installing MySQL from Source](#).

The general steps for building MySQL from source, containing the InnoDB Plugin in place of the standard built-in InnoDB, are as follows:

- Download the MySQL source code.
- Download the InnoDB Plugin source code.
- Replace the source code for the built-in InnoDB with the InnoDB Plugin source tree.
- Compile MySQL as usual, generating a new `mysqld` executable file.

- Configure the MySQL server by editing the configuration file to use InnoDB as the default engine (if desired) and set appropriate configuration parameters to enable use of new InnoDB Plugin features.

The following sections detail these steps for Linux or Unix systems, and for Microsoft Windows.

9.4.1 Building the InnoDB Plugin on Linux or Unix

1. Download the MySQL source code, version 5.1.24 or later from <http://dev.mysql.com/downloads/mysql/5.1.html#source> and extract it.
2. Download the InnoDB Plugin source code from <http://dev.mysql.com/downloads/>.
3. Replace the contents of the `storage/innobase` directory in the MySQL source tree with the InnoDB Plugin source tree.



Note

In MySQL 5.1.38 and up, the MySQL source tree also contains a `storage/innodb_plugin` directory, but that does not affect this procedure. The source that you download from the InnoDB web site may contain additional changes and fixes.

4. Compile and build MySQL. Instead of building a dynamic InnoDB Plugin, it is advisable to build a version of MySQL that contains the InnoDB Plugin. This is because **a dynamic InnoDB Plugin must be built with exactly the same tools and options as the `mysqld` executable**, or spurious errors may occur. Example:

```
% wget ftp://ftp.easynet.be/mysql/Downloads/MySQL-5.1/mysql-5.1.37.tar.gz
% tar -zxf mysql-5.1.37.tar.gz
% cd mysql-5.1.37/storage
% wget http://dev.mysql.com/downloads/innodb_plugin/innodb_plugin-1.0.8.tar.gz
% tar -zxf innodb-1.0.8.tar.gz
% rm -fr innobase
% mv innodb-1.0.8 innobase
% cd ..
% ./configure --with-plugins=innobase
% make
```

5. Reconfigure the MySQL server by editing the `my.cnf` option file to use InnoDB as the default engine (if desired) and set appropriate configuration parameters to enable use of new InnoDB Plugin features, as described in section [Section 9.5, “Configuring the InnoDB Plugin”](#). In particular, we recommend that you set the following specific parameters as follows:

```
default_storage_engine=InnoDB
innodb_file_per_table=1
innodb_file_format=barracuda
innodb_strict_mode=1
```

6. If you build a version of MySQL that contains the InnoDB Plugin (`--with-plugins=innobase`), you do not have to tell MySQL to specify `ignore_builtin_innodb` or specify `plugin-load`, or issue any `INSTALL PLUGIN` statements. The `mysqld` executable that you compiled will contain the new InnoDB Plugin features.

Note: To fully exploit the performance improvements discussed in [Section 7.2, “Faster Locking for Improved Scalability”](#), the InnoDB Plugin source code and build process makes some compile-time tests of platform capabilities to automatically use instructions for atomic memory access where available. If this logic fails, you may need to contact MySQL support.

9.4.2 Building the InnoDB Plugin on Microsoft Windows

[Installing MySQL from Source on Windows](#) includes some information about building from source on Windows. The following discussion is specifically focused on building a version of MySQL containing the InnoDB Plugin.

You need the following tools:

- A compiler environment, one of the following:
 - Microsoft Visual C++ 2003
 - Microsoft Visual C++ 2005
 - Microsoft Visual C++ 2008 (Note: for building MySQL 5.1.31 or later)
 - Microsoft Visual C++ 2005 Express Edition (free of charge)
 - Download and install the [Microsoft Visual C++ 2005 Express Edition](#).
 - Download and install the [Windows Platform SDK](#).
 - Configure the Visual Studio Express Edition to use the Platform SDK according to the [instruction](#).
 - Microsoft Visual C++ 2008 Express Edition (free of charge, for building MySQL 5.1.31 or later)
 - Download and install the [Microsoft Visual C++ 2008 Express Edition](#). The Visual C++ 2008 Express Edition has already been integrated with the Windows SDK.
- [GNU Bison for Windows](#), a general-purpose parser generator that is largely compatible with Berkeley Yacc. This tool is used automatically as part of compiling and building MySQL. For most users, it is sufficient to download and run the “[complete package](#)” to install GNU Bison.
- [CMake 2.6.0 or later](#), a cross-platform make system that can generate MSVC project files.

In addition to installing these tools, you must also set CScript as the default Windows script host by executing the following command in the Command Prompt:

```
cscript //H:CScript
```

After you have installed and configured all the required tools, you may proceed with the compilation.

1. Download the MySQL source code, version 5.1.24 or later from the [MySQL website](#) and extract the source files.
2. Download the InnoDB plugin source code from the [MySQL download site](#).
3. Extract the files from the source code archives.
4. Replace the contents of the `storage\innobase` folder in the MySQL source tree with the InnoDB plugin source tree.

In MySQL 5.1.38 and up, the MySQL source tree also contains a `storage\innodb_plugin` directory, but that does not affect this procedure. The source that you download from the InnoDB web site may contain additional changes and fixes.

5. Compile and build MySQL under the Microsoft Visual Studio Command Prompt as follows:

Visual Studio 2003:

```
win\configure WITH_INNOBASE_STORAGE_ENGINE __NT__
win\build-vs7.bat
devenv mysql /build release /project ALL_BUILD
```

Visual Studio 2005:

```
win\configure WITH_INNOBASE_STORAGE_ENGINE __NT__
win\build-vs8.bat
devenv mysql /build release /project ALL_BUILD
```

Visual Studio 2008:

```
win\configure WITH_INNOBASE_STORAGE_ENGINE __NT__
win\build-vs9.bat
vcbuild mysql.sln "Release"
```

For the 64-bit version, use `win\build-vsN_x64.bat` instead of `win\build-vsN.bat`.

6. Install the compiled `mysqld.exe` from the `sql\release` folder of the source tree by doing one of the following:
 - a. Copy the `mysqld.exe` to the `bin` folder of an earlier MySQL 5.1 installation.
 - b. Make a distribution package and unpack it to the folder where MySQL will be installed. See the MySQL manual section on `make_win_bin_dist` — [Package MySQL Distribution as Zip Archive](#). Note that `scripts\make_win_bin_dist` requires the Cygwin environment.
7. Reconfigure the MySQL server by editing the `my.cnf` or `my.ini` option file to use InnoDB as the default engine (if desired) and set appropriate configuration parameters to enable use of new InnoDB Plugin features, as described in section [Section 9.5, “Configuring the InnoDB Plugin”](#). In particular, we recommend that you set the following specific parameters as follows:

```
default_storage_engine=InnoDB
innodb_file_per_table=1
innodb_file_format=barracuda
innodb_strict_mode=1
```

8. Since you built a version of MySQL that contains the InnoDB Plugin, you do not have to specify `ignore_builtin_innodb` or specify `plugin-load`, or issue any `INSTALL PLUGIN` statements. The `mysqld.exe` that you compiled contains the new InnoDB Plugin features.

9.5 Configuring the InnoDB Plugin

Because the MySQL server as distributed by MySQL includes a built-in copy of InnoDB, if you are using the dynamic InnoDB Plugin and have `INSTALLED` it into the MySQL server, you must always start the server with the option `ignore_builtin_innodb`, either in the option file or on the `mysqld` command line. Also, remember that the startup option `skip_grant_tables` prevents MySQL from loading any plugins. Neither of these options is needed when using a specialized version of MySQL that you build from source.

By default, the InnoDB Plugin does not create tables in a format that is incompatible with the built-in InnoDB in MySQL. Tables in the new format may be compressed, and they may store portions of long columns off-page, outside the B-tree nodes. You may wish to enable the creation of tables in the new format, using one of these techniques:

- Include `innodb_file_per_table=1` and `innodb_file_format=barracuda` in the `[mysqld]` section of the MySQL option file.
- Add `--innodb_file_per_table=1` and `--innodb_file_format=barracuda` to the `mysqld` command line.

- Issue the statements:

```
SET GLOBAL innodb_file_format=barracuda;  
SET GLOBAL innodb_file_per_table=ON;
```

in the MySQL client when running with `SUPER` privileges.

You may also want to enable the new InnoDB strict mode, which guards SQL or certain operational errors that otherwise generate warnings and possible unintended consequences of ignored or incorrect SQL commands or parameters. As described in [Section 8.5, “InnoDB Strict Mode”](#), the `GLOBAL` parameter `innodb_strict_mode` can be set `ON` or `OFF` in the same way as the parameters just mentioned. You can also use the command `SET SESSION innodb_strict_mode=mode` (where `mode` is `ON` or `OFF`) to enable or disable InnoDB strict mode on a per-session basis.

Take care when using new InnoDB configuration parameters or values that apply only when using the InnoDB Plugin. When the MySQL server encounters an unknown option, it fails to start and returns an error: `unknown variable`. This happens, for example, if you include the new parameter `innodb_file_format` when you start the MySQL server with the built-in InnoDB rather than the plugin. This can cause a problem if you accidentally use the built-in InnoDB after a system crash, because InnoDB crash recovery runs before MySQL checks the startup parameters. See [Section 11.4, “Possible Problems”](#) why this can be a problem. One safeguard is to specify the prefix `loose_` before the names of new options, so that if they are not recognized on startup, the server gives a warning instead of a fatal error.

9.6 Frequently Asked Questions about Plugin Installation

9.6.1 Should I use the InnoDB-supplied plugin or the one that is included with MySQL 5.1.38 or higher?

The Plugin that you download from the InnoDB web site should always be at the same level or newer than the shared library that is included with the MySQL distribution starting with version 5.1.38. To pick up the very latest fixes, download from the InnoDB site.

9.6.2 Why doesn't the MySQL service on Windows start after the replacement?

For the types of errors and how to diagnose them, see [Section 9.3.3, “Errors When Installing the InnoDB Plugin on Microsoft Windows”](#). Be especially careful that the `plugin-load` line in the option file does not get split across lines when you copy and paste from the README or this manual, which can produce an “unrecognized option” error in the error log.

9.6.3 The Plugin is installed... now what?

You automatically benefit from the “fast index creation” feature for every index you create on a large InnoDB table. If you switch to the “Barracuda” file format using the `innodb_file_format` option in combination with the `innodb_file_per_table` option, you can take advantage of other features such as table compression. For the full list of features, refer to [Section 1.2, “Features of the InnoDB Plugin”](#).

9.6.4 Once the Plugin is installed, is it permanent?

The Plugin must be loaded whenever the MySQL database server is started. As we saw earlier, there are several ways to configure MySQL to use the Plugin rather than the built-in InnoDB: in the option file, with `mysqld` command-line options, or with `INSTALL` statements after the server starts.

To ensure that you do not accidentally revert to the older InnoDB, be careful to carry any configuration file, command-line options, or post-startup commands forward in the future, such as when transitioning from a

Once the Plugin is installed, is it permanent?

development system to a test system, setting up a replication slave, or when writing new `mysqld` startup scripts.

Chapter 10 Upgrading the InnoDB Plugin

Table of Contents

10.1 Upgrading the Dynamic InnoDB Plugin	71
10.2 Upgrading a Statically Built InnoDB Plugin	71
10.3 Converting Compressed Tables Created Before Version 1.0.2	72

Thanks to the pluggable storage engine architecture of MySQL, upgrading the InnoDB Plugin should be a simple matter of shutting down MySQL, replacing a platform-specific executable file, and restarting the server. If you wish to upgrade and use your existing database, **it is essential to perform a “slow” shutdown**, or the new plugin may fail when merging buffered inserts or purging deleted records. If your database does not contain any compressed tables, you should be able to use your database with the newest InnoDB Plugin without problems after a slow shutdown.

However, if your database contains compressed tables, it may not be compatible with InnoDB Plugin 1.0.8. Because of an incompatible change introduced in InnoDB Plugin version 1.0.2, some compressed tables may need to be rebuilt, as noted in [Section 10.3, “Converting Compressed Tables Created Before Version 1.0.2”](#). Please follow these steps carefully.

You may, of course, rebuild your database using `mysqldump` or other methods. This may be a preferable approach if your database is small or there are many referential constraints among tables.

Note that once you have accessed your database with InnoDB Plugin 1.0.8, you should not try to use it with the Plugin prior to 1.0.2.

10.1 Upgrading the Dynamic InnoDB Plugin

Before shutting down the MySQL server containing the InnoDB Plugin, you must **enable “slow” shutdown**:

```
SET GLOBAL innodb_fast_shutdown=0;
```

For the details of the shutdown procedure, see the MySQL manual on [The Shutdown Process](#).

In the directory where the MySQL server looks for plugins, rename the executable file of the old InnoDB Plugin (`ha_innodb_plugin.so` or `ha_innodb_plugin.dll`), so that you can restore it later if needed. You may remove the file later. The plugin directory is specified by the system variable `plugin_dir`. The default location is usually the `lib/plugin` subdirectory of the directory specified by `basedir`.

Download a suitable package for your server platform, operating system and MySQL version. Extract the contents of the archive using `tar` or a similar tool for Linux and Unix, or Windows Explorer or WinZip or similar utility for Windows. Copy the file `ha_innodb_plugin.so` or `ha_innodb_plugin.dll` to the directory where the MySQL server looks for plugins.

Start the MySQL server. Follow the procedure in [Section 10.3, “Converting Compressed Tables Created Before Version 1.0.2”](#) to convert any compressed tables if needed.

10.2 Upgrading a Statically Built InnoDB Plugin

As with a dynamically installed InnoDB Plugin, you must perform a “slow” shutdown of the MySQL server. If you have built MySQL from source code and replaced the built-in InnoDB in MySQL with the InnoDB

Plugin in the source tree as discussed in [Section 9.4, “Building the InnoDB Plugin from Source Code”](#), you will have a special version of the `mysqld` executable that contains the InnoDB Plugin.

If you intend to upgrade to a dynamically linked InnoDB Plugin, you can follow the advice of [Section 11.3.4, “Uninstalling a Statically Built InnoDB Plugin”](#) and [Section 9.3, “Installing the Precompiled InnoDB Plugin as a Shared Library”](#).

If you intend to upgrade a statically built InnoDB Plugin to another statically built plugin, you will have to rebuild the `mysqld` executable, shut down the server, and replace the `mysqld` executable before starting the server.

Either way, please be sure to follow the instructions of [Section 10.3, “Converting Compressed Tables Created Before Version 1.0.2”](#) if any compressed tables were created.

10.3 Converting Compressed Tables Created Before Version 1.0.2

The InnoDB Plugin version 1.0.2 introduces an **incompatible change to the format of compressed tables**. This means that some compressed tables that were created with an earlier version of the InnoDB Plugin may need to be rebuilt with a bigger `KEY_BLOCK_SIZE` before they can be used.

If you must keep your existing database when you upgrade to InnoDB Plugin 1.0.2 or newer, you will need to perform a “slow” shutdown of MySQL running the previous version of the InnoDB Plugin. Following such a shutdown, and using the newer release of the InnoDB Plugin, you will need to determine which compressed tables need conversion and then follow a procedure to upgrade these tables. Because most users will not have tables where this process is required, this manual does not detail the procedures required. If you have created compressed tables with the InnoDB Plugin prior to release 1.0.2, you may want to contact MySQL support.

Chapter 11 Downgrading from the InnoDB Plugin

Table of Contents

11.1 Overview	73
11.2 The Built-in InnoDB, the Plugin and File Formats	73
11.3 How to Downgrade	74
11.3.1 Converting Tables	74
11.3.2 Adjusting the Configuration	74
11.3.3 Uninstalling a Dynamic Library	74
11.3.4 Uninstalling a Statically Built InnoDB Plugin	75
11.4 Possible Problems	75
11.4.1 Accessing <code>COMPRESSED</code> or <code>DYNAMIC</code> Tables	75
11.4.2 Issues with UNDO and REDO	76
11.4.3 Issues with the Doublewrite Buffer	76
11.4.4 Issues with the Insert Buffer	77

11.1 Overview

There are times when you might want to use the InnoDB Plugin with a given database, and then downgrade to the built-in InnoDB in MySQL. One reason to do this is because you want to take advantage of a new InnoDB Plugin feature (such as “Fast Index Creation”), but revert to the standard built-in InnoDB in MySQL for production operation.

If you have created new tables using the InnoDB Plugin, you may need to convert them to a format that the built-in InnoDB in MySQL can read. Specifically, if you have created tables that use `ROW_FORMAT=COMPRESSED` or `ROW_FORMAT=DYNAMIC` you must convert them to a different format, if you plan to access these tables with the built-in InnoDB in MySQL. If you do not do so, anomalous results may occur.

Although InnoDB checks the format of tables and database files (specifically `*.ibd` files) for compatibility, it is unable to start if there are buffered changes for “too new format” tables in the redo log or in the system tablespace. Thus it is important to carefully follow these procedures when downgrading from the InnoDB Plugin to the built-in InnoDB in MySQL, version 5.1.

This chapter describes the downgrade scenario, and the steps you should follow to ensure correct processing of your database.

11.2 The Built-in InnoDB, the Plugin and File Formats

Starting with version 5.0.21, the built-in InnoDB in MySQL checks the table type before opening a table. Until now, all InnoDB tables have been tagged with the same type, although some changes to the format have been introduced in MySQL versions 4.0, 4.1, and 5.0.

One of the important new features introduced with the InnoDB Plugin is support for identified file formats. This allows the InnoDB Plugin and versions of InnoDB since 5.0.21 to check for file compatibility. It also allows the user to preclude the use of features that would generate downward incompatibilities. By paying attention to the file format used, you can protect your database from corruptions, and ensure a smooth downgrade process.

In general, before using a database file created with the InnoDB Plugin with the built-in InnoDB in MySQL you should verify that the tablespace files (the `*.ibd` files) are compatible with the built-in InnoDB in MySQL. The InnoDB Plugin can read and write tablespaces in both the formats “Antelope” and

“Barracuda”. The built-in InnoDB can only read and write tablespaces in “Antelope” format. To make all tablespaces “legible” to the built-in InnoDB in MySQL, you should follow the instructions in [Section 11.3, “How to Downgrade”](#) to reformat all tablespaces to be in the “Antelope” format.

Generally, after a “slow” shutdown of the InnoDB Plugin (`innodb_fast_shutdown=0`), it should be safe to open the data files with the built-in InnoDB in MySQL. See [Section 11.4, “Possible Problems”](#) for a discussion of possible problems that can arise in this scenario and workarounds for them.

11.3 How to Downgrade

11.3.1 Converting Tables

The built-in InnoDB in MySQL can access only tables in the “Antelope” file format, that is, in the `REDUNDANT` or `COMPACT` row format. If you have created tables in `COMPRESSED` or `DYNAMIC` format, the corresponding tablespaces in the new “Barracuda” file format, and it is necessary to downgrade these tables.

First, identify the tables that require conversion, by executing this command:

```
SELECT table_schema, table_name, row_format
FROM information_schema.tables
WHERE engine='innodb'
AND row_format NOT IN ('Redundant', 'Compact');
```

Next, for each table that requires conversion, run the following command:

```
ALTER TABLE table_name ROW_FORMAT=COMPACT;
```

This command copies the table and its indexes to a new tablespace in the “Antelope” format. See [Chapter 2, *Fast Index Creation in the InnoDB Storage Engine*](#) for a discussion of exactly how such index creation operations are performed.

11.3.2 Adjusting the Configuration

Before you shut down the InnoDB Plugin and start the basic built-in InnoDB in MySQL, review the configuration files. Changes to the startup options do not take effect until the server is restarted, or the InnoDB Plugin is uninstalled and reinstalled.

The InnoDB Plugin introduces several configuration parameters that are not recognized by the built-in InnoDB in MySQL, including: `innodb_file_format`, `innodb_file_format_check`, and `innodb_strict_mode`. See [Section C.1, “New Parameters”](#) for a complete list of new configuration parameters in the InnoDB Plugin. You can include these parameters in the configuration file, only if you use the `loose_` form of the parameter names, so that the built-in InnoDB in MySQL can start.

If the InnoDB Plugin was installed as a dynamic plugin, the startup option `ignore_builtin_innodb` or `skip_innodb` must have been set to disable the built-in InnoDB in MySQL. These options must be removed, so that the built-in InnoDB in MySQL is enabled the next time the server is started.

If the InnoDB Plugin was loaded using `plugin-load` option. This option has to be removed, too.

In MySQL, configuration options can be specified in the `mysqld` command line or the option file (`my.cnf` or `my.ini`). See the MySQL manual on [Using Option Files](#) for more information.

11.3.3 Uninstalling a Dynamic Library

The following applies if the InnoDB Plugin was installed as a dynamic library with the `INSTALL PLUGIN` command, as described in [Section 9.3, “Installing the Precompiled InnoDB Plugin as a Shared Library”](#).

Issue the command `UNINSTALL PLUGIN` for every “plugin” supplied by the library `ha_innodb_plugin.so` (or `ha_innodb_plugin.dll` on Windows). Note that the following commands initiate a shutdown of the InnoDB Plugin:

```
SET GLOBAL innodb_fast_shutdown=0;
UNINSTALL PLUGIN INNODB;
UNINSTALL PLUGIN INNODB_CMP;
UNINSTALL PLUGIN INNODB_CMP_RESET;
UNINSTALL PLUGIN INNODB_CMPMEM;
UNINSTALL PLUGIN INNODB_CMPMEM_RESET;
UNINSTALL PLUGIN INNODB_TRX;
UNINSTALL PLUGIN INNODB_LOCKS;
UNINSTALL PLUGIN INNODB_LOCK_WAITS;
```

Due to MySQL Bug #33731, please ensure that the plugin definitions are actually deleted from the database, so that they are not loaded again:

```
SELECT * FROM mysql.plugin;
DELETE FROM mysql.plugin WHERE name='...';
```

Restart the server. For the details of the shutdown procedure, see the MySQL manual on [The Shutdown Process](#).

11.3.4 Uninstalling a Statically Built InnoDB Plugin

If you have built MySQL from source code and replaced the built-in InnoDB in MySQL with the InnoDB Plugin in the source tree as discussed in [Section 9.4, “Building the InnoDB Plugin from Source Code”](#), you have a special version of the `mysqld` executable that contains the InnoDB Plugin. To “uninstall” the InnoDB Plugin, you replace this executable with something that is built from an unmodified MySQL source code distribution.

Before shutting down the version of the MySQL server with built-in InnoDB Plugin, you must **enable “slow” shutdown**:

```
SET GLOBAL innodb_fast_shutdown=0;
```

For the details of the shutdown procedure, see the MySQL manual on [The Shutdown Process](#).

11.4 Possible Problems

Failure to follow the downgrading procedure described in [Section 11.3, “How to Downgrade”](#) may lead to compatibility issues when files written by the InnoDB Plugin are accessed by the built-in InnoDB in MySQL. This section describes some internal recovery algorithms, to help explain why it is important to follow the downgrade procedure described above. It discusses the issues that may arise, and covers possible ways to fix them.

A general fix is to install the plugin as described in [Chapter 9, *Installing the InnoDB Plugin*](#) and then follow the downgrading procedure described in [Section 11.3, “How to Downgrade”](#).

In the future, the file format management features described in [Chapter 4, *InnoDB File-Format Management*](#) will guard against the types of problems described in this section.

11.4.1 Accessing COMPRESSED or DYNAMIC Tables

The built-in InnoDB in MySQL can only open tables that were created in `REDUNDANT` or `COMPACT` format. Starting with MySQL version 5.0.21, an attempt to open a table in some other format results in `ERROR 1146 (42S02): Table 'test.t' doesn't exist`. Furthermore, a message “unknown table type” appears in the error log.

In the InnoDB Plugin, you may rebuild an incompatible table by issuing a statement `ALTER TABLE table_name ROW_FORMAT=COMPACT`.

11.4.2 Issues with UNDO and REDO

As noted in [Section 11.3, “How to Downgrade”](#), you should ensure a “slow” shutdown is done with the InnoDB Plugin, before running with the built-in InnoDB in MySQL, to clean up all buffers. To initiate a slow shutdown, execute the command `SET GLOBAL innodb_fast_shutdown=0` before initiating the shutdown of the InnoDB Plugin.

We recommend “slow” shutdown (`innodb_fast_shutdown=0`) because the InnoDB Plugin may write special records to the transaction undo log that cause problems if the built-in InnoDB in MySQL attempts to read the log. Specifically, these special records are written when a record in a `COMPRESSED` or `DYNAMIC` table is updated or deleted and the record contains columns stored off-page. The built-in InnoDB in MySQL cannot read these undo log records. Also, the built-in InnoDB in MySQL cannot roll back incomplete transactions that affect tables that it is unable to read (tables in `COMPRESSED` or `DYNAMIC` format).

Note that a “normal” shutdown does not necessarily empty the undo log. A normal shutdown occurs when `innodb_fast_shutdown=1`, the default. When InnoDB is shut down, some active transactions may have uncommitted modifications, or they may be holding a read view that prevents the purging of some version information from the undo log. The next time InnoDB is started after a normal shutdown (`innodb_fast_shutdown=1`), it rolls back any incomplete transactions and purge old version information. Therefore, it is important to perform a “slow” shutdown (`innodb_fast_shutdown=0`) as part of the downgrade process.

In case it is not possible to have the InnoDB Plugin clear the undo log, you can prevent the built-in InnoDB in MySQL from accessing the undo log by setting `innodb_force_recovery=3`. However, this is not a recommended approach, since in addition to preventing the purge of old versions, this recovery mode prevents the rollback of uncommitted transactions. For more information, see the MySQL manual on [Forcing InnoDB Recovery](#).

When it comes to downgrading, there are also considerations with respect to redo log information. For the purpose of crash recovery, InnoDB writes to the log files information about every modification to the data files. When recording changes to tables that were created in `DYNAMIC` or `COMPRESSED` format, the InnoDB Plugin writes redo log entries that cannot be recognized by the built-in InnoDB in MySQL. The built-in InnoDB in MySQL refuses to start if it sees any unknown entries in the redo log.

When InnoDB is *shut down cleanly*, it flushes all unwritten changes from the buffer pool to the data files and makes a checkpoint in the redo log. When InnoDB is subsequently restarted, it scans the redo log starting from the last checkpoint. After a clean shutdown, InnoDB crash recovery only then sees the end-of-log marker in the redo log. In this case, the built-in InnoDB in MySQL would not see any unrecognizable redo log entries. This is a second reason why you should ensure a clean, slow shutdown of MySQL (`innodb_fast_shutdown=0`) before you attempt a downgrade.

In an emergency, you may prevent the redo log scan and the crash recovery from the redo log by setting the parameter `innodb_force_recovery=6`. However, this is **strongly discouraged**, because may lead into severe corruption. See the MySQL manual on [Forcing InnoDB Recovery](#) for more information.

11.4.3 Issues with the Doublewrite Buffer

InnoDB uses a novel file flush technique called doublewrite. Before writing pages to a data file, InnoDB first writes them to a contiguous area called the doublewrite buffer. Only after the write and the flush to the doublewrite buffer have completed does InnoDB write the pages to their proper positions in the data file. If the operating system crashes in the middle of a page write, InnoDB can later find a good copy of the page from the doublewrite buffer during recovery.

The doublewrite buffer may also contain compressed pages. However, the built-in InnoDB in MySQL cannot recognize such pages, and it assumes that compressed pages in the doublewrite buffer are corrupted. It also wrongly assumes that the tablespace (the `.ibd` file) consists of 16K byte pages. Thus, you may find InnoDB warnings in the error log of the form “a page in the doublewrite buffer is not within space bounds”.

The doublewrite buffer is not scanned after a *clean shutdown*. In an emergency, you may prevent crash recovery by setting `innodb_force_recovery=6`. However, this is **strongly discouraged**, because it may lead into severe corruption. For more information, see the MySQL manual on [Forcing InnoDB Recovery](#)

11.4.4 Issues with the Insert Buffer

Secondary indexes are usually nonunique, and insertions into secondary indexes happen in a relatively random order. This would cause a lot of random disk I/O operations without a special mechanism used in InnoDB called the insert buffer.

When a record is inserted into a nonunique secondary index page that is not in the buffer pool, InnoDB inserts the record into a special B-tree: the insert buffer. Periodically, the insert buffer is merged into the secondary index trees in the database. A merge also occurs whenever a secondary index page is loaded to the buffer pool.

A “normal” shutdown does not clear the insert buffer. A normal shutdown occurs when `innodb_fast_shutdown=1`, the default. If the insert buffer is not empty when the InnoDB Plugin is shut down, it may contain changes for tables in `DYNAMIC` or `COMPRESSED` format. Thus, starting the built-in InnoDB in MySQL on the data files may lead into a crash if the insert buffer is not empty.

A “slow” shutdown merges all changes from the insert buffer. To initiate a slow shutdown, execute the command `SET GLOBAL innodb_fast_shutdown=0` before initiating the shutdown of the InnoDB Plugin.

To disable insert buffer merges, you may set `innodb_force_recovery=4` so that you can back up the uncompressed tables with the built-in InnoDB in MySQL. Be sure not to use any `WHERE` conditions that would require access to secondary indexes. For more information, see the MySQL manual on [Forcing InnoDB Recovery](#)

In the InnoDB Plugin 1.0.3 and later, you can disable the buffering of new operations by setting the parameter `innodb_change_buffering`. See [Section 7.4, “Controlling InnoDB Insert Buffering”](#) for details.

Chapter 12 InnoDB Plugin Change History

Table of Contents

12.1 Changes in InnoDB Plugin 1.0.9 and Higher	79
12.2 Changes in InnoDB Plugin 1.0.8 (May, 2010)	79
12.3 Changes in InnoDB Plugin 1.0.7 (April, 2010)	79
12.4 Changes in InnoDB Plugin 1.0.6 (November 27, 2009)	80
12.5 Changes in InnoDB Plugin 1.0.5 (November 18, 2009)	80
12.6 Changes in InnoDB Plugin 1.0.4 (August 11, 2009)	81
12.7 Changes in InnoDB Plugin 1.0.3 (March 11, 2009)	82
12.8 Changes in InnoDB Plugin 1.0.2 (December 1, 2008)	83
12.9 Changes in InnoDB Plugin 1.0.1 (May 8, 2008)	83
12.10 Changes in InnoDB Plugin 1.0.0 (April 15, 2008)	84

The complete change history of the InnoDB Plugin can be viewed in the file [ChangeLog](#) that is included in the source and binary distributions.

12.1 Changes in InnoDB Plugin 1.0.9 and Higher

With the tighter integration of InnoDB into the MySQL server starting in MySQL 5.1, you can find recent InnoDB change log entries in the [MySQL 5.1 Release Notes](#).

12.2 Changes in InnoDB Plugin 1.0.8 (May, 2010)

12.3 Changes in InnoDB Plugin 1.0.7 (April, 2010)

Improved crash recovery performance.

Fixed MySQL Bug #52102: InnoDB Plugin shows performance drop comparing to builtin InnoDB on Windows only. Disabled Windows atomics by default.

Fixed MySQL Bug #51378: Init 'ref_length' to correct value, in case of an out of bound MySQL primary_key.

Made `SHOW ENGINE INNODB MUTEX STATUS` display `SUM(os_waits)` for the buffer pool block mutexes and locks.

Fixed `ALTER TABLE ... IMPORT TABLESPACE` of compressed tables.

Fixed MySQL Bug #49535: Available memory check slows down crash recovery tens of times.

Let the master thread sleep if the amount of work to be done is calibrated as taking less than a second.

Fixed MySQL Bug #49001: SHOW INNODB STATUS deadlock info incorrect when deadlock detection aborts.

Fixed MySQL Bug #35077: Very slow DROP TABLE (ALTER TABLE, OPTIMIZE TABLE) on compressed tables.

Fixed MySQL Bug #49497: Error 1467 (`ER_AUTOINC_READ_FAILED`) on inserting a negative value.

Do not merge buffered inserts to compressed pages before the redo log has been applied in crash recovery.

Do not attempt to access a clustered index record that has been marked for deletion, On the [READ UNCOMMITTED](#) isolation level. In previous versions, the InnoDB would attempt to retrieve a previous version of the record in this case.

Fixed an uninitialized access to `block->is_hashed`, when disabling the adaptive hash index.

Fixed MySQL Bug #46193>: Crash when accessing tables after enabling `innodb_force_recovery` option.

Fixed MySQL Bug #49238: Creating / Dropping a temporary table while at 1023 transactions will cause assert.

Display the `zlib` version number at startup.

12.4 Changes in InnoDB Plugin 1.0.6 (November 27, 2009)

Fixed MySQL Bug #48782: On lock wait timeout, `CREATE INDEX` attempts `DROP TABLE`.

Report duplicate table names to the client connection, not to the error log.

Allow `CREATE INDEX` to be interrupted.

Fixed MySQL Bug #47167: InnoDB Plugin "`set global innodb_file_format_check`" cannot set value by User-Defined Variable.

Fixed MySQL Bug #45992: InnoDB memory not freed after shutdown; and MySQL Bug #46656: InnoDB Plugin memory leaks (`Valgrind`).

12.5 Changes in InnoDB Plugin 1.0.5 (November 18, 2009)

Clean up after a crash during `DROP INDEX`. When InnoDB crashes while dropping an index, ensure that the index will be completely dropped during crash recovery.

When a secondary index exists in the MySQL `.frm` file but not in the InnoDB data dictionary, return an error instead of letting an assertion fail in `index_read`.

Prevent the reuse of tablespace identifiers after InnoDB has crashed during table creation. Also, refuse to start if files with duplicate tablespace identifiers are encountered.

Fixed MySQL Bug #47055: Unconditional exit on `ERROR_WORKING_SET_QUOTA 1453 (0x5AD)` for InnoDB backend.

Fixed MySQL Bug #37232: InnoDB might get too many read locks for `DML` with repeatable-read.

Fixed MySQL Bug #31183: Tablespace full problems not reported in error log; error message unclear.

Modified `innodb-zip.test` so that the test will pass with `zlib 1.2.3.3`. Apparently, the `zlib` function `compressBound()` has been slightly changed, and the maximum record size of a table with 1K compressed page size has been reduced by one byte.

Fixed a regression introduced by the fix for MySQL Bug #26316.

Fixed MySQL Bug #44571: InnoDB Plugin crashes on `ADD INDEX`.

Fixed a bug in the merge sort that can corrupt indexes in fast index creation.

Introduced the settable global variables `innodb_old_blocks_pct` and `innodb_old_blocks_time` for controlling the buffer pool eviction policy, making it possible to tune the buffer pool LRU eviction policy to be more resistant against index scans. See [Section 7.14, “Making Buffer Cache Scan Resistant”](#).

Fixed MySQL Bug #42885: `buf_read_ahead_random`, `buf_read_ahead_linear` counters, thread wakeups. See [Section 8.9, “More Read Ahead Statistics”](#).

Fixed MySQL Bug #46650: InnoDB assertion `autoinc_lock == lock` in `lock_table_remove_low` on `INSERT SELECT`.

Fixed MySQL Bug #46657: InnoDB Plugin: invalid read in `index_merge_innodb` test (`Valgrind`).

Fixed MySQL Bug #42829: binlogging enabled for all schemas regardless of `binlog-db-db` / `binlog-ignore-db`.

12.6 Changes in InnoDB Plugin 1.0.4 (August 11, 2009)

Enabled inlining of functions and prefetch with Sun Studio.

Changed the defaults for `innodb_sync_spin_loops` from 20 to 30 and `innodb_spin_wait_delay` from 5 to 6.

Implemented adaptive flushing of dirty pages, which uses heuristics to avoid I/O bursts at checkpoint. A new parameter `innodb_adaptive_flushing` is added to control whether the new flushing algorithm should be used. See [Section 7.11, “Controlling the Flushing Rate of Dirty Pages”](#).

Implemented I/O capacity tuning. A new parameter `innodb_io_capacity` is added to control the master threads I/O rate. (To preserve the former behavior, set this parameter to a value of 100.) The `ibuf` merge is also changed from synchronous to asynchronous. See [Section 7.10, “Controlling the Master Thread I/O Rate”](#).

Introduced the `PAUSE` instruction inside spin-loop where available. See [Section 7.12, “Using the PAUSE instruction in InnoDB spin loops”](#).

Fixed a crash on `SET GLOBAL innodb_file_format=DEFAULT` or `SET GLOBAL innodb_file_format_check=DEFAULT`.

Changed the default values for `innodb_max_dirty_pages_pct`, `innodb_additional_mem_pool_size`, `innodb_buffer_pool_size`, and `innodb_log_buffer_size`.

Enabled group commit functionality that was broken in 5.0 when distributed transactions were introduced. See [Section 7.9, “Group Commit”](#).

Enabled the functionality of having multiple background threads, with two new configuration parameters, `innodb_read_io_threads` and `innodb_write_io_threads`. The Windows only parameter `innodb_file_io_threads` has been removed. See [Section 7.8, “Multiple Background I/O Threads”](#).

Changed the linear read ahead algorithm and disabled random read ahead. Also introduced a new configuration parameter `innodb_read_ahead_threshold` to control the sensitivity of the linear read ahead. See [Section 7.7, “Changes in the Read Ahead Algorithm”](#).

Standardized comments that allow the extraction of documentation from code base with the `Doxygen` tool.

Fixed a bug that could cause failures in secondary index lookups in consistent reads right after crash recovery.

Corrected the estimation of space needed on a compressed page when performing an update by `delete-and-insert`.

Removed the statically linked copies of the `zlib` and `strings` libraries from the binary Windows plugin. Invoke the copies of these libraries in the `mysqld` executable, like the binary InnoDB Plugin does on other platforms.

Trimmed the output of `SHOW ENGINE INNODB MUTEX`. See [Section 8.8, “More Compact Output of SHOW ENGINE INNODB MUTEX”](#).

On Microsoft Windows, make use of atomic memory access to implement mutexes and rw-locks more efficiently. On Sun Solaris 10, if GCC built-in functions for atomic memory access are unavailable, use library functions instead. See [Section 7.2, “Faster Locking for Improved Scalability”](#).

Fixed MySQL Bug #44032: in `ROW_FORMAT=REDUNDANT`, update UTF-8 `CHAR` to/from `NULL` is not in-place.

Fixed MySQL Bug #43660: `SHOW INDEXES/ANALYZE` does not update cardinality for indexes of InnoDB table.

Made the parameter `innodb_change_buffering` settable by `mysqld` start-up option. Due to a programming mistake, it was only possible to set this parameter by the `SET GLOBAL` command in InnoDB Plugin 1.0.3.

Added a parameter `innodb_spin_wait_delay` for controlling the polling of mutexes and rw-locks. See [Section 7.13, “Control of Spin Lock Polling”](#).

In consistent reads, issue an error message on attempts to use newly created indexes that may lack required history. See [Section 2.6, “Limitations”](#).

12.7 Changes in InnoDB Plugin 1.0.3 (March 11, 2009)

Improved the scalability of InnoDB on multi-core CPUs. See [Section 7.2, “Faster Locking for Improved Scalability”](#).

Added a parameter `innodb_change_buffering` for controlling the insert buffering. See [Section 7.4, “Controlling InnoDB Insert Buffering”](#).

Added a parameter `innodb_use_sys_malloc` for using an operating system memory allocation rather than the InnoDB internal memory allocator. See [Section 7.3, “Using Operating System Memory Allocators”](#).

Made it possible to dynamically enable or disable adaptive hash indexing. See [Section 7.5, “Controlling Adaptive Hash Indexing”](#).

Changed the default value of `innodb_thread_concurrency` from 8 to 0, for unlimited concurrency by default. See [Section 7.6, “Changes Regarding Thread Concurrency”](#).

Fixed an issue that the InnoDB Plugin fails if `innodb_buffer_pool_size` is defined bigger than 4095M on 64-bit Windows.

Fixed MySQL Bug #41676: Table names are case insensitive in locking.

Fixed MySQL Bug #41904: Create unique index problem.

Fixed MySQL Bug #43043: Crash on BLOB delete operation.

Fixed a bug in recovery when dropping incomplete indexes left behind by fast index creation.

Fixed a crash bug when all rows of a compressed table are deleted.

Fixed a corruption bug when a table is dropped on a busy system that contains compressed tables.

Fixed an assertion failure involving the variable `ut_total_allocated_memory` that was caused by unprotected access during fast index creation.

12.8 Changes in InnoDB Plugin 1.0.2 (December 1, 2008)

Implemented the dynamic plugin (`ha_innodb.dll`) on Windows.

Added a parameter `innodb_stats_sample_pages` for controlling the index cardinality estimates.

Made `innodb_stats_on_metadata` a settable global parameter. (MySQL Bug #38189)

Made `innodb_lock_wait_timeout` a settable session parameter. (MySQL Bug #36285)

Fixed bugs related to off-page columns (see [Section 5.3, "DYNAMIC Row Format"](#)).

Fixed various bugs related to compressed tables. This includes MySQL Bug #36172, a possible but rare corruption, and an incompatible file format change relating to very long rows in compressed tables, and to off-page storage of long column values.

Fixed a bug in crash recovery which was a side effect of incorrect implementation of the system tablespace tagging.

Fixed MySQL bugs related to `auto_increment` columns: Bug #26316, Bug #35498, Bug #35602, Bug #36411, Bug #37531, Bug #37788, Bug #38839, Bug #39830, Bug #40224.

Fixed some race conditions, hangs or crashes related to `INFORMATION_SCHEMA` tables, fast index creation, and to the recovery of `PREPARED` transactions.

Fixed crashes on `DROP TABLE` or `CREATE TABLE` when there are `FOREIGN KEY` constraints. (MySQL Bug #38786)

Fixed a crash caused by a conflict between `TRUNCATE TABLE` and `LOCK TABLES`. (MySQL Bug #38231)

Fixed MySQL Bug #39939: `DROP TABLE` or `DISCARD TABLESPACE` takes a long time.

Fixed MySQL Bug #40359: InnoDB plugin error/warning message during shutdown.

Fixed MySQL Bug #40360: Binlog related errors with binlog off.

Applied all changes from MySQL through version 5.1.30.

12.9 Changes in InnoDB Plugin 1.0.1 (May 8, 2008)

Fixed bugs related to the packaging of the InnoDB Plugin: MySQL Bug #36222, Bug #36434.

Fixed crash bugs related to the new features of the InnoDB Plugin: MySQL Bug #36169, Bug #36310.

Implemented the system tablespace tagging discussed in [Section 4.4.1, "Startup File Format Compatibility Checking"](#).

Applied all changes from MySQL through version 5.1.25.

12.10 Changes in InnoDB Plugin 1.0.0 (April 15, 2008)

The initial release of the InnoDB Plugin is based on the built-in InnoDB in MySQL version 5.1. See [Section 1.2, “Features of the InnoDB Plugin”](#) for the main features.

Appendix A Third-Party Software

Table of Contents

A.1 Performance Patches from Google	85
A.2 Multiple Background I/O Threads Patch from Percona	86
A.3 Performance Patches from Sun Microsystems	86

Innobase Oy acknowledges that certain Third Party and Open Source software has been used to develop or is incorporated in InnoDB (including the InnoDB Plugin). This appendix includes required third-party license information.

A.1 Performance Patches from Google

Innobase Oy gratefully acknowledges the following contributions from Google, Inc. to improve InnoDB performance:

- Replacing InnoDB's use of Pthreads mutexes with calls to GCC atomic builtins, as discussed in [Section 7.2, "Faster Locking for Improved Scalability"](#). This change means that InnoDB mutex and rw-lock operations take less CPU time, and improves throughput on those platforms where the atomic operations are available.
- Controlling master thread I/O rate, as discussed in [Section 7.10, "Controlling the Master Thread I/O Rate"](#). The master thread in InnoDB is a thread that performs various tasks in the background. Historically, InnoDB has used a hard coded value as the total I/O capacity of the server. With this change, user can control the number of I/O operations that can be performed per second based on their own workload.

Changes from the Google contributions were incorporated in the following source code files: `btr0cur.c`, `btr0sea.c`, `buf0buf.c`, `buf0buf.ic`, `ha_innodb.cc`, `log0log.c`, `log0log.h`, `os0sync.h`, `row0sel.c`, `srv0srv.c`, `srv0srv.h`, `srv0start.c`, `sync0arr.c`, `sync0rw.c`, `sync0rw.h`, `sync0rw.ic`, `sync0sync.c`, `sync0sync.h`, `sync0sync.ic`, and `univ.i`.

These contributions are incorporated subject to the conditions contained in the file `COPYING.Google`, which are reproduced here.

```
Copyright (c) 2008, 2009, Google Inc.  
All rights reserved.
```

```
Redistribution and use in source and binary forms, with or without  
modification, are permitted provided that the following conditions  
are met:
```

- * Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- * Neither the name of the Google Inc. nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

```
THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS  
"AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT  
LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS
```

```
FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE
COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT,
INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING,
BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES;
LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER
CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN
ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE
POSSIBILITY OF SUCH DAMAGE.
```

A.2 Multiple Background I/O Threads Patch from Percona

Innodb Oy gratefully acknowledges the contribution of Percona, Inc. to improve InnoDB performance by implementing configurable background threads, as discussed in [Section 7.8, "Multiple Background I/O Threads"](#). InnoDB uses background threads to service various types of I/O requests. The change provides another way to make InnoDB more scalable on high end systems.

Changes from the Percona, Inc. contribution were incorporated in the following source code files: [ha_innodb.cc](#), [os0file.c](#), [os0file.h](#), [srv0srv.c](#), [srv0srv.h](#), and [srv0start.c](#).

This contribution is incorporated subject to the conditions contained in the file [COPYING.Percona](#), which are reproduced here.

```
Copyright (c) 2008, 2009, Percona Inc.
All rights reserved.
```

```
Redistribution and use in source and binary forms, with or without
modification, are permitted provided that the following conditions
are met:
```

- * Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- * Neither the name of the Percona Inc. nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

```
THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS
"AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT
LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS
FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE
COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT,
INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING,
BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES;
LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER
CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN
ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE
POSSIBILITY OF SUCH DAMAGE.
```

A.3 Performance Patches from Sun Microsystems

Innodb Oy gratefully acknowledges the following contributions from Sun Microsystems, Inc. to improve InnoDB performance:

- Introducing the PAUSE instruction inside spin loops, as discussed in [Section 7.12, "Using the PAUSE instruction in InnoDB spin loops"](#). This change increases performance in high concurrency, CPU-bound workloads.

- Enabling inlining of functions and prefetch with Sun Studio.

Changes from the Sun Microsystems, Inc. contribution were incorporated in the following source code files: [univ.i](#), [utOut.c](#), and [utOut.h](#).

This contribution is incorporated subject to the conditions contained in the file [COPYING.Sun_Microsystems](#), which are reproduced here.

```
Copyright (c) 2009, Sun Microsystems, Inc.  
All rights reserved.
```

```
Redistribution and use in source and binary forms, with or without  
modification, are permitted provided that the following conditions  
are met:
```

- * Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- * Neither the name of Sun Microsystems, Inc. nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

```
THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS  
"AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT  
LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS  
FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE  
COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT,  
INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING,  
BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES;  
LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER  
CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT  
LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN  
ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE  
POSSIBILITY OF SUCH DAMAGE.
```

Appendix B Using the InnoDB Plugin with MySQL 5.1.30 or Earlier

Up to MySQL 5.1.30, the InnoDB Plugin replaced the built-in InnoDB in MySQL when the server was started with the option `skip_innodb`. Due to [MySQL Bug #42610](#), it was impossible to replace the built-in InnoDB in MySQL with a plugin in MySQL 5.1.31 and 5.1.32. MySQL 5.1.33 introduced the option `ignore_builtin_innodb` to allow InnoDB Plugin installation in the binary release.

Up to MySQL 5.1.30, installing the binary InnoDB Plugin requires that MySQL be shut down and restarted after issuing the `INSTALL PLUGIN` statements. This is because the `INSTALL PLUGIN` statement started the plugin with default options. The options would only be read from the option file (`my.cnf` or `my.ini`) after restarting the server. The InnoDB Plugin worked around this limitation by copying parameters from the internal data structures of the built-in InnoDB in MySQL. Beginning with MySQL 5.1.33, the `INSTALL PLUGIN` statement re-reads the option file and passes all options to the plugin, even those that are not recognized by the built-in InnoDB in MySQL.

To use the binary InnoDB Plugin with MySQL 5.1.30 or earlier, you may follow the instructions given in [Section 9.3, “Installing the Precompiled InnoDB Plugin as a Shared Library”](#), with one change: *Replace the option `ignore_builtin_innodb` with `skip_innodb`*. The general steps for dynamically installing the binary InnoDB Plugin are thus as follows:

- Make sure the MySQL server is not running, using a “slow” shutdown.
- Prepend each InnoDB option with `loose_`, e.g., `loose_innodb_file_per_table` instead of `innodb_file_per_table`, so that MySQL starts even when InnoDB is unavailable.
- Add `skip_innodb` and `default_storage_engine=MyISAM` to the options, to prevent the built-in InnoDB from starting.
- Start the MySQL server.
- `INSTALL` the InnoDB Plugin and the Information Schema tables, using the supplied script or equivalent commands.
- Verify the installation of the plugins.
- Shut down and reconfigure the MySQL server by editing the appropriate configuration file to use InnoDB as the default engine (if desired), and set appropriate configuration parameters to enable use of new InnoDB Plugin features.

This change only affects the binary distributions of MySQL and InnoDB Plugin. The procedure for building from source code is unchanged.

Appendix C List of Parameters Changed in the InnoDB Plugin 1.0

Table of Contents

C.1 New Parameters	91
C.2 Deprecated Parameters	93
C.3 Parameters with New Defaults	93

C.1 New Parameters

Throughout the course of development, the InnoDB Plugin has introduced new configuration parameters. The following table summarizes those parameters:

Table C.1 InnoDB Plugin New Parameter Summary

Name	Cmd-Line	Option File	System Var	Scope	Dynamic	Default
<code>innodb_adaptive_flushing</code>	YES	YES	YES	GLOBAL	YES	TRUE
<code>innodb_change_buffering</code>	YES	YES	YES	GLOBAL	YES	<code>inserts</code>
<code>innodb_file_format</code>	YES	YES	YES	GLOBAL	YES	Antelope
<code>innodb_file_format_check</code>	YES	YES	YES	GLOBAL	YES	ON
<code>innodb_io_capacity</code>	YES	YES	YES	GLOBAL	YES	200
<code>innodb_old_blocks_pct</code>	YES	YES	YES	GLOBAL	YES	37
<code>innodb_old_blocks_time</code>	YES	YES	YES	GLOBAL	YES	0
<code>innodb_read_ahead_threshold</code>	YES	YES	YES	GLOBAL	YES	56
<code>innodb_read_io_threads</code>	YES	YES	YES	GLOBAL	NO	4
<code>innodb_spin_wait_delay</code>	YES	YES	YES	GLOBAL	YES	6
<code>innodb_stats_sample_pages</code>	YES	YES	YES	GLOBAL	YES	8
<code>innodb_strict_mode</code>	YES	YES	YES	GLOBAL SESSION	YES	FALSE
<code>innodb_use_sys_malloc</code>	YES	YES	YES	GLOBAL	NO	TRUE
<code>innodb_write_io_threads</code>	YES	YES	YES	GLOBAL	NO	4

`innodb_adaptive_flushing`

Whether InnoDB uses a new algorithm to estimate the required rate of flushing. The default value is `TRUE`. This parameter was added in InnoDB Plugin 1.0.4. See [Section 7.11, “Controlling the Flushing Rate of Dirty Pages”](#) for more information.

`innodb_change_buffering`

Whether InnoDB performs insert buffering. The default value is `"inserts"` (buffer insert operations). This parameter was added in InnoDB Plugin 1.0.3. See [Section 7.4, “Controlling InnoDB Insert Buffering”](#) for more information.

`innodb_file_format`

Whether to enable the new “Barracuda” file format. The default value is “Antelope”. This parameter was added in InnoDB Plugin 1.0.1. See [Section 4.3, “Enabling File Formats”](#) for more information.

<code>innodb_file_format_check</code>	Whether InnoDB performs file format compatibility checking when opening a database. The default value is <code>ON</code> . This parameter was added in InnoDB Plugin 1.0.1. See Section 4.4.1, “Startup File Format Compatibility Checking” for more information.
<code>innodb_io_capacity</code>	The number of I/O operations that can be performed per second. The allowable value range is any number 100 or greater, and the default value is 200. This parameter was added in InnoDB Plugin 1.0.4. To reproduce the earlier behavior, use a value of 100. See Section 7.10, “Controlling the Master Thread I/O Rate” for more information.
<code>innodb_old_blocks_pct</code>	Controls the desired percentage of “old” blocks in the LRU list of the buffer pool. The default value is 37 and the allowable value range is 5 to 95. This parameter was added in InnoDB Plugin 1.0.5. See Section 7.14, “Making Buffer Cache Scan Resistant” for more information.
<code>innodb_old_blocks_time</code>	The time in milliseconds since the first access to a block during which it can be accessed again without being made “young”. The default value is 0 which means that blocks are moved to the “young” end of the LRU list at the first access. This parameter was added in InnoDB Plugin 1.0.5. See Section 7.14, “Making Buffer Cache Scan Resistant” for more information.
<code>innodb_read_ahead_threshold</code>	Control the sensitivity of the linear read ahead. The allowable value range is 0 to 64 and the default value is 56. This parameter was added in InnoDB Plugin 1.0.4. See Section 7.7, “Changes in the Read Ahead Algorithm” for more information.
<code>innodb_read_io_threads</code>	The number of background I/O threads used for reads. The allowable value range is 1 to 64 and the default value is 4. This parameter was added in InnoDB Plugin 1.0.4. See Section 7.8, “Multiple Background I/O Threads” for more information.
<code>innodb_spin_wait_delay</code>	Maximum delay between polling for a spin lock. The allowable value range is 0 (meaning unlimited) or positive integers and the default value is 6. This parameter was added in InnoDB Plugin 1.0.4. See Section 7.13, “Control of Spin Lock Polling” for more information.
<code>innodb_stats_sample_pages</code>	The number of index pages to sample when calculating statistics. The allowable value range is 1-unlimited and the default value is 8. This parameter was added in InnoDB Plugin 1.0.2. See Section 8.6, “Controlling Optimizer Statistics Estimation” for more information.
<code>innodb_strict_mode</code>	Whether InnoDB raises error conditions in certain cases, rather than issuing a warning. The default value is <code>OFF</code> . This parameter was added in InnoDB Plugin 1.0.2. See Section 8.5, “InnoDB Strict Mode” for more information.
<code>innodb_use_sys_malloc</code>	Whether InnoDB uses its own memory allocator or an allocator of the operating system. The default value is <code>ON</code> (use an allocator of the underlying system). This parameter was added

in InnoDB Plugin 1.0.3. See [Section 7.3, “Using Operating System Memory Allocators”](#) for more information.

`innodb_write_io_threads`

The number of background I/O threads used for writes. The allowable value range is 1 to 64 and the default value is 4. This parameter was added in InnoDB Plugin 1.0.4. See [Section 7.8, “Multiple Background I/O Threads”](#) for more information.

C.2 Deprecated Parameters

Beginning in InnoDB Plugin 1.0.4 the following configuration parameter has been removed:

`innodb_file_io_threads`

This parameter has been replaced by two new parameters `innodb_read_io_threads` and `innodb_write_io_threads`. See [Section 7.8, “Multiple Background I/O Threads”](#) for more information.

C.3 Parameters with New Defaults

For better out-of-the-box performance, InnoDB Plugin 1.0.4 changes the default values for the following configuration parameters:

Table C.2 InnoDB Plugin Parameters with New Defaults

Name	Old Default	New Default
<code>innodb_additional_mem_pool_size</code>	1MB	8MB
<code>innodb_buffer_pool_size</code>	8MB	128MB
<code>innodb_log_buffer_size</code>	1MB	8MB
<code>innodb_max_dirty_pages_pct</code>	90	75
<code>innodb_sync_spin_loops</code>	20	30
<code>innodb_thread_concurrency</code>	8	0

InnoDB Glossary

These terms are commonly used in information about the InnoDB storage engine.

A

.ARM file

Metadata for ARCHIVE tables. Contrast with **.ARZ file**. Files with this extension are always included in backups produced by the `mysqlbackup` command of the **MySQL Enterprise Backup** product. See Also [.ARZ file](#), [MySQL Enterprise Backup](#), [mysqlbackup command](#).

.ARZ file

Data for ARCHIVE tables. Contrast with **.ARM file**. Files with this extension are always included in backups produced by the `mysqlbackup` command of the **MySQL Enterprise Backup** product. See Also [.ARM file](#), [MySQL Enterprise Backup](#), [mysqlbackup command](#).

ACID

An acronym standing for atomicity, consistency, isolation, and durability. These properties are all desirable in a database system, and are all closely tied to the notion of a **transaction**. The transactional features of InnoDB adhere to the ACID principles.

Transactions are **atomic** units of work that can be **committed** or **rolled back**. When a transaction makes multiple changes to the database, either all the changes succeed when the transaction is committed, or all the changes are undone when the transaction is rolled back.

The database remains in a consistent state at all times -- after each commit or rollback, and while transactions are in progress. If related data is being updated across multiple tables, queries see either all old values or all new values, not a mix of old and new values.

Transactions are protected (isolated) from each other while they are in progress; they cannot interfere with each other or see each other's uncommitted data. This isolation is achieved through the **locking** mechanism. Experienced users can adjust the **isolation level**, trading off less protection in favor of increased performance and **concurrency**, when they can be sure that the transactions really do not interfere with each other.

The results of transactions are durable: once a commit operation succeeds, the changes made by that transaction are safe from power failures, system crashes, race conditions, or other potential dangers that many non-database applications are vulnerable to. Durability typically involves writing to disk storage, with a certain amount of redundancy to protect against power failures or software crashes during write operations. (In InnoDB, the **doublewrite buffer** assists with durability.)

See Also [atomic](#), [commit](#), [concurrency](#), [doublewrite buffer](#), [isolation level](#), [locking](#), [rollback](#), [transaction](#).

adaptive flushing

An algorithm for **InnoDB** tables that smooths out the I/O overhead introduced by **checkpoints**. Instead of **flushing** all modified **pages** from the **buffer pool** to the **data files** at once, MySQL periodically flushes small sets of modified pages. The adaptive flushing algorithm extends this process by estimating the optimal rate to perform these periodic flushes, based on the rate of flushing and how fast **redo** information is generated. First introduced in MySQL 5.1, in the InnoDB Plugin.

See Also [buffer pool](#), [checkpoint](#), [data files](#), [flush](#), [InnoDB](#), [page](#), [redo log](#).

adaptive hash index

An optimization for InnoDB tables that can speed up lookups using `=` and `IN` operators, by constructing a **hash index** in memory. MySQL monitors index searches for InnoDB tables, and if queries could benefit from a hash index, it builds one automatically for index **pages** that are frequently accessed. In a sense, the adaptive hash index configures MySQL at runtime to take advantage of ample main memory, coming closer to the architecture

of main-memory databases. This feature is controlled by the [innodb_adaptive_hash_index](#) configuration option. Because this feature benefits some workloads and not others, and the memory used for the hash index is reserved in the **buffer pool**, typically you should benchmark with this feature both enabled and disabled.

The hash index is always built based on an existing InnoDB **secondary index**, which is organized as a **B-tree** structure. MySQL can build a hash index on a prefix of any length of the key defined for the B-tree, depending on the pattern of searches against the index. A hash index can be partial; the whole B-tree index does not need to be cached in the buffer pool.

In MySQL 5.6 and higher, another way to take advantage of fast single-value lookups with InnoDB tables is to use the **memcached** interface to InnoDB. See [InnoDB Integration with memcached](#) for details. See Also [B-tree](#), [buffer pool](#), [hash index](#), [memcached](#), [page](#), [secondary index](#).

AHI

Acronym for **adaptive hash index**.
See Also [adaptive hash index](#).

AIO

Acronym for **asynchronous I/O**. You might see this acronym in InnoDB messages or keywords.
See Also [asynchronous I/O](#).

Antelope

The code name for the original InnoDB **file format**. It supports the **redundant** and **compact** row formats, but not the newer **dynamic** and **compressed** row formats available in the **Barracuda** file format.

If your application could benefit from InnoDB table **compression**, or uses BLOBs or large text columns that could benefit from the dynamic row format, you might switch some tables to Barracuda format. You select the file format to use by setting the [innodb_file_format](#) option before creating the table.

See Also [Barracuda](#), [compact row format](#), [compressed row format](#), [dynamic row format](#), [file format](#), [innodb_file_format](#), [redundant row format](#).

application programming interface (API)

A set of functions or procedures. An API provides a stable set of names and types for functions, procedures, parameters, and return values.

apply

When a backup produced by the **MySQL Enterprise Backup** product does not include the most recent changes that occurred while the backup was underway, the process of updating the backup files to include those changes is known as the **apply** step. It is specified by the [apply-log](#) option of the [mysqlbackup](#) command.

Before the changes are applied, we refer to the files as a **raw backup**. After the changes are applied, we refer to the files as a **prepared backup**. The changes are recorded in the [ibbackup_logfile](#) file; once the apply step is finished, this file is no longer necessary.

See Also [hot backup](#), [ibbackup_logfile](#), [MySQL Enterprise Backup](#), [prepared backup](#), [raw backup](#).

asynchronous I/O

A type of I/O operation that allows other processing to proceed before the I/O is completed. Also known as **non-blocking I/O** and abbreviated as **AIO**. InnoDB uses this type of I/O for certain operations that can run in parallel without affecting the reliability of the database, such as reading pages into the **buffer pool** that have not actually been requested, but might be needed soon.

Historically, InnoDB has used asynchronous I/O on Windows systems only. Starting with the InnoDB Plugin 1.1, InnoDB uses asynchronous I/O on Linux systems. This change introduces a dependency on [libaio](#). On other Unix-like systems, InnoDB uses synchronous I/O only.

See Also [buffer pool](#), [non-blocking I/O](#).

atomic

In the SQL context, **transactions** are units of work that either succeed entirely (when **committed**) or have no effect at all (when **rolled back**). The indivisible ("atomic") property of transactions is the "A" in the acronym **ACID**. See Also [ACID](#), [commit](#), [rollback](#), [transaction](#).

atomic instruction

Special instructions provided by the CPU, to ensure that critical low-level operations cannot be interrupted.

auto-increment

A property of a table column (specified by the [AUTO_INCREMENT](#) keyword) that automatically adds an ascending sequence of values in the column. InnoDB supports auto-increment only for **primary key** columns.

It saves work for the developer, not to have to produce new unique values when inserting new rows. It provides useful information for the query optimizer, because the column is known to be not null and with unique values. The values from such a column can be used as lookup keys in various contexts, and because they are auto-generated there is no reason to ever change them; for this reason, primary key columns are often specified as auto-incrementing.

Auto-increment columns can be problematic with statement-based replication, because replaying the statements on a slave might not produce the same set of column values as on the master, due to timing issues. When you have an auto-incrementing primary key, you can use statement-based replication only with the setting [innodb_autoinc_lock_mode=1](#). If you have [innodb_autoinc_lock_mode=2](#), which allows higher concurrency for insert operations, use **row-based replication** rather than **statement-based replication**. The setting [innodb_autoinc_lock_mode=0](#) is the previous (traditional) default setting and should not be used except for compatibility purposes.

See Also [auto-increment locking](#), [innodb_autoinc_lock_mode](#), [primary key](#), [row-based replication](#), [statement-based replication](#).

auto-increment locking

The convenience of an **auto-increment** primary key involves some tradeoff with concurrency. In the simplest case, if one transaction is inserting values into the table, any other transactions must wait to do their own inserts into that table, so that rows inserted by the first transaction receive consecutive primary key values. InnoDB includes optimizations, and the [innodb_autoinc_lock_mode](#) option, so that you can choose how to trade off between predictable sequences of auto-increment values and maximum **concurrency** for insert operations.

See Also [auto-increment](#), [concurrency](#), [innodb_autoinc_lock_mode](#).

autocommit

A setting that causes a **commit** operation after each **SQL** statement. This mode is not recommended for working with InnoDB tables with **transactions** that span several statements. It can help performance for **read-only transactions** on InnoDB tables, where it minimizes overhead from **locking** and generation of **undo** data, especially in MySQL 5.6.4 and up. It is also appropriate for working with MyISAM tables, where transactions are not applicable.

See Also [commit](#), [locking](#), [read-only transaction](#), [SQL](#), [transaction](#), [undo](#).

availability

The ability to cope with, and if necessary recover from, failures on the host, including failures of MySQL, the operating system, or the hardware and maintenance activity that may otherwise cause downtime. Often paired with **scalability** as critical aspects of a large-scale deployment.

See Also [scalability](#).

B

B-tree

A tree data structure that is popular for use in database indexes. The structure is kept sorted at all times, enabling fast lookup for exact matches (equals operator) and ranges (for example, greater than, less than, and [BETWEEN](#) operators). This type of index is available for most storage engines, such as InnoDB and MyISAM.

Because B-tree nodes can have many children, a B-tree is not the same as a binary tree, which is limited to 2 children per node.

Contrast with **hash index**, which is only available in the MEMORY storage engine. The MEMORY storage engine can also use B-tree indexes, and you should choose B-tree indexes for MEMORY tables if some queries use range operators.

See Also [hash index](#).

backticks

Identifiers within MySQL SQL statements must be quoted using the backtick character (```) if they contain special characters or reserved words. For example, to refer to a table named `FOO#BAR` or a column named `SELECT`, you would specify the identifiers as ``FOO#BAR`` and ``SELECT``. Since the backticks provide an extra level of safety, they are used extensively in program-generated SQL statements, where the identifier names might not be known in advance.

Many other database systems use double quotation marks (`"`) around such special names. For portability, you can enable `ANSI_QUOTES` mode in MySQL and use double quotation marks instead of backticks to qualify identifier names.

See Also [SQL](#).

backup

The process of copying some or all table data and metadata from a MySQL instance, for safekeeping. Can also refer to the set of copied files. This is a crucial task for DBAs. The reverse of this process is the **restore** operation.

With MySQL, **physical backups** are performed by the **MySQL Enterprise Backup** product, and **logical backups** are performed by the `mysqldump` command. These techniques have different characteristics in terms of size and representation of the backup data, and speed (especially speed of the restore operation).

Backups are further classified as **hot**, **warm**, or **cold** depending on how much they interfere with normal database operation. (Hot backups have the least interference, cold backups the most.)

See Also [cold backup](#), [hot backup](#), [logical backup](#), [MySQL Enterprise Backup](#), [mysqldump](#), [physical backup](#), [warm backup](#).

Barracuda

The code name for an InnoDB **file format** that supports compression for table data. This file format was first introduced in the InnoDB Plugin. It supports the **compressed** row format that enables InnoDB table compression, and the **dynamic** row format that improves the storage layout for BLOB and large text columns. You can select it through the `innodb_file_format` option.

Because the InnoDB **system tablespace** is stored in the original **Antelope** file format, to use the Barracuda file format you must also enable the **file-per-table** setting, which puts newly created tables in their own tablespaces separate from the system tablespace.

The **MySQL Enterprise Backup** product version 3.5 and above supports backing up tablespaces that use the Barracuda file format.

See Also [Antelope](#), [compact row format](#), [compressed row format](#), [dynamic row format](#), [file format](#), [file-per-table](#), [innodb_file_format](#), [MySQL Enterprise Backup](#), [row format](#), [system tablespace](#).

beta

An early stage in the life of a software product, when it is available only for evaluation, typically without a definite release number or a number less than 1. InnoDB does not use the beta designation, preferring an **early adopter** phase that can extend over several point releases, leading to a **GA** release.

See Also [early adopter](#), [GA](#).

binary log

A file containing a record of all statements that attempt to change table data. These statements can be replayed to bring slave servers up to date in a **replication** scenario, or to bring a database up to date after restoring table

data from a backup. The binary logging feature can be turned on and off, although Oracle recommends always enabling it if you use replication or perform backups.

You can examine the contents of the binary log, or replay those statements during replication or recovery, by using the `mysqlbinlog` command. For full information about the binary log, see [The Binary Log](#). For MySQL configuration options related to the binary log, see [Binary Log Options and Variables](#).

For the **MySQL Enterprise Backup** product, the file name of the binary log and the current position within the file are important details. To record this information for the master server when taking a backup in a replication context, you can specify the `--slave-info` option.

Prior to MySQL 5.0, a similar capability was available, known as the update log. In MySQL 5.0 and higher, the binary log replaces the update log.

See Also [binlog](#), [MySQL Enterprise Backup](#), [replication](#).

binlog

An informal name for the **binary log** file. For example, you might see this abbreviation used in e-mail messages or forum discussions.

See Also [binary log](#).

blind query expansion

A special mode of **full-text search** enabled by the `WITH QUERY EXPANSION` clause. It performs the search twice, where the search phrase for the second search is the original search phrase concatenated with the few most highly relevant documents from the first search. This technique is mainly applicable for short search phrases, perhaps only a single word. It can uncover relevant matches where the precise search term does not occur in the document.

See Also [full-text search](#).

bottleneck

A portion of a system that is constrained in size or capacity, that has the effect of limiting overall throughput. For example, a memory area might be smaller than necessary; access to a single required resource might prevent multiple CPU cores from running simultaneously; or waiting for disk I/O to complete might prevent the CPU from running at full capacity. Removing bottlenecks tends to improve **concurrency**. For example, the ability to have multiple InnoDB **buffer pool** instances reduces contention when multiple sessions read from and write to the buffer pool simultaneously.

See Also [buffer pool](#), [concurrency](#).

bounce

A **shutdown** operation immediately followed by a restart. Ideally with a relatively short **warmup** period so that performance and throughput quickly return to a high level.

See Also [shutdown](#).

buddy allocator

A mechanism for managing different-sized **pages** in the InnoDB **buffer pool**.

See Also [buffer pool](#), [page](#), [page size](#).

buffer

A memory or disk area used for temporary storage. Data is buffered in memory so that it can be written to disk efficiently, with a few large I/O operations rather than many small ones. Data is buffered on disk for greater reliability, so that it can be recovered even when a **crash** or other failure occurs at the worst possible time. The main types of buffers used by InnoDB are the **buffer pool**, the **doublewrite buffer**, and the **insert buffer**.

See Also [buffer pool](#), [crash](#), [doublewrite buffer](#), [insert buffer](#).

buffer pool

The memory area that holds cached InnoDB data for both tables and indexes. For efficiency of high-volume read operations, the buffer pool is divided into **pages** that can potentially hold multiple rows. For efficiency of cache

management, the buffer pool is implemented as a linked list of pages; data that is rarely used is aged out of the cache, using a variation of the **LRU** algorithm. On systems with large memory, you can improve concurrency by dividing the buffer pool into multiple **buffer pool instances**.

Several **InnoDB** status variables, `information_schema` tables, and `performance_schema` tables help to monitor the internal workings of the buffer pool. Starting in MySQL 5.6, you can also dump and restore the contents of the buffer pool, either automatically during shutdown and restart, or manually at any time, through a set of **InnoDB** configuration variables such as `innodb_buffer_pool_dump_at_shutdown` and `innodb_buffer_pool_load_at_startup`.

See Also [buffer pool instance](#), [LRU](#), [page](#), [warm up](#).

buffer pool instance

Any of the multiple regions into which the **buffer pool** can be divided, controlled by the `innodb_buffer_pool_instances` configuration option. The total memory size specified by the `innodb_buffer_pool_size` is divided among all the instances. Typically, multiple buffer pool instances are appropriate for systems devoting multiple gigabytes to the **InnoDB** buffer pool, with each instance 1 gigabyte or larger. On systems loading or looking up large amounts of data in the buffer pool from many concurrent sessions, having multiple instances reduces the contention for exclusive access to the data structures that manage the buffer pool.

See Also [buffer pool](#).

built-in

The built-in **InnoDB** storage engine within MySQL is the original form of distribution for the storage engine. Contrast with the **InnoDB Plugin**. Starting with MySQL 5.5, the **InnoDB Plugin** is merged back into the MySQL code base as the built-in **InnoDB** storage engine (known as **InnoDB 1.1**).

This distinction is important mainly in MySQL 5.1, where a feature or bug fix might apply to the **InnoDB Plugin** but not the built-in **InnoDB**, or vice versa.

See Also [InnoDB](#), [plugin](#).

business rules

The relationships and sequences of actions that form the basis of business software, used to run a commercial company. Sometimes these rules are dictated by law, other times by company policy. Careful planning ensures that the relationships encoded and enforced by the database, and the actions performed through application logic, accurately reflect the real policies of the company and can handle real-life situations.

For example, an employee leaving a company might trigger a sequence of actions from the human resources department. The human resources database might also need the flexibility to represent data about a person who has been hired, but not yet started work. Closing an account at an online service might result in data being removed from a database, or the data might be moved or flagged so that it could be recovered if the account is re-opened. A company might establish policies regarding salary maximums, minimums, and adjustments, in addition to basic sanity checks such as the salary not being a negative number. A retail database might not allow a purchase with the same serial number to be returned more than once, or might not allow credit card purchases above a certain value, while a database used to detect fraud might allow these kinds of things.

See Also [relational](#).

C

.cfg file

A metadata file used with the **InnoDB transportable tablespace** feature. It is produced by the command `FLUSH TABLES ... FOR EXPORT`, puts one or more tables in a consistent state that can be copied to another server. The `.cfg` file is copied along with the corresponding `.ibd` file, and used to adjust the internal values of the `.ibd` file, such as the **space ID**, during the `ALTER TABLE ... IMPORT TABLESPACE` step.

See Also [.ibd file](#), [space ID](#), [transportable tablespace](#).

cache

The general term for any memory area that stores copies of data for frequent or high-speed retrieval. In InnoDB, the primary kind of cache structure is the **buffer pool**.

See Also [buffer](#), [buffer pool](#).

cardinality

The number of different values in a table **column**. When queries refer to columns that have an associated **index**, the cardinality of each column influences which access method is most efficient. For example, for a column with a **unique constraint**, the number of different values is equal to the number of rows in the table. If a table has a million rows but only 10 different values for a particular column, each value occurs (on average) 100,000 times. A query such as `SELECT c1 FROM t1 WHERE c1 = 50;` thus might return 1 row or a huge number of rows, and the database server might process the query differently depending on the cardinality of `c1`.

If the values in a column have a very uneven distribution, the cardinality might not be a good way to determine the best query plan. For example, `SELECT c1 FROM t1 WHERE c1 = x;` might return 1 row when `x=50` and a million rows when `x=30`. In such a case, you might need to use **index hints** to pass along advice about which lookup method is more efficient for a particular query.

Cardinality can also apply to the number of distinct values present in multiple columns, as in a **composite index**.

For InnoDB, the process of estimating cardinality for indexes is influenced by the [innodb_stats_sample_pages](#) and the [innodb_stats_on_metadata](#) configuration options. The estimated values are more stable when the **persistent statistics** feature is enabled (in MySQL 5.6 and higher).

See Also [column](#), [composite index](#), [index](#), [index hint](#), [persistent statistics](#), [random dive](#), [selectivity](#), [unique constraint](#).

change buffer

A special data structure that records changes to **pages in secondary indexes**. These values could result from SQL `INSERT`, `UPDATE`, or `DELETE` statements (**DML**). The set of features involving the change buffer is known collectively as **change buffering**, consisting of **insert buffering**, **delete buffering**, and **purge buffering**.

Changes are only recorded in the change buffer when the relevant page from the secondary index is not in the **buffer pool**. When the relevant index page is brought into the buffer pool while associated changes are still in the change buffer, the changes for that page are applied in the buffer pool (**merged**) using the data from the change buffer. Periodically, the **purge** operation that runs during times when the system is mostly idle, or during a slow shutdown, writes the new index pages to disk. The purge operation can write the disk blocks for a series of index values more efficiently than if each value were written to disk immediately.

Physically, the change buffer is part of the **system tablespace**, so that the index changes remain buffered across database restarts. The changes are only applied (**merged**) when the pages are brought into the buffer pool due to some other read operation.

The kinds and amount of data stored in the change buffer are governed by the [innodb_change_buffering](#) and [innodb_change_buffer_max_size](#) configuration options. To see information about the current data in the change buffer, issue the `SHOW ENGINE INNODB STATUS` command.

Formerly known as the **insert buffer**.

See Also [buffer pool](#), [change buffering](#), [delete buffering](#), [DML](#), [insert buffer](#), [insert buffering](#), [merge](#), [page](#), [purge](#), [purge buffering](#), [secondary index](#), [system tablespace](#).

change buffering

The general term for the features involving the **change buffer**, consisting of **insert buffering**, **delete buffering**, and **purge buffering**. Index changes resulting from SQL statements, which could normally involve random I/O operations, are held back and performed periodically by a background **thread**. This sequence of operations can write the disk blocks for a series of index values more efficiently than if each value were written to disk

immediately. Controlled by the [innodb_change_buffering](#) and [innodb_change_buffer_max_size](#) configuration options.

See Also [change buffer](#), [delete buffering](#), [insert buffering](#), [purge buffering](#).

checkpoint

As changes are made to data pages that are cached in the **buffer pool**, those changes are written to the **data files** sometime later, a process known as **flushing**. The checkpoint is a record of the latest changes (represented by an **LSN** value) that have been successfully written to the data files.

See Also [buffer pool](#), [data files](#), [flush](#), [fuzzy checkpointing](#), [LSN](#).

checksum

In **InnoDB**, a validation mechanism to detect corruption when a **page** in a **tablespace** is read from disk into the InnoDB **buffer pool**. This feature is turned on and off by the [innodb_checksums](#) configuration option. In MySQL 5.6, you can enable a faster checksum algorithm by also specifying the configuration option [innodb_checksum_algorithm=crc32](#).

The [innochecksum](#) command helps to diagnose corruption problems by testing the checksum values for a specified **tablespace** file while the MySQL server is shut down.

MySQL also uses checksums for replication purposes. For details, see the configuration options [binlog_checksum](#), [master_verify_checksum](#), and [slave_sql_verify_checksum](#).

See Also [buffer pool](#), [page](#), [tablespace](#).

child table

In a **foreign key** relationship, a child table is one whose rows refer (or point) to rows in another table with an identical value for a specific column. This is the table that contains the `FOREIGN KEY ... REFERENCES` clause and optionally `ON UPDATE` and `ON DELETE` clauses. The corresponding row in the **parent table** must exist before the row can be created in the child table. The values in the child table can prevent delete or update operations on the parent table, or can cause automatic deletion or updates in the child table, based on the `ON CASCADE` option used when creating the foreign key.

See Also [foreign key](#), [parent table](#).

clean page

A **page** in the InnoDB **buffer pool** where all changes made in memory have also been written (**flushed**) to the **data files**. The opposite of a **dirty page**.

See Also [buffer pool](#), [data files](#), [dirty page](#), [flush](#), [page](#).

clean shutdown

A **shutdown** that completes without errors and applies all changes to InnoDB tables before finishing, as opposed to a **crash** or a **fast shutdown**. Synonym for **slow shutdown**.

See Also [crash](#), [fast shutdown](#), [shutdown](#), [slow shutdown](#).

client

A type of program that sends requests to a server, and interprets or processes the results. The client software might run only some of the time (such as a mail or chat program), and might run interactively (such as the [mysql](#) command processor).

See Also [mysql](#), [server](#).

clustered index

The InnoDB term for a **primary key** index. InnoDB table storage is organized based on the values of the primary key columns, to speed up queries and sorts involving the primary key columns. For best performance, choose the primary key columns carefully based on the most performance-critical queries. Because modifying the columns of the clustered index is an expensive operation, choose primary columns that are rarely or never updated.

In the Oracle Database product, this type of table is known as an **index-organized table**.

See Also [index](#), [primary key](#), [secondary index](#).

cold backup

A **backup** taken while the database is shut down. For busy applications and web sites, this might not be practical, and you might prefer a **warm backup** or a **hot backup**.

See Also [backup](#), [hot backup](#), [warm backup](#).

column

A data item within a **row**, whose storage and semantics are defined by a data type. Each **table** and **index** is largely defined by the set of columns it contains.

Each column has a **cardinality** value. A column can be the **primary key** for its table, or part of the primary key. A column can be subject to a **unique constraint**, a **NOT NULL constraint**, or both. Values in different columns, even across different tables, can be linked by a **foreign key** relationship.

In discussions of MySQL internal operations, sometimes **field** is used as a synonym.

See Also [cardinality](#), [foreign key](#), [index](#), [primary key](#), [row](#), [SQL](#), [table](#), [unique constraint](#).

column index

An **index** on a single column.

See Also [composite index](#), [index](#).

column prefix

When an index is created with a length specification, such as `CREATE INDEX idx ON t1 (c1(N))`, only the first N characters of the column value are stored in the index. Keeping the index prefix small makes the index compact, and the memory and disk I/O savings help performance. (Although making the index prefix too small can hinder query optimization by making rows with different values appear to the query optimizer to be duplicates.)

For columns containing binary values or long text strings, where sorting is not a major consideration and storing the entire value in the index would waste space, the index automatically uses the first N (typically 768) characters of the value to do lookups and sorts.

See Also [index](#).

commit

A **SQL** statement that ends a **transaction**, making permanent any changes made by the transaction. It is the opposite of **rollback**, which undoes any changes made in the transaction.

InnoDB uses an **optimistic** mechanism for commits, so that changes can be written to the data files before the commit actually occurs. This technique makes the commit itself faster, with the tradeoff that more work is required in case of a rollback.

By default, MySQL uses the **autocommit** setting, which automatically issues a commit following each SQL statement.

See Also [autocommit](#), [optimistic](#), [rollback](#), [SQL](#), [transaction](#).

compact row format

The default **InnoDB row format** since MySQL 5.0.3. Available for tables that use the **Antelope file format**. It has a more compact representation for nulls and variable-length fields than the prior default (**redundant row format**).

Because of the **B-tree** indexes that make row lookups so fast in InnoDB, there is little if any performance benefit to keeping all rows the same size.

For additional information about **InnoDB COMPACT** row format, see [Section 5.2, “COMPACT and REDUNDANT Row Formats”](#).

See Also [Antelope](#), [file format](#), [redundant row format](#), [row format](#).

composite index

An **index** that includes multiple columns.

See Also [index](#), [index prefix](#).

compressed backup

The compression feature of the **MySQL Enterprise Backup** product makes a compressed copy of each tablespace, changing the extension from `.ibd` to `.ibz`. Compressing the backup data allows you to keep more backups on hand, and reduces the time to transfer backups to a different server. The data is uncompressed during the restore operation. When a compressed backup operation processes a table that is already compressed, it skips the compression step for that table, because compressing again would result in little or no space savings.

A set of files produced by the **MySQL Enterprise Backup** product, where each **tablespace** is compressed. The compressed files are renamed with a `.ibz` file extension.

Applying **compression** right at the start of the backup process helps to avoid storage overhead during the compression process, and to avoid network overhead when transferring the backup files to another server. The process of **applying the binary log** takes longer, and requires uncompressing the backup files.

See Also [apply](#), [binary log](#), [compression](#), [hot backup](#), [MySQL Enterprise Backup](#), [tablespace](#).

compressed row format

A **row format** that enables data and index **compression** for **InnoDB** tables. It was introduced in the **InnoDB** Plugin, available as part of the **Barracuda** file format. Large fields are stored away from the page that holds the rest of the row data, as in **dynamic row format**. Both index pages and the large fields are compressed, yielding memory and disk savings. Depending on the structure of the data, the decrease in memory and disk usage might or might not outweigh the performance overhead of uncompressing the data as it is used. See [Chapter 3, InnoDB Data Compression](#) for usage details.

For additional information about **InnoDB COMPRESSED** row format, see [Section 5.3, “DYNAMIC Row Format”](#).

See Also [Barracuda](#), [compression](#), [dynamic row format](#), [row format](#).

compression

A feature with wide-ranging benefits from using less disk space, performing less I/O, and using less memory for caching. **InnoDB** table and index data can be kept in a compressed format during database operation.

The data is uncompressed when needed for queries, and re-compressed when changes are made by **DML** operations. After you enable compression for a table, this processing is transparent to users and application developers. DBAs can consult **information_schema** tables to monitor how efficiently the compression parameters work for the MySQL instance and for particular compressed tables.

When **InnoDB** table data is compressed, the compression applies to the **table** itself, any associated **index** data, and the pages loaded into the **buffer pool**. Compression does not apply to pages in the **undo buffer**.

The table compression feature requires using MySQL 5.5 or higher, or the **InnoDB** Plugin in MySQL 5.1 or earlier, and creating the table using the **Barracuda** file format and **compressed row format**, with the **innodb_file_per_table** setting turned on. The compression for each table is influenced by the **KEY_BLOCK_SIZE** clause of the **CREATE TABLE** and **ALTER TABLE** statements. In MySQL 5.6 and higher, compression is also affected by the server-wide configuration options **innodb_compression_failure_threshold_pct**, **innodb_compression_level**, and **innodb_compression_pad_pct_max**. See [Chapter 3, InnoDB Data Compression](#) for usage details.

Another type of compression is the **compressed backup** feature of the **MySQL Enterprise Backup** product. See Also [Barracuda](#), [buffer pool](#), [compressed row format](#), [DML](#), [hot backup](#), [index](#), [INFORMATION_SCHEMA](#), [innodb_file_per_table](#), [plugin](#), [table](#), [undo buffer](#).

compression failure

Not actually an error, rather an expensive operation that can occur when using **compression** in combination with **DML** operations. It occurs when: updates to a compressed **page** overflow the area on the page

reserved for recording modifications; the page is compressed again, with all changes applied to the table data; the re-compressed data does not fit on the original page, requiring MySQL to split the data into two new pages and compress each one separately. To check the frequency of this condition, query the table `INFORMATION_SCHEMA.INNODB_CMP` and check how much the value of the `COMPRESS_OPS` column exceeds the value of the `COMPRESS_OPS_OK` column. Ideally, compression failures do not occur often; when they do, you can adjust the configuration options `innodb_compression_level`, `innodb_compression_failure_threshold_pct`, and `innodb_compression_pad_pct_max`. See Also [compression](#), [DML](#), [page](#).

concatenated index

See [composite index](#).

concurrency

The ability of multiple operations (in database terminology, **transactions**) to run simultaneously, without interfering with each other. Concurrency is also involved with performance, because ideally the protection for multiple simultaneous transactions works with a minimum of performance overhead, using efficient mechanisms for **locking**.

See Also [ACID](#), [locking](#), [transaction](#).

configuration file

The file that holds the **option** values used by MySQL at startup. Traditionally, on Linux and UNIX this file is named `my.cnf`, and on Windows it is named `my.ini`. You can set a number of options related to InnoDB under the `[mysqld]` section of the file.

Typically, this file is searched for in the locations `/etc/my.cnf` `/etc/mysql/my.cnf` `/usr/local/mysql/etc/my.cnf` and `~/my.cnf`. See [Using Option Files](#) for details about the search path for this file.

When you use the **MySQL Enterprise Backup** product, you typically use two configuration files: one that specifies where the data comes from and how it is structured (which could be the original configuration file for your real server), and a stripped-down one containing only a small set of options that specify where the backup data goes and how it is structured. The configuration files used with the **MySQL Enterprise Backup** product must contain certain options that are typically left out of regular configuration files, so you might need to add some options to your existing configuration file for use with **MySQL Enterprise Backup**.

See Also [my.cnf](#), [option file](#).

consistent read

A read operation that uses snapshot information to present query results based on a point in time, regardless of changes performed by other transactions running at the same time. If queried data has been changed by another transaction, the original data is reconstructed based on the contents of the **undo log**. This technique avoids some of the **locking** issues that can reduce **concurrency** by forcing transactions to wait for other transactions to finish.

With the **repeatable read** isolation level, the snapshot is based on the time when the first read operation is performed. With the **read committed** isolation level, the snapshot is reset to the time of each consistent read operation.

Consistent read is the default mode in which InnoDB processes `SELECT` statements in **READ COMMITTED** and **REPEATABLE READ** isolation levels. Because a consistent read does not set any locks on the tables it accesses, other sessions are free to modify those tables while a consistent read is being performed on the table.

For technical details about the applicable isolation levels, see [Consistent Nonlocking Reads](#).

See Also [ACID](#), [concurrency](#), [isolation level](#), [locking](#), [MVCC](#), [READ COMMITTED](#), [READ UNCOMMITTED](#), [REPEATABLE READ](#), [SERIALIZABLE](#), [transaction](#), [undo log](#).

constraint

An automatic test that can block database changes to prevent data from becoming inconsistent. (In computer science terms, a kind of assertion related to an invariant condition.) Constraints are a crucial component of

the **ACID** philosophy, to maintain data consistency. Constraints supported by MySQL include **FOREIGN KEY constraints** and **unique constraints**.

See Also [ACID](#), [foreign key](#), [relational](#), [unique constraint](#).

counter

A value that is incremented by a particular kind of [InnoDB](#) operation. Useful for measuring how busy a server is, troubleshooting the sources of performance issues, and testing whether changes (for example, to configuration settings or indexes used by queries) have the desired low-level effects. Different kinds of counters are available through **performance_schema** tables and **information_schema** tables, particularly [information_schema.innodb_metrics](#).

See Also [INFORMATION_SCHEMA](#), [metrics counter](#), [Performance Schema](#).

covering index

An **index** that includes all the columns retrieved by a query. Instead of using the index values as pointers to find the full table rows, the query returns values from the index structure, saving disk I/O. InnoDB can apply this optimization technique to more indexes than MyISAM can, because InnoDB **secondary indexes** also include the primary key columns. InnoDB cannot apply this technique for queries against tables modified by a transaction, until that transaction ends.

Any **column index** or **composite index** could act as a covering index, given the right query. Design your indexes and queries to take advantage of this optimization technique wherever possible.

See Also [column index](#), [composite index](#), [index](#), [secondary index](#).

crash

MySQL uses the term "crash" to refer generally to any unexpected [shutdown](#) operation where the server cannot do its normal cleanup. For example, a crash could happen due to a hardware fault on the database server machine or storage device; a power failure; a potential data mismatch that causes the MySQL server to halt; a **fast shutdown** initiated by the DBA; or many other reasons. The robust, automatic **crash recovery** for **InnoDB** tables ensures that data is made consistent when the server is restarted, without any extra work for the DBA.

See Also [crash recovery](#), [fast shutdown](#), [InnoDB](#), [redo log](#), [shutdown](#).

crash recovery

The cleanup activities that occur when MySQL is started again after a **crash**. For **InnoDB** tables, changes from incomplete transactions are replayed using data from the **redo log**. Changes that were **committed** before the crash, but not yet written into the [data files](#), are reconstructed from the **doublewrite buffer**. When the database is shut down normally, this type of activity is performed during shutdown by the **purge** operation.

During normal operation, committed data can be stored in the **change buffer** for a period of time before being written to the data files. There is always a tradeoff between keeping the data files up-to-date, which introduces performance overhead during normal operation, and buffering the data, which can make shutdown and crash recovery take longer.

See Also [change buffer](#), [commit](#), [crash](#), [data files](#), [doublewrite buffer](#), [InnoDB](#), [purge](#), [redo log](#).

CRUD

Acronym for "create, read, update, delete", a common sequence of operations in database applications. Often denotes a class of applications with relatively simple database usage (basic **DDL**, **DML** and **query** statements in **SQL**) that can be implemented quickly in any language.

See Also [DDL](#), [DML](#), [query](#), [SQL](#).

cursor

An internal data structure that is used to represent the result set of a **query**, or other operation that performs a search using an SQL [WHERE](#) clause. It works like an iterator in other high-level languages, producing each value from the result set as requested.

Although usually SQL handles the processing of cursors for you, you might delve into the inner workings when dealing with performance-critical code.

See Also [query](#).

D

data definition language

See [DDL](#).

data dictionary

Metadata that keeps track of InnoDB-related objects such as **tables**, **indexes**, and table **columns**. This metadata is physically located in the InnoDB **system tablespace**. For historical reasons, it overlaps to some degree with information stored in the **.frm files**.

Because the **MySQL Enterprise Backup** product always backs up the system tablespace, all backups include the contents of the data dictionary.

See Also [column](#), [.frm file](#), [hot backup](#), [index](#), [MySQL Enterprise Backup](#), [system tablespace](#), [table](#).

data directory

The directory under which each MySQL **instance** keeps the **data files** for InnoDB and the directories representing individual databases. Controlled by the `datadir` configuration option.

See Also [data files](#), [instance](#).

data files

The files that physically contain the InnoDB **table** and **index** data. There can be a one-to-many relationship between data files and tables, as in the case of the **system tablespace**, which can hold multiple InnoDB tables as well as the **data dictionary**. There can also be a one-to-one relationship between data files and tables, as when the **file-per-table** setting is enabled, causing each newly created table to be stored in a separate **tablespace**.

See Also [data dictionary](#), [file-per-table](#), [index](#), [system tablespace](#), [table](#), [tablespace](#).

data manipulation language

See [DML](#).

data warehouse

A database system or application that primarily runs large **queries**. The read-only or read-mostly data might be organized in **denormalized** form for query efficiency. Can benefit from the optimizations for **read-only transactions** in MySQL 5.6 and higher. Contrast with **OLTP**.

See Also [denormalized](#), [OLTP](#), [query](#), [read-only transaction](#).

database

Within the MySQL **data directory**, each database is represented by a separate directory. The InnoDB **system tablespace**, which can hold table data from multiple databases within a MySQL **instance**, is kept in its **data files** that reside outside the individual database directories. When **file-per-table** mode is enabled, the **.ibd files** representing individual InnoDB tables are stored inside the database directories.

For long-time MySQL users, a database is a familiar notion. Users coming from an Oracle Database background will find that the MySQL meaning of a database is closer to what Oracle Database calls a **schema**.

See Also [data files](#), [file-per-table](#), [.ibd file](#), [instance](#), [schema](#), [system tablespace](#).

DCL

Data control language, a set of **SQL** statements for managing privileges. In MySQL, consists of the [GRANT](#) and [REVOKE](#) statements. Contrast with **DDL** and **DML**.

See Also [DDL](#), [DML](#), [SQL](#).

DDL

Data definition language, a set of **SQL** statements for manipulating the database itself rather than individual table rows. Includes all forms of the [CREATE](#), [ALTER](#), and [DROP](#) statements. Also includes the [TRUNCATE](#) statement,

because it works differently than a `DELETE FROM table_name` statement, even though the ultimate effect is similar.

DDL statements automatically **commit** the current **transaction**; they cannot be **rolled back**.

The InnoDB-related aspects of DDL include speed improvements for `CREATE INDEX` and `DROP INDEX` statements, and the way the **file-per-table** setting affects the behavior of the `TRUNCATE TABLE` statement.

Contrast with **DML** and **DCL**.

See Also [commit](#), [DCL](#), [DML](#), [file-per-table](#), [rollback](#), [SQL](#), [transaction](#).

deadlock

A situation where different **transactions** are unable to proceed, because each holds a **lock** that the other needs. Because both transactions are waiting for a resource to become available, neither will ever release the locks it holds.

A deadlock can occur when the transactions lock rows in multiple tables (through statements such as `UPDATE` or `SELECT ... FOR UPDATE`), but in the opposite order. A deadlock can also occur when such statements lock ranges of index records and **gaps**, with each transaction acquiring some locks but not others due to a timing issue.

To reduce the possibility of deadlocks, use transactions rather than `LOCK TABLE` statements; keep transactions that insert or update data small enough that they do not stay open for long periods of time; when different transactions update multiple tables or large ranges of rows, use the same order of operations (such as `SELECT ... FOR UPDATE`) in each transaction; create indexes on the columns used in `SELECT ... FOR UPDATE` and `UPDATE ... WHERE` statements. The possibility of deadlocks is not affected by the **isolation level**, because the isolation level changes the behavior of read operations, while deadlocks occur because of write operations.

If a deadlock does occur, InnoDB detects the condition and **rolls back** one of the transactions (the **victim**). Thus, even if your application logic is perfectly correct, you must still handle the case where a transaction must be retried. To see the last deadlock in an InnoDB user transaction, use the command `SHOW ENGINE INNODB STATUS`. If frequent deadlocks highlight a problem with transaction structure or application error handling, run with the `innodb_print_all_deadlocks` setting enabled to print information about all deadlocks to the `mysqld` error log.

For background information on how deadlocks are automatically detected and handled, see [Deadlock Detection and Rollback](#). For tips on avoiding and recovering from deadlock conditions, see [How to Cope with Deadlocks](#). See Also [concurrency](#), [gap](#), [isolation level](#), [lock](#), [locking](#), [rollback](#), [transaction](#), [victim](#).

deadlock detection

A mechanism that automatically detects when a **deadlock** occurs, and automatically **rolls back** one of the **transactions** involved (the **victim**).

See Also [deadlock](#), [rollback](#), [transaction](#), [victim](#).

delete

When InnoDB processes a `DELETE` statement, the rows are immediately marked for deletion and no longer are returned by queries. The storage is reclaimed sometime later, during the periodic garbage collection known as the **purge** operation, performed by a separate thread. For removing large quantities of data, related operations with their own performance characteristics are **truncate** and **drop**.

See Also [drop](#), [purge](#), [truncate](#).

delete buffering

The technique of storing index changes due to `DELETE` operations in the **insert buffer** rather than writing them immediately, so that the physical writes can be performed to minimize random I/O. (Because delete operations are a two-step process, this operation buffers the write that normally marks an index record for deletion.) It is one of the types of **change buffering**; the others are **insert buffering** and **purge buffering**.

See Also [change buffer](#), [change buffering](#), [insert buffer](#), [insert buffering](#), [purge buffering](#).

denormalized

A data storage strategy that duplicates data across different tables, rather than linking the tables with **foreign keys** and **join** queries. Typically used in **data warehouse** applications, where the data is not updated after loading. In such applications, query performance is more important than making it simple to maintain consistent data during updates. Contrast with **normalized**.

See Also [data warehouse](#), [normalized](#).

descending index

A type of index available with some database systems, where index storage is optimized to process [ORDER BY column DESC](#) clauses. Currently, although MySQL allows the [DESC](#) keyword in the [CREATE TABLE](#) statement, it does not use any special storage layout for the resulting index.

See Also [index](#).

dirty page

A **page** in the InnoDB **buffer pool** that has been updated in memory, where the changes are not yet written (**flushed**) to the [data files](#). The opposite of a **clean page**.

See Also [buffer pool](#), [clean page](#), [data files](#), [flush](#), [page](#).

dirty read

An operation that retrieves unreliable data, data that was updated by another transaction but not yet **committed**. It is only possible with the **isolation level** known as **read uncommitted**.

This kind of operation does not adhere to the **ACID** principle of database design. It is considered very risky, because the data could be **rolled back**, or updated further before being committed; then, the transaction doing the dirty read would be using data that was never confirmed as accurate.

Its polar opposite is **consistent read**, where InnoDB goes to great lengths to ensure that a transaction does not read information updated by another transaction, even if the other transaction commits in the meantime.

See Also [ACID](#), [commit](#), [consistent read](#), [isolation level](#), [READ COMMITTED](#), [READ UNCOMMITTED](#), [rollback](#).

disk-based

A kind of database that primarily organizes data on disk storage (hard drives or equivalent). Data is brought back and forth between disk and memory to be operated upon. It is the opposite of an **in-memory database**. Although InnoDB is disk-based, it also contains features such as **the buffer pool**, multiple buffer pool instances, and the **adaptive hash index** that allow certain kinds of workloads to work primarily from memory.

See Also [adaptive hash index](#), [buffer pool](#), [in-memory database](#).

disk-bound

A type of **workload** where the primary **bottleneck** is CPU operations in memory. Typically involves read-intensive operations where the results can all be cached in the **buffer pool**.

See Also [bottleneck](#), [buffer pool](#), [disk-bound](#), [workload](#).

disk-bound

A type of **workload** where the primary **bottleneck** is disk I/O. (Also known as **I/O-bound**.) Typically involves frequent writes to disk, or random reads of more data than can fit into the **buffer pool**.

See Also [bottleneck](#), [buffer pool](#), [disk-bound](#), [workload](#).

DML

Data manipulation language, a set of **SQL** statements for performing insert, update, and delete operations. The [SELECT](#) statement is sometimes considered as a DML statement, because the [SELECT . . . FOR UPDATE](#) form is subject to the same considerations for **locking** as [INSERT](#), [UPDATE](#), and [DELETE](#).

DML statements for an InnoDB table operate in the context of a **transaction**, so their effects can be **committed** or **rolled back** as a single unit.

Contrast with **DDL** and **DCL**.

See Also [commit](#), [DCL](#), [DDL](#), [locking](#), [rollback](#), [SQL](#), [transaction](#).

document id

In the InnoDB **full-text search** feature, a special column in the table containing the **FULLTEXT index**, to uniquely identify the document associated with each **ilist** value. Its name is `FTS_DOC_ID` (uppercase required). The column itself must be of `BIGINT UNSIGNED NOT NULL` type, with a unique index named `FTS_DOC_ID_INDEX`. Preferably, you define this column when creating the table. If InnoDB must add the column to the table while creating a **FULLTEXT** index, the indexing operation is considerably more expensive.

See Also [full-text search](#), [FULLTEXT index](#), [ilist](#).

doublewrite buffer

InnoDB uses a novel file flush technique called doublewrite. Before writing **pages** to the **data files**, InnoDB first writes them to a contiguous area called the doublewrite buffer. Only after the write and the flush to the doublewrite buffer have completed, does InnoDB write the pages to their proper positions in the data file. If the operating system crashes in the middle of a page write, InnoDB can later find a good copy of the page from the doublewrite buffer during **crash recovery**.

Although data is always written twice, the doublewrite buffer does not require twice as much I/O overhead or twice as many I/O operations. Data is written to the buffer itself as a large sequential chunk, with a single `fsync()` call to the operating system.

To turn off the doublewrite buffer, specify the option `innodb_doublewrite=0`.

See Also [crash recovery](#), [data files](#), [page](#), [purge](#).

drop

A kind of **DDL** operation that removes a schema object, through a statement such as `DROP TABLE` or `DROP INDEX`. It maps internally to an `ALTER TABLE` statement. From an InnoDB perspective, the performance considerations of such operations involve the time that the **data dictionary** is locked to ensure that interrelated objects are all updated, and the time to update memory structures such as the **buffer pool**. For a **table**, the drop operation has somewhat different characteristics than a **truncate** operation (`TRUNCATE TABLE` statement).

See Also [buffer pool](#), [data dictionary](#), [DDL](#), [table](#), [truncate](#).

dynamic row format

A row format introduced in the **InnoDB Plugin**, available as part of the **Barracuda file format**. Because **TEXT** and **BLOB** fields are stored outside of the rest of the page that holds the row data, it is very efficient for rows that include large objects. Since the large fields are typically not accessed to evaluate query conditions, they are not brought into the **buffer pool** as often, resulting in fewer I/O operations and better utilization of cache memory.

For additional information about **InnoDB DYNAMIC** row format, see [Section 5.3, “DYNAMIC Row Format”](#).

See Also [Barracuda](#), [buffer pool](#), [file format](#), [row format](#).

E

early adopter

A stage similar to beta, when a software product is typically evaluated for performance, functionality, and compatibility in a non-mission-critical setting. InnoDB uses the **early adopter** designation rather than **beta**, through a succession of point releases leading up to a **GA** release.

See Also [beta](#), [GA](#).

error log

A type of **log** showing information about MySQL startup and critical runtime errors and **crash** information. For details, see [The Error Log](#).

See Also [crash](#), [log](#).

eviction

The process of removing an item from a cache or other temporary storage area, such as the InnoDB **buffer pool**. Often, but not always, uses the **LRU** algorithm to determine which item to remove. When a **dirty page** is evicted, its contents are **flushed** to disk, and any **dirty neighbor** pages might be flushed also. See Also [buffer pool](#), [dirty page](#), [flush](#), [LRU](#).

exclusive lock

A kind of **lock** that prevents any other **transaction** from locking the same row. Depending on the transaction **isolation level**, this kind of lock might block other transactions from writing to the same row, or might also block other transactions from reading the same row. The default InnoDB isolation level, **REPEATABLE READ**, enables higher **concurrency** by allowing transactions to read rows that have exclusive locks, a technique known as **consistent read**.

See Also [concurrency](#), [consistent read](#), [isolation level](#), [lock](#), [REPEATABLE READ](#), [shared lock](#), [transaction](#).

extent

A group of **pages** within a **tablespace** totaling 1 megabyte. With the default **page size** of 16KB, an extent contains 64 pages. In MySQL 5.6, the page size can also be 4KB or 8KB, in which case an extent contains more pages, still adding up to 1MB.

InnoDB features such as **segments**, **read-ahead** requests and the **doublewrite buffer** use I/O operations that read, write, allocate, or free data one extent at a time.

See Also [doublewrite buffer](#), [neighbor page](#), [page](#), [page size](#), [read-ahead](#), [segment](#), [tablespace](#).

F

.frm file

A file containing the metadata, such as the table definition, of a MySQL table.

For backups, you must always keep the full set of `.frm` files along with the backup data to be able to restore tables that are altered or dropped after the backup.

Although each InnoDB table has a `.frm` file, InnoDB maintains its own table metadata in the system tablespace; the `.frm` files are not needed for InnoDB to operate on InnoDB tables.

These files are backed up by the **MySQL Enterprise Backup** product. These files must not be modified by an `ALTER TABLE` operation while the backup is taking place, which is why backups that include non-InnoDB tables perform a `FLUSH TABLES WITH READ LOCK` operation to freeze such activity while backing up the `.frm` files. Restoring a backup can result in `.frm` files being created, changed, or removed to match the state of the database at the time of the backup.

See Also [MySQL Enterprise Backup](#).

Fast Index Creation

A capability first introduced in the InnoDB Plugin, now part of the MySQL server in 5.5 and higher, that speeds up creation of InnoDB **secondary indexes** by avoiding the need to completely rewrite the associated table. The speedup applies to dropping secondary indexes also.

Because index maintenance can add performance overhead to many data transfer operations, consider doing operations such as `ALTER TABLE ... ENGINE=INNODB` or `INSERT INTO ... SELECT * FROM ...` without any secondary indexes in place, and creating the indexes afterward.

In MySQL 5.6, this feature becomes more general: you can read and write to tables while an index is being created, and many more kinds of `ALTER TABLE` operations can be performed without copying the table, without blocking **DML** operations, or both. Thus in MySQL 5.6 and higher, we typically refer to this set of features as **online DDL** rather than Fast Index Creation.

See Also [DML](#), [index](#), [online DDL](#), [secondary index](#).

fast shutdown

The default **shutdown** procedure for InnoDB, based on the configuration setting `innodb_fast_shutdown=1`. To save time, certain **flush** operations are skipped. This type of shutdown is safe during normal usage, because the flush operations are performed during the next startup, using the same mechanism as in **crash recovery**. In cases where the database is being shut down for an upgrade or downgrade, do a **slow shutdown** instead to ensure that all relevant changes are applied to the **data files** during the shutdown. See Also [crash recovery](#), [data files](#), [flush](#), [shutdown](#), [slow shutdown](#).

file format

The format used by InnoDB for each table, typically with the **file-per-table** setting enabled so that each table is stored in a separate **.ibd file**. Currently, the file formats available in InnoDB are known as **Antelope** and **Barracuda**. Each file format supports one or more **row formats**. The row formats available for Barracuda tables, **COMPRESSED** and **DYNAMIC**, enable important new storage features for InnoDB tables. See Also [Antelope](#), [Barracuda](#), [file-per-table](#), [.ibd file](#), [ibdata file](#), [row format](#).

file-per-table

A general name for the setting controlled by the `innodb_file_per_table` option. That is a very important configuration option that affects many aspects of InnoDB file storage, availability of features, and I/O characteristics. In MySQL 5.6.7 and higher, it is enabled by default. Prior to MySQL 5.6.7, it is disabled by default.

For each table created while this setting is in effect, the data is stored in a separate **.ibd file** rather than in the **ibdata files** of the **system tablespace**. When table data is stored in individual files, you have more flexibility to choose nondefault **file formats** and **row formats**, which are required for features such as data **compression**. The `TRUNCATE TABLE` operation is also much faster, and the reclaimed space can be used by the operating system rather than remaining reserved for InnoDB.

The **MySQL Enterprise Backup** product is more flexible for tables that are in their own files. For example, tables can be excluded from a backup, but only if they are in separate files. Thus, this setting is suitable for tables that are backed up less frequently or on a different schedule.

See Also [compressed row format](#), [compression](#), [file format](#), [.ibd file](#), [ibdata file](#), [innodb_file_per_table](#), [row format](#), [system tablespace](#).

fill factor

In an InnoDB **index**, the proportion of a **page** that is taken up by index data before the page is split. The unused space when index data is first divided between pages allows for rows to be updated with longer string values without requiring expensive index maintenance operations. If the fill factor is too low, the index consumes more space than needed, causing extra I/O overhead when reading the index. If the fill factor is too high, any update that increases the length of column values can cause extra I/O overhead for index maintenance. See [Physical Structure of an InnoDB Index](#) for more information.

See Also [index](#), [page](#).

fixed row format

This row format is used by the MyISAM storage engine, not by InnoDB. If you create an InnoDB table with the option `row_format=fixed`, InnoDB translates this option to use the **compact row format** instead, although the `fixed` value might still show up in output such as `SHOW TABLE STATUS` reports.

See Also [compact row format](#), [row format](#).

flush

To write changes to the database files, that had been buffered in a memory area or a temporary disk storage area. The InnoDB storage structures that are periodically flushed include the **redo log**, the **undo log**, and the **buffer pool**.

Flushing can happen because a memory area becomes full and the system needs to free some space, because a **commit** operation means the changes from a transaction can be finalized, or because a **slow shutdown** operation means that all outstanding work should be finalized. When it is not critical to flush all the buffered data

at once, [InnoDB](#) can use a technique called **fuzzy checkpointing** to flush small batches of pages to spread out the I/O overhead.

See Also [buffer pool](#), [commit](#), [fuzzy checkpointing](#), [neighbor page](#), [redo log](#), [slow shutdown](#), [undo log](#).

flush list

An internal InnoDB data structure that tracks **dirty pages** in the **buffer pool**: that is, **pages** that have been changed and need to be written back out to disk. This data structure is updated frequently by InnoDB's internal **mini-transactions**, and so is protected by its own **mutex** to allow concurrent access to the buffer pool.

See Also [buffer pool](#), [dirty page](#), [LRU](#), [mini-transaction](#), [mutex](#), [page](#), [page cleaner](#).

foreign key

A type of pointer relationship, between rows in separate InnoDB tables. The foreign key relationship is defined on one column in both the **parent table** and the **child table**.

In addition to enabling fast lookup of related information, foreign keys help to enforce **referential integrity**, by preventing any of these pointers from becoming invalid as data is inserted, updated, and deleted. This enforcement mechanism is a type of **constraint**. A row that points to another table cannot be inserted if the associated foreign key value does not exist in the other table. If a row is deleted or its foreign key value changed, and rows in another table point to that foreign key value, the foreign key can be set up to prevent the deletion, cause the corresponding column values in the other table to become **null**, or automatically delete the corresponding rows in the other table.

One of the stages in designing a **normalized** database is to identify data that is duplicated, separate that data into a new table, and set up a foreign key relationship so that the multiple tables can be queried like a single table, using a **join** operation.

See Also [child table](#), [FOREIGN KEY constraint](#), [join](#), [normalized](#), [NULL](#), [parent table](#), [referential integrity](#), [relational](#).

FOREIGN KEY constraint

The type of **constraint** that maintains database consistency through a **foreign key** relationship. Like other kinds of constraints, it can prevent data from being inserted or updated if data would become inconsistent; in this case, the inconsistency being prevented is between data in multiple tables. Alternatively, when a **DML** operation is performed, [FOREIGN KEY](#) constraints can cause data in **child rows** to be deleted, changed to different values, or set to **null**, based on the [ON CASCADE](#) option specified when creating the foreign key.

See Also [child table](#), [constraint](#), [DML](#), [foreign key](#), [NULL](#).

FTS

In most contexts, an acronym for **full-text search**. Sometimes in performance discussions, an acronym for **full table scan**.

See Also [full table scan](#), [full-text search](#).

full backup

A **backup** that includes all the **tables** in each MySQL **database**, and all the databases in a MySQL **instance**.

Contrast with **partial backup**.

See Also [backup](#), [database](#), [instance](#), [partial backup](#), [table](#).

full table scan

An operation that requires reading the entire contents of a table, rather than just selected portions using an index. Typically performed either with small lookup tables, or in data warehousing situations with large tables where all available data is aggregated and analyzed. How frequently these operations occur, and the sizes of the tables relative to available memory, have implications for the algorithms used in query optimization and managing the buffer pool.

The purpose of **indexes** is to allow lookups for specific values or ranges of values within a large table, thus avoiding full table scans when practical.

See Also [buffer pool](#), [index](#), [LRU](#).

full-text search

The MySQL feature for finding words, phrases, Boolean combinations of words, and so on within table data, in a faster, more convenient, and more flexible way than using the SQL [LIKE](#) operator or writing your own application-level search algorithm. It uses the SQL function [MATCH\(\)](#) and **FULLTEXT indexes**.

See Also [FULLTEXT index](#).

FULLTEXT index

The special kind of **index** that holds the **search index** in the MySQL **full-text search** mechanism. Represents the words from values of a column, omitting any that are specified as **stopwords**. Originally, only available for [MyISAM](#) tables. Starting in MySQL 5.6.4, it is also available for **InnoDB** tables.

See Also [full-text search](#), [index](#), [InnoDB](#), [search index](#), [stopword](#).

fuzzy checkpointing

A technique that **flushes** small batches of **dirty pages** from the **buffer pool**, rather than flushing all dirty pages at once which would disrupt database processing.

See Also [buffer pool](#), [dirty page](#), [flush](#).

G

GA

"Generally available", the stage when a software product leaves beta and is available for sale, official support, and production use.

See Also [beta](#), [early adopter](#).

gap

A place in an InnoDB **index** data structure where new values could be inserted. When you lock a set of rows with a statement such as [SELECT ... FOR UPDATE](#), InnoDB can create locks that apply to the gaps as well as the actual values in the index. For example, if you select all values greater than 10 for update, a gap lock prevents another transaction from inserting a new value that is greater than 10. The **supremum record** and **infimum record** represent the gaps containing all values greater than or less than all the current index values.

See Also [concurrency](#), [gap lock](#), [index](#), [infimum record](#), [isolation level](#), [supremum record](#).

gap lock

A **lock** on a **gap** between index records, or a lock on the gap before the first or after the last index record. For example, [SELECT c1 FOR UPDATE FROM t WHERE c1 BETWEEN 10 and 20;](#) prevents other transactions from inserting a value of 15 into the column `t.c1`, whether or not there was already any such value in the column, because the gaps between all existing values in the range are locked. Contrast with **record lock** and **next-key lock**.

Gap locks are part of the tradeoff between performance and **concurrency**, and are used in some transaction **isolation levels** and not others.

See Also [gap](#), [infimum record](#), [lock](#), [next-key lock](#), [record lock](#), [supremum record](#).

general log

See [general query log](#).

general query log

A type of **log** used for diagnosis and troubleshooting of SQL statements processed by the MySQL server. Can be stored in a file or in a database table. You must enable this feature through the [general_log](#) configuration option to use it. You can disable it for a specific connection through the [sql_log_off](#) configuration option.

Records a broader range of queries than the **slow query log**. Unlike the **binary log**, which is used for replication, the general query log contains [SELECT](#) statements and does not maintain strict ordering. For more information, see [The General Query Log](#).

See Also [binary log](#), [general query log](#), [log](#).

global_transaction

A type of **transaction** involved in **XA** operations. It consists of several actions that are transactional in themselves, but that all must either complete successfully as a group, or all be rolled back as a group. In essence, this extends **ACID** properties "up a level" so that multiple ACID transactions can be executed in concert as components of a global operation that also has ACID properties. For this type of distributed transaction, you must use the **SERIALIZABLE** isolation level to achieve ACID properties.

See Also [ACID](#), [SERIALIZABLE](#), [transaction](#), [XA](#).

group commit

An **InnoDB** optimization that performs some low-level I/O operations (**log write**) once for a set of **commit** operations, rather than flushing and syncing separately for each commit.

When the binlog is enabled, you typically also set the configuration option `sync_binlog=0`, because group commit for the binary log is only supported if it is set to 0.

See Also [commit](#), [plugin](#), [XA](#).

H

hash index

A type of **index** intended for queries that use equality operators, rather than range operators such as greater-than or **BETWEEN**. It is available for MEMORY tables. Although hash indexes are the default for MEMORY tables for historic reasons, that storage engine also supports **B-tree** indexes, which are often a better choice for general-purpose queries.

MySQL includes a variant of this index type, the **adaptive hash index**, that is constructed automatically for InnoDB tables if needed based on runtime conditions.

See Also [adaptive hash index](#), [B-tree](#), [index](#), [InnoDB](#).

HDD

Acronym for "hard disk drive". Refers to storage media using spinning platters, usually when comparing and contrasting with **SSD**. Its performance characteristics can influence the throughput of a **disk-based** workload.

See Also [disk-based](#), [SSD](#).

heartbeat

A periodic message that is sent to indicate that a system is functioning properly. In a **replication** context, if the **master** stops sending such messages, one of the **slaves** can take its place. Similar techniques can be used between the servers in a cluster environment, to confirm that all of them are operating properly.

See Also [replication](#).

high-water mark

A value representing an upper limit, either a hard limit that should not be exceeded at runtime, or a record of the maximum value that was actually reached. Contrast with **low-water mark**.

See Also [low-water mark](#).

history list

A list of **transactions** with delete-marked records scheduled to be processed by the **InnoDB purge** operation. Recorded in the **undo log**. The length of the history list is reported by the command `SHOW ENGINE INNODB STATUS`. If the history list grows longer than the value of the `innodb_max_purge_lag` configuration option, each **DML** operation is delayed slightly to allow the purge operation to finish **flushing** the deleted records.

Also known as **purge lag**.

See Also [flush](#), [purge](#), [purge lag](#), [rollback segment](#), [transaction](#), [undo log](#).

hot

A condition where a row, table, or internal data structure is accessed so frequently, requiring some form of locking or mutual exclusion, that it results in a performance or scalability issue.

Although "hot" typically indicates an undesirable condition, a **hot backup** is the preferred type of backup. See Also [hot backup](#).

hot backup

A backup taken while the database is running and applications are reading and writing to it. The backup involves more than simply copying data files: it must include any data that was inserted or updated while the backup was in process; it must exclude any data that was deleted while the backup was in process; and it must ignore any changes that were not committed.

The Oracle product that performs hot backups, of InnoDB tables especially but also tables from MyISAM and other storage engines, is known as **MySQL Enterprise Backup**.

The hot backup process consists of two stages. The initial copying of the data files produces a **raw backup**. The **apply** step incorporates any changes to the database that happened while the backup was running. Applying the changes produces a **prepared** backup; these files are ready to be restored whenever necessary.

See Also [apply](#), [MySQL Enterprise Backup](#), [prepared backup](#), [raw backup](#).

I

.ibd file

Each InnoDB **table** created using the **file-per-table** mode goes into its own **tablespace** file, with a **.ibd** extension, inside the **database** directory. This file contains the table data and any **indexes** for the table. File-per-table mode, controlled by the **innodb_file_per_table** option, affects many aspects of InnoDB storage usage and performance, and is enabled by default in MySQL 5.6.7 and higher.

This extension does not apply to the **system tablespace**, which consists of the **ibdata files**.

When a **.ibd** file is included in a compressed backup by the **MySQL Enterprise Backup** product, the compressed equivalent is a **.ibz** file.

If a table is created with the **DATA DIRECTORY =** clause in MySQL 5.6 and higher, the **.ibd** file is located outside the normal database directory, and is pointed to by a **.isl file**.

See Also [database](#), [file-per-table](#), [ibdata file](#), [.ibz file](#), [index](#), [innodb_file_per_table](#), [.isl file](#), [MySQL Enterprise Backup](#), [system tablespace](#), [table](#), [tablespace](#).

.ibz file

When the **MySQL Enterprise Backup** product performs a **compressed backup**, it transforms each **tablespace** file that is created using the **file-per-table** setting from a **.ibd** extension to a **.ibz** extension.

The compression applied during backup is distinct from the **compressed row format** that keeps table data compressed during normal operation. A compressed backup operation skips the compression step for a tablespace that is already in compressed row format, as compressing a second time would slow down the backup but produce little or no space savings.

See Also [compressed backup](#), [compressed row format](#), [file-per-table](#), [.ibd file](#), [MySQL Enterprise Backup](#), [tablespace](#).

.isl file

A file that specifies the location of a **.ibd file** for an InnoDB table created with the **DATA DIRECTORY =** clause in MySQL 5.6 and higher. It functions like a symbolic link, without the platform restrictions of the actual symbolic link mechanism. You can store InnoDB **tablespaces** outside the **database** directory, for example, on an especially large or fast storage device depending on the usage of the table. For details, see [Specifying the Location of a Tablespace](#).

See Also [database](#), [.ibd file](#), [table](#), [tablespace](#).

I/O-bound

See [disk-bound](#).

ib-file set

The set of files managed by InnoDB within a MySQL database: the **system tablespace**, any **file-per-table** tablespaces, and the (typically 2) **redo log** files. Used sometimes in detailed discussions of InnoDB file structures and formats, to avoid ambiguity between the meanings of **database** between different DBMS products, and the non-InnoDB files that may be part of a MySQL database.

See Also [database](#), [file-per-table](#), [redo log](#), [system tablespace](#).

ibbackup_logfile

A supplemental backup file created by the **MySQL Enterprise Backup** product during a **hot backup** operation. It contains information about any data changes that occurred while the backup was running. The initial backup files, including [ibbackup_logfile](#), are known as a **raw backup**, because the changes that occurred during the backup operation are not yet incorporated. After you perform the **apply** step to the raw backup files, the resulting files do include those final data changes, and are known as a **prepared backup**. At this stage, the [ibbackup_logfile](#) file is no longer necessary.

See Also [apply](#), [hot backup](#), [MySQL Enterprise Backup](#), [prepared backup](#), [raw backup](#).

ibdata file

A set of files with names such as [ibdata1](#), [ibdata2](#), and so on, that make up the InnoDB **system tablespace**. These files contain metadata about InnoDB tables, (the **data dictionary**), and the storage areas for the **undo log**, the **change buffer**, and the **doublewrite buffer**. They also can contain some or all of the table data also (depending on whether the **file-per-table** mode is in effect when each table is created). When the **innodb_file_per_table** option is enabled, data and indexes for newly created tables are stored in separate **.ibd files** rather than in the system tablespace.

The growth of the [ibdata](#) files is influenced by the [innodb_autoextend_increment](#) configuration option.

See Also [change buffer](#), [data dictionary](#), [doublewrite buffer](#), [file-per-table](#), [.ibd file](#), [innodb_file_per_table](#), [system tablespace](#), [undo log](#).

ibtmp file

The InnoDB temporary tablespace data file for non-compressed InnoDB temporary tables and related objects. The configuration file option, [innodb_temp_data_file_path](#), allows users to define a relative path for the temporary data file. If [innodb_temp_data_file_path](#) is not specified, the default behavior is to create a single auto- extending 12MB data file named [ibtmp1](#) in the data directory, alongside [ibdata1](#).

See Also [temporary tablespace](#).

ib_logfile

A set of files, typically named [ib_logfile0](#) and [ib_logfile1](#), that form the **redo log**. Also sometimes referred to as the **log group**. These files record statements that attempt to change data in InnoDB tables. These statements are replayed automatically to correct data written by incomplete transactions, on startup following a crash.

This data cannot be used for manual recovery; for that type of operation, use the **binary log**.

See Also [binary log](#), [log group](#), [redo log](#).

ilist

Within an InnoDB **FULLTEXT index**, the data structure consisting of a document ID and positional information for a token (that is, a particular word).

See Also [FULLTEXT index](#).

implicit row lock

A row lock that InnoDB acquires to ensure consistency, without you specifically requesting it.

See Also [row lock](#).

in-memory database

A type of database system that maintains data in memory, to avoid overhead due to disk I/O and translation between disk blocks and memory areas. Some in-memory databases sacrifice durability (the "D" in the **ACID**

design philosophy) and are vulnerable to hardware, power, and other types of failures, making them more suitable for read-only operations. Other in-memory databases do use durability mechanisms such as logging changes to disk or using non-volatile memory.

MySQL features that address the same kinds of memory-intensive processing include the InnoDB **buffer pool**, **adaptive hash index**, and **read-only transaction** optimization, the MEMORY storage engine, the MyISAM key cache, and the MySQL **query cache**.

See Also [ACID](#), [adaptive hash index](#), [buffer pool](#), [disk-based](#), [read-only transaction](#).

incremental backup

A type of **hot backup**, performed by the **MySQL Enterprise Backup** product, that only saves data changed since some point in time. Having a full backup and a succession of incremental backups lets you reconstruct backup data over a long period, without the storage overhead of keeping several full backups on hand. You can restore the full backup and then apply each of the incremental backups in succession, or you can keep the full backup up-to-date by applying each incremental backup to it, then perform a single restore operation.

The granularity of changed data is at the **page** level. A page might actually cover more than one row. Each changed page is included in the backup.

See Also [hot backup](#), [MySQL Enterprise Backup](#), [page](#).

index

A data structure that provides a fast lookup capability for **rows** of a **table**, typically by forming a tree structure (**B-tree**) representing all the values of a particular **column** or set of columns.

InnoDB tables always have a **clustered index** representing the **primary key**. They can also have one or more **secondary indexes** defined on one or more columns. Depending on their structure, secondary indexes can be classified as **partial**, **column**, or **composite** indexes.

Indexes are a crucial aspect of **query** performance. Database architects design tables, queries, and indexes to allow fast lookups for data needed by applications. The ideal database design uses a **covering index** where practical; the query results are computed entirely from the index, without reading the actual table data. Each **foreign key** constraint also requires an index, to efficiently check whether values exist in both the **parent** and **child** tables.

Although a B-tree index is the most common, a different kind of data structure is used for **hash indexes**, as in the [MEMORY](#) storage engine and the InnoDB **adaptive hash index**.

See Also [adaptive hash index](#), [B-tree](#), [child table](#), [clustered index](#), [column index](#), [composite index](#), [covering index](#), [foreign key](#), [hash index](#), [parent table](#), [partial index](#), [primary key](#), [query](#), [row](#), [secondary index](#), [table](#).

index cache

A memory area that holds the token data for InnoDB **full-text search**. It buffers the data to minimize disk I/O when data is inserted or updated in columns that are part of a **FULLTEXT index**. The token data is written to disk when the index cache becomes full. Each InnoDB **FULLTEXT** index has its own separate index cache, whose size is controlled by the configuration option `innodb_ft_cache_size`.

See Also [full-text search](#), [FULLTEXT index](#).

index hint

Extended SQL syntax for overriding the **indexes** recommended by the optimizer. For example, the [FORCE INDEX](#), [USE INDEX](#), and [IGNORE INDEX](#) clauses. Typically used when indexed columns have unevenly distributed values, resulting in inaccurate **cardinality** estimates.

See Also [cardinality](#), [index](#).

index prefix

In an **index** that applies to multiple columns (known as a **composite index**), the initial or leading columns of the index. A query that references the first 1, 2, 3, and so on columns of a composite index can use the index, even if the query does not reference all the columns in the index.

See Also [composite index](#), [index](#).

index statistics

See [statistics](#).

infimum record

A **pseudo-record** in an **index**, representing the **gap** below the smallest value in that index. If a transaction has a statement such as `SELECT ... FOR UPDATE ... WHERE col < 10;`, and the smallest value in the column is 5, it is a lock on the infimum record that prevents other transactions from inserting even smaller values such as 0, -10, and so on.

See Also [gap](#), [index](#), [pseudo-record](#), [supremum record](#).

INFORMATION_SCHEMA

The name of the **database** that provides a query interface to the MySQL **data dictionary**. (This name is defined by the ANSI SQL standard.) To examine information (metadata) about the database, you can query tables such as `INFORMATION_SCHEMA.TABLES` and `INFORMATION_SCHEMA.COLUMNS`, rather than using `SHOW` commands that produce unstructured output.

The information schema contains some tables that are specific to **InnoDB**, such as `INNODB_LOCKS` and `INNODB_TRX`. You use these tables not to see how the database is structured, but to get real-time information about the workings of InnoDB tables to help with performance monitoring, tuning, and troubleshooting. In particular, these tables provide information about MySQL features related to **compression**, and **transactions** and their associated **locks**.

See Also [compression](#), [data dictionary](#), [database](#), [InnoDB](#), [lock](#), [transaction](#).

InnoDB

A MySQL component that combines high performance with **transactional** capability for reliability, robustness, and concurrent access. It embodies the **ACID** design philosophy. Represented as a **storage engine**; it handles tables created or altered with the `ENGINE=INNODB` clause. See [The InnoDB Storage Engine](#) for architectural details and administration procedures, and [Optimizing for InnoDB Tables](#) for performance advice.

In MySQL 5.5 and higher, InnoDB is the default storage engine for new tables and the `ENGINE=INNODB` clause is not required. In MySQL 5.1 only, many of the advanced InnoDB features require enabling the component known as the InnoDB Plugin. See [InnoDB as the Default MySQL Storage Engine](#) for the considerations involved in transitioning to recent releases where InnoDB tables are the default.

InnoDB tables are ideally suited for **hot backups**. See [MySQL Enterprise Backup](#) for information about the **MySQL Enterprise Backup** product for backing up MySQL servers without interrupting normal processing. See Also [ACID](#), [hot backup](#), [storage engine](#), [transaction](#).

innodb_autoinc_lock_mode

The `innodb_autoinc_lock_mode` option controls the algorithm used for **auto-increment locking**. When you have an auto-incrementing **primary key**, you can use statement-based replication only with the setting `innodb_autoinc_lock_mode=1`. This setting is known as **consecutive** lock mode, because multi-row inserts within a transaction receive consecutive auto-increment values. If you have `innodb_autoinc_lock_mode=2`, which allows higher concurrency for insert operations, use row-based replication rather than statement-based replication. This setting is known as **interleaved** lock mode, because multiple multi-row insert statements running at the same time can receive autoincrement values that are interleaved. The setting `innodb_autoinc_lock_mode=0` is the previous (traditional) default setting and should not be used except for compatibility purposes.

See Also [auto-increment locking](#), [mixed-mode insert](#), [primary key](#).

innodb_file_format

The `innodb_file_format` option determines the **file format** for all InnoDB **tablespaces** created after you specify a value for this option. To create tablespaces other than the **system tablespace**, you must also use the **file-per-table** option. Currently, you can specify the **Antelope** and **Barracuda** file formats.

See Also [Antelope](#), [Barracuda](#), [file format](#), [file-per-table](#), [innodb_file_per_table](#), [system tablespace](#), [tablespace](#).

innodb_file_per_table

A very important configuration option that affects many aspects of InnoDB file storage, availability of features, and I/O characteristics. In MySQL 5.6.7 and higher, it is enabled by default. Prior to MySQL 5.6.7, it is disabled by default. The [innodb_file_per_table](#) option turns on **file-per-table** mode, which stores each newly created InnoDB table and its associated index in its own **.ibd file**, outside the **system tablespace**.

This option affects the performance and storage considerations for a number of SQL statements, such as [DROP TABLE](#) and [TRUNCATE TABLE](#).

This option is needed to take full advantage of many other InnoDB features, such as such as table **compression**, or backups of named tables in **MySQL Enterprise Backup**.

This option was once static, but can now be set using the [SET GLOBAL](#) command.

For reference information, see [innodb_file_per_table](#). For usage information, see [Using Per-Table Tablespaces](#).

See Also [compression](#), [file-per-table](#), [.ibd file](#), [MySQL Enterprise Backup](#), [system tablespace](#).

innodb_lock_wait_timeout

The [innodb_lock_wait_timeout](#) option sets the balance between **waiting** for shared resources to become available, or giving up and handling the error, retrying, or doing alternative processing in your application. Rolls back any InnoDB transaction that waits more than a specified time to acquire a **lock**. Especially useful if **deadlocks** are caused by updates to multiple tables controlled by different storage engines; such deadlocks are not **detected** automatically.

See Also [deadlock](#), [deadlock detection](#), [lock](#), [wait](#).

innodb_strict_mode

The [innodb_strict_mode](#) option controls whether InnoDB operates in **strict mode**, where conditions that are normally treated as warnings, cause errors instead (and the underlying statements fail).

This mode is the default setting in MySQL 5.5.5 and higher.

See Also [strict mode](#).

insert

One of the primary **DML** operations in **SQL**. The performance of inserts is a key factor in **data warehouse** systems that load millions of rows into tables, and **OLTP** systems where many concurrent connections might insert rows into the same table, in arbitrary order. If insert performance is important to you, you should learn about **InnoDB** features such as the **insert buffer** used in **change buffering**, and **auto-increment** columns.

See Also [auto-increment](#), [change buffering](#), [data warehouse](#), [DML](#), [InnoDB](#), [insert buffer](#), [OLTP](#), [SQL](#).

insert buffer

Former name for the **change buffer**. Now that **change buffering** includes delete and update operations as well as inserts, "change buffer" is the preferred term.

See Also [change buffer](#), [change buffering](#).

insert buffering

The technique of storing secondary index changes due to [INSERT](#) operations in the **insert buffer** rather than writing them immediately, so that the physical writes can be performed to minimize random I/O. It is one of the types of **change buffering**; the others are **delete buffering** and **purge buffering**.

Insert buffering is not used if the secondary index is **unique**, because the uniqueness of new values cannot be verified before the new entries are written out. Other kinds of change buffering do work for unique indexes.

See Also [change buffer](#), [change buffering](#), [delete buffering](#), [insert buffer](#), [purge buffering](#), [unique index](#).

instance

A single **mysqld** daemon managing a **data directory** representing one or more **databases** with a set of **tables**. It is common in development, testing, and some **replication** scenarios to have multiple instances on the same **server** machine, each managing its own data directory and listening on its own port or socket. With one instance running a **disk-bound** workload, the server might still have extra CPU and memory capacity to run additional instances.

See Also [data directory](#), [database](#), [disk-bound](#), [mysqld](#), [replication](#), [server](#).

instrumentation

Modifications at the source code level to collect performance data for tuning and debugging. In MySQL, data collected by instrumentation is exposed through a SQL interface using the [INFORMATION_SCHEMA](#) and [PERFORMANCE_SCHEMA](#) databases.

See Also [INFORMATION_SCHEMA](#), [Performance Schema](#).

intention exclusive lock

See [intention lock](#).

intention lock

A kind of **lock** that applies to the table level, used to indicate what kind of lock the transaction intends to acquire on rows in the table. Different transactions can acquire different kinds of intention locks on the same table, but the first transaction to acquire an **intention exclusive** (IX) lock on a table prevents other transactions from acquiring any S or X locks on the table. Conversely, the first transaction to acquire an **intention shared** (IS) lock on a table prevents other transactions from acquiring any X locks on the table. The two-phase process allows the lock requests to be resolved in order, without blocking locks and corresponding operations that are compatible. For more details on this locking mechanism, see [InnoDB Lock Modes](#).

See Also [lock](#), [lock mode](#), [locking](#).

intention shared lock

See [intention lock](#).

inverted index

A data structure optimized for document retrieval systems, used in the implementation of InnoDB **full-text search**. The InnoDB **FULLTEXT index**, implemented as an inverted index, records the position of each word within a document, rather than the location of a table row. A single column value (a document stored as a text string) is represented by many entries in the inverted index.

See Also [full-text search](#), [FULLTEXT index](#), [iist](#).

IOPS

Acronym for **I/O operations per second**. A common measurement for busy systems, particularly **OLTP** applications. If this value is near the maximum that the storage devices can handle, the application can become **disk-bound**, limiting **scalability**.

See Also [disk-bound](#), [OLTP](#), [scalability](#).

isolation level

One of the foundations of database processing. Isolation is the **I** in the acronym **ACID**; the isolation level is the setting that fine-tunes the balance between performance and reliability, consistency, and reproducibility of results when multiple **transactions** are making changes and performing queries at the same time.

From highest amount of consistency and protection to the least, the isolation levels supported by InnoDB are: **SERIALIZABLE**, **REPEATABLE READ**, **READ COMMITTED**, and **READ UNCOMMITTED**.

With InnoDB tables, many users can keep the default isolation level (**REPEATABLE READ**) for all operations. Expert users might choose the **read committed** level as they push the boundaries of scalability with OLTP processing, or during data warehousing operations where minor inconsistencies do not affect the aggregate results of large amounts of data. The levels on the edges (**SERIALIZABLE** and **READ UNCOMMITTED**) change the processing behavior to such an extent that they are rarely used.

See Also [ACID](#), [READ COMMITTED](#), [READ UNCOMMITTED](#), [REPEATABLE READ](#), [SERIALIZABLE](#), [transaction](#).

J

join

A **query** that retrieves data from more than one table, by referencing columns in the tables that hold identical values. Ideally, these columns are part of an InnoDB **foreign key** relationship, which ensures **referential integrity** and that the join columns are **indexed**. Often used to save space and improve query performance by replacing repeated strings with numeric IDs, in a **normalized** data design.

See Also [foreign key](#), [index](#), [normalized](#), [query](#), [referential integrity](#).

K

KEY_BLOCK_SIZE

An option to specify the size of data pages within an InnoDB table that uses **compressed row format**. The default is 8 kilobytes. Lower values risk hitting internal limits that depend on the combination of row size and compression percentage.

See Also [compressed row format](#).

L

latch

A lightweight structure used by InnoDB to implement a **lock** for its own internal memory structures, typically held for a brief time measured in milliseconds or microseconds. A general term that includes both **mutexes** (for exclusive access) and **rw-locks** (for shared access). Certain latches are the focus of InnoDB performance tuning, such as the **data dictionary** mutex. Statistics about latch use and contention are available through the **Performance Schema** interface.

See Also [data dictionary](#), [lock](#), [locking](#), [mutex](#), [Performance Schema](#), [rw-lock](#).

list

The InnoDB **buffer pool** is represented as a list of memory **pages**. The list is reordered as new pages are accessed and enter the buffer pool, as pages within the buffer pool are accessed again and are considered newer, and as pages that are not accessed for a long time are **evicted** from the buffer pool. The buffer pool is actually divided into **sublists**, and the replacement policy is a variation of the familiar **LRU** technique.

See Also [buffer pool](#), [eviction](#), [LRU](#), [sublist](#).

lock

The high-level notion of an object that controls access to a resource, such as a table, row, or internal data structure, as part of a **locking** strategy. For intensive performance tuning, you might delve into the actual structures that implement locks, such as **mutexes** and **latches**.

See Also [latch](#), [lock mode](#), [locking](#), [mutex](#).

lock escalation

An operation used in some database systems that converts many row locks into a single table lock, saving memory space but reducing concurrent access to the table. InnoDB uses a space-efficient representation for row locks, so that lock escalation is not needed.

See Also [locking](#), [row lock](#), [table lock](#).

lock mode

A shared (S) lock allows a transaction to read a row. Multiple transactions can acquire an S lock on that same row at the same time.

An exclusive (X) lock allows a transaction to update or delete a row. No other transaction can acquire any kind of lock on that same row at the same time.

Intention locks apply to the table level, and are used to indicate what kind of lock the transaction intends to acquire on rows in the table. Different transactions can acquire different kinds of intention locks on the same table, but the first transaction to acquire an intention exclusive (IX) lock on a table prevents other transactions from acquiring any S or X locks on the table. Conversely, the first transaction to acquire an intention shared (IS) lock on a table prevents other transactions from acquiring any X locks on the table. The two-phase process allows the lock requests to be resolved in order, without blocking locks and corresponding operations that are compatible.

See Also [intention lock](#), [lock](#), [locking](#).

locking

The system of protecting a **transaction** from seeing or changing data that is being queried or changed by other transactions. The locking strategy must balance reliability and consistency of database operations (the principles of the **ACID** philosophy) against the performance needed for good **concurrency**. Fine-tuning the locking strategy often involves choosing an **isolation level** and ensuring all your database operations are safe and reliable for that isolation level.

See Also [ACID](#), [concurrency](#), [isolation level](#), [latch](#), [lock](#), [mutex](#), [transaction](#).

locking read

A **SELECT** statement that also performs a **locking** operation on an InnoDB table. Either `SELECT ... FOR UPDATE` or `SELECT ... LOCK IN SHARE MODE`. It has the potential to produce a **deadlock**, depending on the **isolation level** of the transaction. The opposite of a **non-locking read**. Not allowed for global tables in a **read-only transaction**.

See Also [deadlock](#), [isolation level](#), [locking](#), [non-locking read](#), [read-only transaction](#).

log

In the InnoDB context, "log" or "log files" typically refers to the **redo log** represented by the `ib_logfile*` files. Another log area, which is physically part of the **system tablespace**, is the **undo log**.

Other kinds of logs that are important in MySQL are the **error log** (for diagnosing startup and runtime problems), **binary log** (for working with replication and performing point-in-time restores), the **general query log** (for diagnosing application problems), and the **slow query log** (for diagnosing performance problems).

See Also [binary log](#), [error log](#), [general query log](#), [ib_logfile](#), [redo log](#), [slow query log](#), [system tablespace](#), [undo log](#).

log buffer

The memory area that holds data to be written to the **log files** that make up the **redo log**. It is controlled by the `innodb_log_buffer_size` configuration option.

See Also [log file](#), [redo log](#).

log file

One of the `ib_logfileN` files that make up the **redo log**. Data is written to these files from the **log buffer** memory area.

See Also [ib_logfile](#), [log buffer](#), [redo log](#).

log group

The set of files that make up the **redo log**, typically named `ib_logfile0` and `ib_logfile1`. (For that reason, sometimes referred to collectively as **ib_logfile**.)

See Also [ib_logfile](#), [redo log](#).

logical

A type of operation that involves high-level, abstract aspects such as tables, queries, indexes, and other SQL concepts. Typically, logical aspects are important to make database administration and application development convenient and usable. Contrast with **physical**.

See Also [logical backup](#), [physical](#).

logical backup

A **backup** that reproduces table structure and data, without copying the actual data files. For example, the [mysqldump](#) command produces a logical backup, because its output contains statements such as [CREATE TABLE](#) and [INSERT](#) that can re-create the data. Contrast with **physical backup**. A logical backup offers flexibility (for example, you could edit table definitions or insert statements before restoring), but can take substantially longer to **restore** than a physical backup.

See Also [backup](#), [mysqldump](#), [physical backup](#), [restore](#).

loose_

In MySQL 5.1, a prefix added to InnoDB configuration options when installing the InnoDB **Plugin** after server startup, so any new configuration options not recognized by the current level of MySQL do not cause a startup failure. MySQL processes configuration options that start with this prefix, but gives a warning rather than a failure if the part after the prefix is not a recognized option.

See Also [plugin](#).

low-water mark

A value representing a lower limit, typically a threshold value at which some corrective action begins or becomes more aggressive. Contrast with **high-water mark**.

See Also [high-water mark](#).

LRU

An acronym for "least recently used", a common method for managing storage areas. The items that have not been used recently are **evicted** when space is needed to cache newer items. InnoDB uses the LRU mechanism by default to manage the **pages** within the **buffer pool**, but makes exceptions in cases where a page might be read only a single time, such as during a **full table scan**. This variation of the LRU algorithm is called the **midpoint insertion strategy**. The ways in which the buffer pool management differs from the traditional LRU algorithm is fine-tuned by the options [innodb_old_blocks_pct](#), [innodb_old_blocks_time](#), and the new MySQL 5.6 options [innodb_lru_scan_depth](#) and [innodb_flush_neighbors](#).

See Also [buffer pool](#), [eviction](#), [full table scan](#), [midpoint insertion strategy](#), [page](#).

LSN

Acronym for "log sequence number". This arbitrary, ever-increasing value represents a point in time corresponding to operations recorded in the **redo log**. (This point in time is regardless of **transaction** boundaries; it can fall in the middle of one or more transactions.) It is used internally by InnoDB during **crash recovery** and for managing the buffer pool.

Prior to MySQL 5.6.3, the LSN was a 4-byte unsigned integer. The LSN became an 8-byte unsigned integer in MySQL 5.6.3 when the redo log file size limit increased from 4GB to 512GB, as additional bytes were required to store extra size information. Applications built on MySQL 5.6.3 or later that use LSN values should use 64-bit rather than 32-bit variables to store and compare LSN values.

In the **MySQL Enterprise Backup** product, you can specify an LSN to represent the point in time from which to take an **incremental backup**. The relevant LSN is displayed by the output of the [mysqlbackup](#) command. Once you have the LSN corresponding to the time of a full backup, you can specify that value to take a subsequent incremental backup, whose output contains another LSN for the next incremental backup.

See Also [crash recovery](#), [incremental backup](#), [MySQL Enterprise Backup](#), [redo log](#), [transaction](#).

M

.MRG file

A file containing references to other tables, used by the [MERGE](#) storage engine. Files with this extension are always included in backups produced by the [mysqlbackup](#) command of the **MySQL Enterprise Backup** product.

See Also [MySQL Enterprise Backup](#), [mysqlbackup command](#).

.MYD file

A file that MySQL uses to store data for a MyISAM table.

See Also [.MYI file](#), [MySQL Enterprise Backup](#), [mysqlbackup command](#).

.MYI file

A file that MySQL uses to store indexes for a MyISAM table.

See Also [.MYD file](#), [MySQL Enterprise Backup](#), [mysqlbackup command](#).

master server

Frequently shortened to "master". A database server machine in a **replication** scenario that processes the initial insert, update, and delete requests for data. These changes are propagated to, and repeated on, other servers known as **slave servers**.

See Also [replication](#), [slave server](#).

master thread

An InnoDB **thread** that performs various tasks in the background. Most of these tasks are I/O related, such as writing changes from the **insert buffer** to the appropriate secondary indexes.

To improve **concurrency**, sometimes actions are moved from the master thread to separate background threads. For example, in MySQL 5.6 and higher, **dirty pages** are **flushed** from the **buffer pool** by the **page cleaner** thread rather than the master thread.

See Also [buffer pool](#), [dirty page](#), [flush](#), [insert buffer](#), [page cleaner](#), [thread](#).

MDL

Acronym for "metadata lock".

See Also [metadata lock](#).

memcached

A popular component of many MySQL and **NoSQL** software stacks, allowing fast reads and writes for single values and caching the results entirely in memory. Traditionally, applications required extra logic to write the same data to a MySQL database for permanent storage, or to read data from a MySQL database when it was not cached yet in memory. Now, applications can use the simple [memcached](#) protocol, supported by client libraries for many languages, to communicate directly with MySQL servers using **InnoDB** or MySQL Cluster tables. These NoSQL interfaces to MySQL tables allow applications to achieve higher read and write performance than by issuing SQL commands directly, and can simplify application logic and deployment configurations for systems that already incorporated [memcached](#) for in-memory caching.

The [memcached](#) interface to InnoDB tables is available in MySQL 5.6 and higher; see [InnoDB Integration with memcached](#) for details. The [memcached](#) interface to MySQL Cluster tables is available in MySQL Cluster 7.2; see <http://dev.mysql.com/doc/ndbapi/en/ndbmemcache.html> for details.

See Also [InnoDB](#), [NoSQL](#).

merge

To apply changes to data cached in memory, such as when a page is brought into the **buffer pool**, and any applicable changes recorded in the **change buffer** are incorporated into the page in the buffer pool. The updated data is eventually written to the **tablespace** by the **flush** mechanism.

See Also [buffer pool](#), [change buffer](#), [flush](#), [tablespace](#).

metadata lock

A type of **lock** that prevents **DDL** operations on a table that is being used at the same time by another **transaction**. For details, see [Metadata Locking](#).

Enhancements to **online** operations, particularly in MySQL 5.6 and higher, are focused on reducing the amount of metadata locking. The objective is for DDL operations that do not change the table structure (such as [CREATE](#)

[INDEX](#) and [DROP INDEX](#) for [InnoDB](#) tables) to proceed while the table is being queried, updated, and so on by other transactions.

See Also [DDL](#), [lock](#), [online](#), [transaction](#).

metrics counter

A feature implemented by the [innodb_metrics](#) table in the [information_schema](#), in MySQL 5.6 and higher. You can query **counts** and totals for low-level InnoDB operations, and use the results for performance tuning in combination with data from the [performance_schema](#).

See Also [counter](#), [INFORMATION_SCHEMA](#), [Performance Schema](#).

midpoint insertion strategy

The technique of initially bringing **pages** into the InnoDB **buffer pool** not at the "newest" end of the list, but instead somewhere in the middle. The exact location of this point can vary, based on the setting of the [innodb_old_blocks_pct](#) option. The intent is that blocks that are only read once, such as during a **full table scan**, can be aged out of the buffer pool sooner than with a strict **LRU** algorithm.

See Also [buffer pool](#), [full table scan](#), [LRU](#), [page](#).

mini-transaction

An internal phase of InnoDB processing, when making changes at the **physical** level to internal data structures during **DML** operations. A mini-transaction has no notion of **rollback**; multiple mini-transactions can occur within a single **transaction**. Mini-transactions write information to the **redo log** that is used during **crash recovery**. A mini-transaction can also happen outside the context of a regular transaction, for example during **purge** processing by background threads.

See Also [commit](#), [crash recovery](#), [DML](#), [physical](#), [purge](#), [redo log](#), [rollback](#), [transaction](#).

mixed-mode insert

An [INSERT](#) statement where **auto-increment** values are specified for some but not all of the new rows. For example, a multi-value [INSERT](#) could specify a value for the auto-increment column in some cases and [NULL](#) in other cases. [InnoDB](#) generates auto-increment values for the rows where the column value was specified as [NULL](#). Another example is an [INSERT ... ON DUPLICATE KEY UPDATE](#) statement, where auto-increment values might be generated but not used, for any duplicate rows that are processed as [UPDATE](#) rather than [INSERT](#) statements.

Can cause consistency issues between **master** and **slave** servers in a **replication** configuration. Can require adjusting the value of the [innodb_autoinc_lock_mode](#) configuration option.

See Also [auto-increment](#), [innodb_autoinc_lock_mode](#), [master server](#), [replication](#), [slave server](#).

multi-core

A type of processor that can take advantage of multi-threaded programs, such as the MySQL server.

multiversion concurrency control

See [MVCC](#).

mutex

Informal abbreviation for "mutex variable". (Mutex itself is short for "mutual exclusion".) The low-level object that InnoDB uses to represent and enforce exclusive-access **locks** to internal in-memory data structures. Once the lock is acquired, any other process, thread, and so on is prevented from acquiring the same lock. Contrast with **rw-locks**, which allow shared access. Mutexes and rw-locks are known collectively as **latches**.

See Also [latch](#), [lock](#), [Performance Schema](#), [Pthreads](#), [rw-lock](#).

MVCC

Acronym for "multiversion concurrency control". This technique lets InnoDB **transactions** with certain **isolation levels** to perform **consistent read** operations; that is, to query rows that are being updated by other transactions, and see the values from before those updates occurred. This is a powerful technique to increase **concurrency**, by allowing queries to proceed without waiting due to **locks** held by the other transactions.

This technique is not universal in the database world. Some other database products, and some other MySQL storage engines, do not support it.

See Also [ACID](#), [concurrency](#), [consistent read](#), [isolation level](#), [lock](#), [transaction](#).

my.cnf

The name, on UNIX or Linux systems, of the MySQL option file.

See Also [my.ini](#), [option file](#).

my.ini

The name, on Windows systems, of the MySQL option file.

See Also [my.cnf](#), [option file](#).

mysql

The `mysql` program is the command-line interpreter for the MySQL database. It processes **SQL** statements, and also MySQL-specific commands such as `SHOW TABLES`, by passing requests to the `mysqld` daemon.

See Also [mysqld](#), [SQL](#).

MySQL Enterprise Backup

A licensed product that performs **hot backups** of MySQL databases. It offers the most efficiency and flexibility when backing up **InnoDB** tables, but can also back up MyISAM and other kinds of tables.

See Also [hot backup](#), [InnoDB](#).

mysqlbackup command

A command-line tool of the **MySQL Enterprise Backup** product. It performs a **hot backup** operation for InnoDB tables, and a **warm backup** for MyISAM and other kinds of tables. See [MySQL Enterprise Backup](#) for more information about this command.

See Also [hot backup](#), [MySQL Enterprise Backup](#), [warm backup](#).

mysqld

The `mysqld` program is the database engine for the MySQL database. It runs as a UNIX daemon or Windows service, constantly waiting for requests and performing maintenance work in the background.

See Also [mysql](#).

mysqldump

A command that performs a **logical backup** of some combination of databases, tables, and table data. The results are SQL statements that reproduce the original schema objects, data, or both. For substantial amounts of data, a **physical backup** solution such as **MySQL Enterprise Backup** is faster, particularly for the **restore** operation.

See Also [logical backup](#), [MySQL Enterprise Backup](#), [physical backup](#), [restore](#).

N

natural key

A indexed column, typically a **primary key**, where the values have some real-world significance. Usually advised against because:

- If the value should ever change, there is potentially a lot of index maintenance to re-sort the **clustered index** and update the copies of the primary key value that are repeated in each **secondary index**.
- Even seemingly stable values can change in unpredictable ways that are difficult to represent correctly in the database. For example, one country can change into two or several, making the original country code obsolete. Or, rules about unique values might have exceptions. For example, even if taxpayer IDs are intended to be unique to a single person, a database might have to handle records that violate that rule, such as in cases of identity theft. Taxpayer IDs and other sensitive ID numbers also make poor primary keys, because they may need to be secured, encrypted, and otherwise treated differently than other columns.

Thus, it is typically better to use arbitrary numeric values to form a **synthetic key**, for example using an **auto-increment** column.

See Also [auto-increment](#), [primary key](#), [secondary index](#), [synthetic key](#).

neighbor page

Any **page** in the same **extent** as a particular page. When a page is selected to be **flushed**, any neighbor pages that are **dirty** are typically flushed as well, as an I/O optimization for traditional hard disks. In MySQL 5.6 and up, this behavior can be controlled by the configuration variable `innodb_flush_neighbors`; you might turn that setting off for SSD drives, which do not have the same overhead for writing smaller batches of data at random locations.

See Also [dirty page](#), [extent](#), [flush](#), [page](#).

next-key lock

A combination of a **record lock** on the index record and a [gap lock](#) on the gap before the index record.

See Also [gap lock](#), [locking](#), [record lock](#).

non-blocking I/O

An industry term that means the same as **asynchronous I/O**.

See Also [asynchronous I/O](#).

non-locking read

A **query** that does not use the `SELECT ... FOR UPDATE` or `SELECT ... LOCK IN SHARE MODE` clauses.

The only kind of query allowed for global tables in a **read-only transaction**. The opposite of a **locking read**.

See Also [locking read](#), [query](#), [read-only transaction](#).

non-repeatable read

The situation when a query retrieves data, and a later query within the same **transaction** retrieves what should be the same data, but the queries return different results (changed by another transaction committing in the meantime).

This kind of operation goes against the **ACID** principle of database design. Within a transaction, data should be consistent, with predictable and stable relationships.

Among different **isolation levels**, non-repeatable reads are prevented by the **serializable read** and **repeatable read** levels, and allowed by the **consistent read**, and **read uncommitted** levels.

See Also [ACID](#), [consistent read](#), [isolation level](#), [READ UNCOMMITTED](#), [REPEATABLE READ](#), [SERIALIZABLE](#), [transaction](#).

normalized

A database design strategy where data is split into multiple tables, and duplicate values condensed into single rows represented by an ID, to avoid storing, querying, and updating redundant or lengthy values. It is typically used in **OLTP** applications.

For example, an address might be given a unique ID, so that a census database could represent the relationship **lives at this address** by associating that ID with each member of a family, rather than storing multiple copies of a complex value such as **123 Main Street, Anytown, USA**.

For another example, although a simple address book application might store each phone number in the same table as a person's name and address, a phone company database might give each phone number a special ID, and store the numbers and IDs in a separate table. This normalized representation could simplify large-scale updates when area codes split apart.

Normalization is not always recommended. Data that is primarily queried, and only updated by deleting entirely and reloading, is often kept in fewer, larger tables with redundant copies of duplicate values. This data representation is referred to as **denormalized**, and is frequently found in data warehousing applications.

See Also [denormalized](#), [foreign key](#), [OLTP](#), [relational](#).

NoSQL

A broad term for a set of data access technologies that do not use the **SQL** language as their primary mechanism for reading and writing data. Some NoSQL technologies act as key-value stores, only accepting single-value reads and writes; some relax the restrictions of the **ACID** methodology; still others do not require a pre-planned **schema**. MySQL users can combine NoSQL-style processing for speed and simplicity with SQL operations for flexibility and convenience, by using the **memcached** API to directly access some kinds of MySQL tables. The [memcached](#) interface to InnoDB tables is available in MySQL 5.6 and higher; see [InnoDB Integration with memcached](#) for details. The [memcached](#) interface to MySQL Cluster tables is available in MySQL Cluster 7.2; see <http://dev.mysql.com/doc/ndbapi/en/ndbmemcache.html> for details.

See Also [ACID](#), [InnoDB](#), [memcached](#), [schema](#), [SQL](#).

NOT NULL constraint

A type of **constraint** that specifies that a **column** cannot contain any **NULL** values. It helps to preserve **referential integrity**, as the database server can identify data with erroneous missing values. It also helps in the arithmetic involved in query optimization, allowing the optimizer to predict the number of entries in an index on that column.

See Also [column](#), [constraint](#), [NULL](#), [primary key](#), [referential integrity](#).

NULL

A special value in **SQL**, indicating the absence of data. Any arithmetic operation or equality test involving a [NULL](#) value, in turn produces a [NULL](#) result. (Thus it is similar to the IEEE floating-point concept of NaN, "not a number".) Any aggregate calculation such as [AVG\(\)](#) ignores rows with [NULL](#) values, when determining how many rows to divide by. The only test that works with [NULL](#) values uses the SQL idioms [IS NULL](#) or [IS NOT NULL](#).

[NULL](#) values play a part in index operations, because for performance a database must minimize the overhead of keeping track of missing data values. Typically, [NULL](#) values are not stored in an index, because a query that tests an indexed column using a standard comparison operator could never match a row with a [NULL](#) value for that column. For the same reason, unique indexes do not prevent [NULL](#) values; those values simply are not represented in the index. Declaring a [NOT NULL](#) constraint on a column provides reassurance that there are no rows left out of the index, allowing for better query optimization (accurate counting of rows and estimation of whether to use the index).

Because the **primary key** must be able to uniquely identify every row in the table, a single-column primary key cannot contain any [NULL](#) values, and a multi-column primary key cannot contain any rows with [NULL](#) values in all columns.

Although the Oracle database allows a [NULL](#) value to be concatenated with a string, InnoDB treats the result of such an operation as [NULL](#).

See Also [index](#), [primary key](#), [SQL](#).

O

.OPT file

A file containing database configuration information. Files with this extension are always included in backups produced by the [mysqlbackup](#) command of the **MySQL Enterprise Backup** product.

See Also [MySQL Enterprise Backup](#), [mysqlbackup command](#).

off-page column

A column containing variable-length data (such as [BLOB](#) and [VARCHAR](#)) that is too long to fit on a **B-tree** page. The data is stored in **overflow pages**. The [DYNAMIC](#) row format in the InnoDB **Barracuda** file format is more efficient for such storage than the older [COMPACT](#) row format.

See Also [B-tree](#), [Barracuda](#), [overflow page](#).

OLTP

Acronym for "Online Transaction Processing". A database system, or a database application, that runs a workload with many **transactions**, with frequent writes as well as reads, typically affecting small amounts of data at a time. For example, an airline reservation system or an application that processes bank deposits. The data might be organized in **normalized** form for a balance between **DML** (insert/update/delete) efficiency and **query** efficiency. Contrast with **data warehouse**.

With its **row-level locking** and **transactional** capability, **InnoDB** is the ideal storage engine for MySQL tables used in OLTP applications.

See Also [data warehouse](#), [DML](#), [InnoDB](#), [query](#), [row lock](#), [transaction](#).

online

A type of operation that involves no downtime, blocking, or restricted operation for the database. Typically applied to **DDL**. Operations that shorten the periods of restricted operation, such as **fast index creation**, have evolved into a wider set of **online DDL** operations in MySQL 5.6.

In the context of backups, a **hot backup** is an online operation and a **warm backup** is partially an online operation.

See Also [DDL](#), [Fast Index Creation](#), [hot backup](#), [online DDL](#), [warm backup](#).

online DDL

A feature that improves the performance, concurrency, and availability of InnoDB tables during **DDL** (primarily [ALTER TABLE](#)) operations. See [InnoDB and Online DDL](#) for details.

The details vary according to the type of operation. In some cases, the table can be modified concurrently while the [ALTER TABLE](#) is in progress. The operation might be able to be performed without doing a table copy, or using a specially optimized type of table copy. Space usage is controlled by the [innodb_online_alter_log_max_size](#) configuration option.

This feature is an enhancement of the **Fast Index Creation** feature in MySQL 5.5 and the InnoDB Plugin for MySQL 5.1.

See Also [DDL](#), [Fast Index Creation](#), [online](#).

optimistic

A methodology that guides low-level implementation decisions for a relational database system. The requirements of performance and **concurrency** in a relational database mean that operations must be started or dispatched quickly. The requirements of consistency and **referential integrity** mean that any operation could fail: a transaction might be rolled back, a **DML** operation could violate a constraint, a request for a lock could cause a deadlock, a network error could cause a timeout. An optimistic strategy is one that assumes most requests or attempts will succeed, so that relatively little work is done to prepare for the failure case. When this assumption is true, the database does little unnecessary work; when requests do fail, extra work must be done to clean up and undo changes.

InnoDB uses optimistic strategies for operations such as **locking** and **commits**. For example, data changed by a transaction can be written to the data files before the commit occurs, making the commit itself very fast, but requiring more work to undo the changes if the transaction is rolled back.

The opposite of an optimistic strategy is a **pessimistic** one, where a system is optimized to deal with operations that are unreliable and frequently unsuccessful. This methodology is rare in a database system, because so much care goes into choosing reliable hardware, networks, and algorithms.

See Also [commit](#), [concurrency](#), [DML](#), [locking](#), [pessimistic](#).

optimizer

The MySQL component that determines the best **indexes** and **join** order to use for a **query**, based on characteristics and data distribution of the relevant **tables**.

See Also [index](#), [join](#), [query](#), [table](#).

option

A configuration parameter for MySQL, either stored in the **option file** or passed on the command line.

For the options that apply to **InnoDB** tables, each option name starts with the prefix `innodb_`.

See Also [InnoDB](#), [option file](#).

option file

The file that holds the configuration **options** for the MySQL instance. Traditionally, on Linux and UNIX this file is named `my.cnf`, and on Windows it is named `my.ini`.

See Also [configuration file](#), [my.cnf](#), [option](#).

overflow page

Separately allocated disk **pages** that hold variable-length columns (such as [BLOB](#) and [VARCHAR](#)) that are too long to fit on a **B-tree** page. The associated columns are known as **off-page columns**.

See Also [B-tree](#), [off-page column](#), [page](#).

P

.PAR file

A file containing partition definitions. Files with this extension are always included in backups produced by the `mysqlbackup` command of the **MySQL Enterprise Backup** product.

See Also [MySQL Enterprise Backup](#), [mysqlbackup command](#).

page

A unit representing how much data InnoDB transfers at any one time between disk (the **data files**) and memory (the **buffer pool**). A page can contain one or more **rows**, depending on how much data is in each row. If a row does not fit entirely into a single page, InnoDB sets up additional pointer-style data structures so that the information about the row can be stored in one page.

One way to fit more data in each page is to use **compressed row format**. For tables that use BLOBs or large text fields, **compact row format** allows those large columns to be stored separately from the rest of the row, reducing I/O overhead and memory usage for queries that do not reference those columns.

When InnoDB reads or writes sets of pages as a batch to increase I/O throughput, it reads or writes an **extent** at a time.

All the InnoDB disk data structures within a MySQL instance share the same **page size**.

See Also [buffer pool](#), [compact row format](#), [compressed row format](#), [data files](#), [extent](#), [page size](#), [row](#).

page cleaner

An InnoDB background **thread** that **flushes dirty pages** from the **buffer pool**. Prior to MySQL 5.6, this activity was performed by the **master thread**

See Also [buffer pool](#), [dirty page](#), [flush](#), [master thread](#), [thread](#).

page size

For releases up to and including MySQL 5.5, the size of each InnoDB **page** is fixed at 16 kilobytes. This value represents a balance: large enough to hold the data for most rows, yet small enough to minimize the performance overhead of transferring unneeded data to memory. Other values are not tested or supported.

Starting in MySQL 5.6, the page size for an InnoDB **instance** can be either 4KB, 8KB, or 16KB, controlled by the `innodb_page_size` configuration option. You set the size when creating the MySQL instance, and it remains constant afterwards. The same page size applies to all InnoDB **tablespaces**, both the **system tablespace** and any separate tablespaces created in **file-per-table** mode.

Smaller page sizes can help performance with storage devices that use small block sizes, particularly for **SSD** devices in **disk-bound** workloads, such as for **OLTP** applications. As individual rows are updated, less data is copied into memory, written to disk, reorganized, locked, and so on.

See Also [disk-bound](#), [file-per-table](#), [instance](#), [OLTP](#), [page](#), [SSD](#), [system tablespace](#), [tablespace](#).

parent table

The table in a **foreign key** relationship that holds the initial column values pointed to from the **child table**. The consequences of deleting, or updating rows in the parent table depend on the [ON UPDATE](#) and [ON DELETE](#) clauses in the foreign key definition. Rows with corresponding values in the child table could be automatically deleted or updated in turn, or those columns could be set to [NULL](#), or the operation could be prevented.

See Also [child table](#), [foreign key](#).

partial backup

A **backup** that contains some of the **tables** in a MySQL database, or some of the databases in a MySQL instance. Contrast with **full backup**.

See Also [backup](#), [full backup](#), [table](#).

partial index

An **index** that represents only part of a column value, typically the first N characters (the **prefix**) of a long [VARCHAR](#) value.

See Also [index](#), [index prefix](#).

Performance Schema

The [performance_schema](#) schema, in MySQL 5.5 and up, presents a set of tables that you can query to get detailed information about the performance characteristics of many internal parts of the MySQL server.

See Also [latch](#), [mutex](#), [rw-lock](#).

persistent statistics

A feature in MySQL 5.6 that stores **index** statistics for InnoDB **tables** on disk, providing better **plan stability** for **queries**.

See Also [index](#), [optimizer](#), [plan stability](#), [query](#), [table](#).

pessimistic

A methodology that sacrifices performance or concurrency in favor of safety. It is appropriate if a high proportion of requests or attempts might fail, or if the consequences of a failed request are severe. InnoDB uses what is known as a pessimistic **locking** strategy, to minimize the chance of **deadlocks**. At the application level, you might avoid deadlocks by using a pessimistic strategy of acquiring all locks needed by a transaction at the very beginning.

Many built-in database mechanisms use the opposite **optimistic** methodology.

See Also [deadlock](#), [locking](#), [optimistic](#).

phantom

A row that appears in the result set of a query, but not in the result set of an earlier query. For example, if a query is run twice within a **transaction**, and in the meantime, another transaction commits after inserting a new row or updating a row so that it matches the [WHERE](#) clause of the query.

This occurrence is known as a phantom read. It is harder to guard against than a **non-repeatable read**, because locking all the rows from the first query result set does not prevent the changes that cause the phantom to appear.

Among different **isolation levels**, phantom reads are prevented by the **serializable read** level, and allowed by the **repeatable read**, **consistent read**, and **read uncommitted** levels.

See Also [consistent read](#), [isolation level](#), [non-repeatable read](#), [READ UNCOMMITTED](#), [REPEATABLE READ](#), [SERIALIZABLE](#), [transaction](#).

physical

A type of operation that involves hardware-related aspects such as disk blocks, memory pages, files, bits, disk reads, and so on. Typically, physical aspects are important during expert-level performance tuning and problem diagnosis. Contrast with **logical**.

See Also [logical](#), [physical backup](#).

physical backup

A **backup** that copies the actual data files. For example, the [mysqlbackup](#) command of the **MySQL Enterprise Backup** product produces a physical backup, because its output contains data files that can be used directly by the [mysqld](#) server, resulting in a faster **restore** operation. Contrast with **logical backup**.

See Also [backup](#), [logical backup](#), [MySQL Enterprise Backup](#), [restore](#).

PITR

Acronym for **point-in-time recovery**.

See Also [point-in-time recovery](#).

plan stability

A property of a **query execution plan**, where the optimizer makes the same choices each time for a given **query**, so that performance is consistent and predictable.

See Also [query](#), [query execution plan](#).

plugin

In MySQL 5.1 and earlier, a separately installable form of the **InnoDB** storage engine that includes features and performance enhancements not included in the **built-in** InnoDB for those releases.

For MySQL 5.5 and higher, the MySQL distribution includes the very latest InnoDB features and performance enhancements, known as InnoDB 1.1, and there is no longer a separate InnoDB Plugin.

This distinction is important mainly in MySQL 5.1, where a feature or bug fix might apply to the InnoDB Plugin but not the built-in InnoDB, or vice versa.

See Also [built-in](#), [InnoDB](#).

point-in-time recovery

The process of restoring a **backup** to recreate the state of the database at a specific date and time. Commonly abbreviated **PITR**. Because it is unlikely that the specified time corresponds exactly to the time of a backup, this technique usually requires a combination of a **physical backup** and a **logical backup**. For example, with the **MySQL Enterprise Backup** product, you restore the last backup that you took before the specified point in time, then replay changes from the **binary log** between the time of the backup and the PITR time.

See Also [backup](#), [logical backup](#), [MySQL Enterprise Backup](#), [physical backup](#), [PITR](#).

prefix

See [index prefix](#).

prepared backup

A set of backup files, produced by the **MySQL Enterprise Backup** product, after all the stages of applying **binary logs** and **incremental backups** are finished. The resulting files are ready to be **restored**. Prior to the apply steps, the files are known as a **raw backup**.

See Also [binary log](#), [hot backup](#), [incremental backup](#), [MySQL Enterprise Backup](#), [raw backup](#), [restore](#).

primary key

A set of columns -- and by implication, the index based on this set of columns -- that can uniquely identify every row in a table. As such, it must be a unique index that does not contain any [NULL](#) values.

InnoDB requires that every table has such an index (also called the **clustered index** or **cluster index**), and organizes the table storage based on the column values of the primary key.

When choosing primary key values, consider using arbitrary values (a **synthetic key**) rather than relying on values derived from some other source (a **natural key**).

See Also [clustered index](#), [index](#), [natural key](#), [synthetic key](#).

process

An instance of an executing program. The operating system switches between multiple running processes, allowing for a certain degree of **concurrency**. On most operating systems, processes can contain multiple **threads** of execution that share resources. Context-switching between threads is faster than the equivalent switching between processes.

See Also [concurrency](#), [thread](#).

pseudo-record

An artificial record in an index, used for **locking** key values or ranges that do not currently exist.

See Also [infimum record](#), [locking](#), [supremum record](#).

Pthreads

The POSIX threads standard, which defines an API for threading and locking operations on UNIX and Linux systems. On UNIX and Linux systems, InnoDB uses this implementation for **mutexes**.

See Also [mutex](#).

purge

A type of garbage collection performed by a separate thread, running on a periodic schedule. The purge includes these actions: removing obsolete values from indexes; physically removing rows that were marked for deletion by previous `DELETE` statements.

See Also [crash recovery](#), [delete](#), [doublewrite buffer](#).

purge buffering

The technique of storing index changes due to `DELETE` operations in the **insert buffer** rather than writing them immediately, so that the physical writes can be performed to minimize random I/O. (Because delete operations are a two-step process, this operation buffers the write that normally purges an index record that was previously marked for deletion.) It is one of the types of **change buffering**; the others are **insert buffering** and **delete buffering**.

See Also [change buffer](#), [change buffering](#), [delete buffering](#), [insert buffer](#), [insert buffering](#).

purge lag

Another name for the **InnoDB history list**. Related to the `innodb_max_purge_lag` configuration option.

See Also [history list](#), [purge](#).

purge thread

A **thread** within the InnoDB process that is dedicated to performing the periodic **purge** operation. In MySQL 5.6 and higher, multiple purge threads are enabled by the `innodb_purge_threads` configuration option.

See Also [purge](#), [thread](#).

Q

query

In **SQL**, an operation that reads information from one or more **tables**. Depending on the organization of data and the parameters of the query, the lookup might be optimized by consulting an **index**. If multiple tables are involved, the query is known as a **join**.

For historical reasons, sometimes discussions of internal processing for statements use "query" in a broader sense, including other types of MySQL statements such as **DDL** and **DML** statements.

See Also [DDL](#), [DML](#), [index](#), [join](#), [SQL](#), [table](#).

query execution plan

The set of decisions made by the optimizer about how to perform a **query** most efficiently, including which **index** or indexes to use, and the order in which to **join** tables. **Plan stability** involves the same choices being made consistently for a given query.

See Also [index](#), [join](#), [plan stability](#), [query](#).

query log

See [general query log](#).

quiesce

To reduce the amount of database activity, often in preparation for an operation such as an [ALTER TABLE](#), a **backup**, or a **shutdown**. Might or might not involve doing as much **flushing** as possible, so that **InnoDB** does not continue doing background I/O.

In MySQL 5.6 and higher, the syntax [FLUSH TABLES ... FOR EXPORT](#) writes some data to disk for **InnoDB** tables that make it simpler to back up those tables by copying the data files.

See Also [backup](#), [flush](#), [InnoDB](#), [shutdown](#).

R

RAID

Acronym for "Redundant Array of Inexpensive Drives". Spreading I/O operations across multiple drives enables greater **concurrency** at the hardware level, and improves the efficiency of low-level write operations that otherwise would be performed in sequence.

See Also [concurrency](#).

random dive

A technique for quickly estimating the number of different values in a column (the column's cardinality). **InnoDB** samples pages at random from the index and uses that data to estimate the number of different values. This operation occurs when each table is first opened.

Originally, the number of sampled pages was fixed at 8; now, it is determined by the setting of the [innodb_stats_sample_pages](#) parameter.

The way the random pages are picked depends on the setting of the [innodb_use_legacy_cardinality_algorithm](#) parameter. The default setting (OFF) has better randomness than in older releases.

See Also [cardinality](#).

raw backup

The initial set of backup files produced by the **MySQL Enterprise Backup** product, before the changes reflected in the **binary log** and any **incremental backups** are applied. At this stage, the files are not ready to **restore**. After these changes are applied, the files are known as a **prepared backup**.

See Also [binary log](#), [hot backup](#), [ibbackup_logfile](#), [incremental backup](#), [MySQL Enterprise Backup](#), [prepared backup](#), [restore](#).

READ COMMITTED

An **isolation level** that uses a **locking** strategy that relaxes some of the protection between **transactions**, in the interest of performance. Transactions cannot see uncommitted data from other transactions, but they can see data that is committed by another transaction after the current transaction started. Thus, a transaction never sees any bad data, but the data that it does see may depend to some extent on the timing of other transactions.

When a transaction with this isolation level performs [UPDATE ... WHERE](#) or [DELETE ... WHERE](#) operations, other transactions might have to wait. The transaction can perform [SELECT ... FOR UPDATE](#), and [LOCK IN SHARE MODE](#) operations without making other transactions wait.

See Also [ACID](#), [isolation level](#), [locking](#), [REPEATABLE READ](#), [SERIALIZABLE](#), [transaction](#).

READ UNCOMMITTED

The **isolation level** that provides the least amount of protection between transactions. Queries employ a **locking** strategy that allows them to proceed in situations where they would normally wait for another transaction. However, this extra performance comes at the cost of less reliable results, including data that has been changed by other transactions and not committed yet (known as **dirty read**). Use this isolation level only with great caution,

and be aware that the results might not be consistent or reproducible, depending on what other transactions are doing at the same time. Typically, transactions with this isolation level do only queries, not insert, update, or delete operations.

See Also [ACID](#), [dirty read](#), [isolation level](#), [locking](#), [transaction](#).

read view

An internal snapshot used by the **MVCC** mechanism of InnoDB. Certain **transactions**, depending on their **isolation level**, see the data values as they were at the time the transaction (or in some cases, the statement) started. Isolation levels that use a read view are **REPEATABLE READ**, **READ COMMITTED**, and **READ UNCOMMITTED**.

See Also [isolation level](#), [MVCC](#), [READ COMMITTED](#), [READ UNCOMMITTED](#), [REPEATABLE READ](#), [transaction](#).

read-ahead

A type of I/O request that prefetches a group of **pages** (an entire **extent**) into the **buffer pool** asynchronously, in anticipation that these pages will be needed soon. The linear read-ahead technique prefetches all the pages of one extent based on access patterns for pages in the preceding extent, and is part of all MySQL versions starting with the InnoDB Plugin for MySQL 5.1. The random read-ahead technique prefetches all the pages for an extent once a certain number of pages from the same extent are in the buffer pool. Random read-ahead is not part of MySQL 5.5, but is re-introduced in MySQL 5.6 under the control of the `innodb_random_read_ahead` configuration option.

See Also [buffer pool](#), [extent](#), [page](#).

read-only transaction

A type of transaction that can be optimized for InnoDB tables by eliminating some of the bookkeeping involved with creating a **read view** for each transaction. Can only perform **non-locking read** queries. It can be started explicitly with the syntax `START TRANSACTION READ ONLY`, or automatically under certain conditions. See [Optimizations for Read-Only Transactions](#) for details.

See Also [non-locking read](#), [read view](#), [transaction](#).

record lock

A **lock** on an index record. For example, `SELECT c1 FOR UPDATE FROM t WHERE c1 = 10;` prevents any other transaction from inserting, updating, or deleting rows where the value of `t.c1` is 10. Contrast with **gap lock** and **next-key lock**.

See Also [gap lock](#), [lock](#), [next-key lock](#).

redo

The data, in units of records, recorded in the **redo log** when **DML** statements make changes to InnoDB tables. It is used during **crash recovery** to correct data written by incomplete **transactions**. The ever-increasing **LSN** value represents the cumulative amount of redo data that has passed through the redo log.

See Also [crash recovery](#), [DML](#), [LSN](#), [redo log](#), [transaction](#).

redo log

A disk-based data structure used during **crash recovery**, to correct data written by incomplete **transactions**. During normal operation, it encodes requests to change InnoDB table data, which result from SQL statements or low-level API calls through NoSQL interfaces. Modifications that did not finish updating the **data files** before an unexpected **shutdown** are replayed automatically.

The redo log is physically represented as a set of files, typically named `ib_logfile0` and `ib_logfile1`. The data in the redo log is encoded in terms of records affected; this data is collectively referred to as **redo**. The passage of data through the redo logs is represented by the ever-increasing **LSN** value. The original 4GB limit on maximum size for the redo log is raised to 512GB in MySQL 5.6.3.

The disk layout of the redo log is influenced by the configuration options `innodb_log_file_size`, `innodb_log_group_home_dir`, and (rarely) `innodb_log_files_in_group`. The performance of redo

log operations is also affected by the **log buffer**, which is controlled by the `innodb_log_buffer_size` configuration option.

See Also [crash recovery](#), [data files](#), [ib_logfile](#), [log buffer](#), [LSN](#), [redo](#), [shutdown](#), [transaction](#).

redundant row format

The oldest [InnoDB](#) row format, available for tables using the **Antelope file format**. Prior to MySQL 5.0.3, it was the only row format available in [InnoDB](#). In MySQL 5.0.3 and later, the default is **compact row format**. You can still specify redundant row format for compatibility with older [InnoDB](#) tables.

For additional information about [InnoDB REDUNDANT](#) row format, see [Section 5.2, “COMPACT and REDUNDANT Row Formats”](#).

See Also [Antelope](#), [compact row format](#), [file format](#), [row format](#).

referential integrity

The technique of maintaining data always in a consistent format, part of the **ACID** philosophy. In particular, data in different tables is kept consistent through the use of **foreign key constraints**, which can prevent changes from happening or automatically propagate those changes to all related tables. Related mechanisms include the **unique constraint**, which prevents duplicate values from being inserted by mistake, and the **NOT NULL constraint**, which prevents blank values from being inserted by mistake.

See Also [ACID](#), [FOREIGN KEY constraint](#), [NOT NULL constraint](#), [unique constraint](#).

relational

An important aspect of modern database systems. The database server encodes and enforces relationships such as one-to-one, one-to-many, many-to-one, and uniqueness. For example, a person might have zero, one, or many phone numbers in an address database; a single phone number might be associated with several family members. In a financial database, a person might be required to have exactly one taxpayer ID, and any taxpayer ID could only be associated with one person.

The database server can use these relationships to prevent bad data from being inserted, and to find efficient ways to look up information. For example, if a value is declared to be unique, the server can stop searching as soon as the first match is found, and it can reject attempts to insert a second copy of the same value.

At the database level, these relationships are expressed through SQL features such as **columns** within a table, **unique** and **NOT NULL constraints**, **foreign keys**, and different kinds of join operations. Complex relationships typically involve data split between more than one table. Often, the data is **normalized**, so that duplicate values in one-to-many relationships are stored only once.

In a mathematical context, the relations within a database are derived from set theory. For example, the **OR** and **AND** operators of a **WHERE** clause represent the notions of union and intersection.

See Also [ACID](#), [constraint](#), [foreign key](#), [normalized](#).

relevance

In the **full-text search** feature, a number signifying the similarity between the search string and the data in the **FULLTEXT index**. For example, when you search for a single word, that word is typically more relevant for a row where it occurs several times in the text than a row where it appears only once.

See Also [full-text search](#), [FULLTEXT index](#).

REPEATABLE READ

The default **isolation level** for [InnoDB](#). It prevents any rows that are queried from being changed by other transactions, thus blocking **non-repeatable reads** but not **phantom** reads. It uses a moderately strict **locking** strategy so that all queries within a transaction see data from the same snapshot, that is, the data as it was at the time the transaction started.

When a transaction with this isolation level performs `UPDATE ... WHERE`, `DELETE ... WHERE`, `SELECT ... FOR UPDATE`, and `LOCK IN SHARE MODE` operations, other transactions might have to wait.

See Also [ACID](#), [consistent read](#), [isolation level](#), [locking](#), [phantom](#), [SERIALIZABLE](#), [transaction](#).

replication

The practice of sending changes from a **master database**, to one or more **slave databases**, so that all databases have the same data. This technique has a wide range of uses, such as load-balancing for better scalability, disaster recovery, and testing software upgrades and configuration changes. The changes can be sent between the database by methods called **row-based replication** and **statement-based replication**.

See Also [row-based replication](#), [statement-based replication](#).

restore

The process of putting a set of backup files from the **MySQL Enterprise Backup** product in place for use by MySQL. This operation can be performed to fix a corrupted database, to return to some earlier point in time, or (in a **replication** context) to set up a new **slave database**. In the **MySQL Enterprise Backup** product, this operation is performed by the `copy-back` option of the `mysqlbackup` command.

See Also [hot backup](#), [MySQL Enterprise Backup](#), [mysqlbackup command](#), [prepared backup](#), [replication](#).

rollback

A **SQL** statement that ends a **transaction**, undoing any changes made by the transaction. It is the opposite of **commit**, which makes permanent any changes made in the transaction.

By default, MySQL uses the **autocommit** setting, which automatically issues a commit following each SQL statement. You must change this setting before you can use the rollback technique.

See Also [ACID](#), [commit](#), [transaction](#).

rollback segment

The storage area containing the **undo log**, part of the **system tablespace**.

See Also [system tablespace](#), [undo log](#).

row

The logical data structure defined by a set of **columns**. A set of rows makes up a **table**. Within InnoDB **data files**, each **page** can contain one or more rows.

Although InnoDB uses the term **row format** for consistency with MySQL syntax, the row format is a property of each table and applies to all rows in that table.

See Also [column](#), [data files](#), [page](#), [row format](#), [table](#).

row format

The disk storage format for a **row** from an InnoDB **table**. As InnoDB gains new capabilities such as compression, new row formats are introduced to support the resulting improvements in storage efficiency and performance.

Each table has its own row format, specified through the `ROW_FORMAT` option. To see the row format for each InnoDB table, issue the command `SHOW TABLE STATUS`. Because all the tables in the system tablespace share the same row format, to take advantage of other row formats typically requires setting the `innodb_file_per_table` option, so that each table is stored in a separate tablespace.

See Also [compact row format](#), [compressed row format](#), [dynamic row format](#), [fixed row format](#), [redundant row format](#), [row](#), [table](#).

row lock

A **lock** that prevents a row from being accessed in an incompatible way by another **transaction**. Other rows in the same table can be freely written to by other transactions. This is the type of **locking** done by **DML** operations on **InnoDB** tables.

Contrast with **table locks** used by MyISAM, or during **DDL** operations on InnoDB tables that cannot be done with **online DDL**; those locks block concurrent access to the table.

See Also [DDL](#), [DML](#), [InnoDB](#), [lock](#), [locking](#), [online DDL](#), [table lock](#), [transaction](#).

row-based replication

A form of **replication** where events are propagated from the **master** server specifying how to change individual rows on the **slave** server. It is safe to use for all settings of the `innodb_autoinc_lock_mode` option.

See Also [auto-increment locking](#), [innodb_autoinc_lock_mode](#), [master server](#), [replication](#), [slave server](#), [statement-based replication](#).

row-level locking

The **locking** mechanism used for **InnoDB** tables, relying on **row locks** rather than **table locks**. Multiple **transactions** can modify the same table concurrently. Only if two transactions try to modify the same row does one of the transactions wait for the other to complete (and release its row locks).

See Also [InnoDB](#), [locking](#), [row lock](#), [table lock](#), [transaction](#).

rw-lock

The low-level object that InnoDB uses to represent and enforce shared-access **locks** to internal in-memory data structures. Once the lock is acquired, any other process, thread, and so on can read the data structure, but no one else can write to it. Contrast with **mutexes**, which enforce exclusive access. Mutexes and rw-locks are known collectively as **latches**.

See Also [latch](#), [lock](#), [mutex](#), [Performance Schema](#).

S

savepoint

Savepoints help to implement nested **transactions**. They can be used to provide scope to operations on tables that are part of a larger transaction. For example, scheduling a trip in a reservation system might involve booking several different flights; if a desired flight is unavailable, you might **roll back** the changes involved in booking that one leg, without rolling back the earlier flights that were successfully booked.

See Also [rollback](#), [transaction](#).

scalability

The ability to add more work and issue more simultaneous requests to a system, without a sudden drop in performance due to exceeding the limits of system capacity. Software architecture, hardware configuration, application coding, and type of workload all play a part in scalability. When the system reaches its maximum capacity, popular techniques for increasing scalability are **scale up** (increasing the capacity of existing hardware or software) and **scale out** (adding new servers and more instances of MySQL). Often paired with **availability** as critical aspects of a large-scale deployment.

See Also [availability](#), [scale out](#), [scale up](#).

scale out

A technique for increasing **scalability** by adding new servers and more instances of MySQL. For example, setting up replication, MySQL Cluster, connection pooling, or other features that spread work across a group of servers. Contrast with **scale up**.

See Also [scalability](#), [scale up](#).

scale up

A technique for increasing **scalability** by increasing the capacity of existing hardware or software.

For example, increasing the memory on a server and adjusting memory-related parameters such as `innodb_buffer_pool_size` and `innodb_buffer_pool_instances`. Contrast with **scale out**.

See Also [scalability](#), [scale out](#).

schema

Conceptually, a schema is a set of interrelated database objects, such as tables, table columns, data types of the columns, indexes, foreign keys, and so on. These objects are connected through SQL syntax, because the columns make up the tables, the foreign keys refer to tables and columns, and so on. Ideally, they are also connected logically, working together as part of a unified application or flexible framework. For example, the

information_schema and **performance_schema** databases use "schema" in their names to emphasize the close relationships between the tables and columns they contain.

In MySQL, physically, a **schema** is synonymous with a **database**. You can substitute the keyword `SCHEMA` instead of `DATABASE` in MySQL SQL syntax, for example using `CREATE SCHEMA` instead of `CREATE DATABASE`.

Some other database products draw a distinction. For example, in the Oracle Database product, a **schema** represents only a part of a database: the tables and other objects owned by a single user.
See Also [database](#), [ib-file set](#), [INFORMATION_SCHEMA](#), [Performance Schema](#).

search index

In MySQL, **full-text search** queries use a special kind of index, the **FULLTEXT** index. In MySQL 5.6.4 and up, [InnoDB](#) and [MyISAM](#) tables both support **FULLTEXT** indexes; formerly, these indexes were only available for [MyISAM](#) tables.

See Also [full-text search](#), [FULLTEXT index](#).

secondary index

A type of InnoDB **index** that represents a subset of table columns. An InnoDB table can have zero, one, or many secondary indexes. (Contrast with the **clustered index**, which is required for each InnoDB table, and stores the data for all the table columns.)

A secondary index can be used to satisfy queries that only require values from the indexed columns. For more complex queries, it can be used to identify the relevant rows in the table, which are then retrieved through lookups using the clustered index.

Creating and dropping secondary indexes has traditionally involved significant overhead from copying all the data in the InnoDB table. The **fast index creation** feature of the InnoDB Plugin makes both `CREATE INDEX` and `DROP INDEX` statements much faster for InnoDB secondary indexes.

See Also [clustered index](#), [Fast Index Creation](#), [index](#).

segment

A division within an InnoDB **tablespace**. If a tablespace is analogous to a directory, the segments are analogous to files within that directory. A segment can grow. New segments can be created.

For example, within a **file-per-table** tablespace, the table data is in one segment and each associated index is in its own segment. The **system tablespace** contains many different segments, because it can hold many tables and their associated indexes. The system tablespace also includes up to 128 **rollback segments** making up the **undo log**.

Segments grow and shrink as data is inserted and deleted. When a segment needs more room, it is extended by one **extent** (1 megabyte) at a time. Similarly, a segment releases one extent's worth of space when all the data in that extent is no longer needed.

See Also [extent](#), [file-per-table](#), [rollback segment](#), [system tablespace](#), [tablespace](#), [undo log](#).

selectivity

A property of data distribution, the number of distinct values in a column (its **cardinality**) divided by the number of records in the table. High selectivity means that the column values are relatively unique, and can be retrieved efficiently through an index. If you (or the query optimizer) can predict that a test in a `WHERE` clause only matches a small number (or proportion) of rows in a table, the overall **query** tends to be efficient if it evaluates that test first, using an index.

See Also [cardinality](#), [query](#).

semi-consistent read

A type of read operation used for `UPDATE` statements, that is a combination of **read committed** and **consistent read**. When an `UPDATE` statement examines a row that is already locked, InnoDB returns the latest committed

version to MySQL so that MySQL can determine whether the row matches the [WHERE](#) condition of the [UPDATE](#). If the row matches (must be updated), MySQL reads the row again, and this time InnoDB either locks it or waits for a lock on it. This type of read operation can only happen when the transaction has the read committed **isolation level**, or when the [innodb_locks_unsafe_for_binlog](#) option is enabled.
See Also [consistent read](#), [isolation level](#), [READ COMMITTED](#).

SERIALIZABLE

The **isolation level** that uses the most conservative locking strategy, to prevent any other transactions from inserting or changing data that was read by this transaction, until it is finished. This way, the same query can be run over and over within a transaction, and be certain to retrieve the same set of results each time. Any attempt to change data that was committed by another transaction since the start of the current transaction, cause the current transaction to wait.

This is the default isolation level specified by the SQL standard. In practice, this degree of strictness is rarely needed, so the default isolation level for InnoDB is the next most strict, **repeatable read**.
See Also [ACID](#), [consistent read](#), [isolation level](#), [locking](#), [REPEATABLE READ](#), [transaction](#).

server

A type of program that runs continuously, waiting to receive and act upon requests from another program (the client). Because often an entire computer is dedicated to running one or more server programs (such as a database server, a web server, an application server, or some combination of these), the term **server** can also refer to the computer that runs the server software.
See Also [client](#), [mysqld](#).

shared lock

A kind of **lock** that allows other **transactions** to read the locked object, and to also acquire other shared locks on it, but not to write to it. The opposite of **exclusive lock**.
See Also [exclusive lock](#), [lock](#), [transaction](#).

shared tablespace

Another way of referring to the **system tablespace**.
See Also [system tablespace](#).

sharp checkpoint

The process of **flushing** to disk all **dirty** buffer pool pages whose redo entries are contained in certain portion of the **redo log**. Occurs before InnoDB reuses a portion of a log file; the log files are used in a circular fashion. Typically occurs with write-intensive **workloads**.
See Also [dirty page](#), [flush](#), [redo log](#), [workload](#).

shutdown

The process of stopping the MySQL server. By default, this process does cleanup operations for **InnoDB** tables, so it can **slow** to shut down, but fast to start up later. If you skip the cleanup operations, it is **fast** to shut down but must do the cleanup during the next restart.

The shutdown mode is controlled by the [innodb_fast_shutdown](#) option.
See Also [fast shutdown](#), [InnoDB](#), [slow shutdown](#), [startup](#).

slave server

Frequently shortened to "slave". A database **server** machine in a **replication** scenario that receives changes from another server (the **master**) and applies those same changes. Thus it maintains the same contents as the master, although it might lag somewhat behind.

In MySQL, slave servers are commonly used in disaster recovery, to take the place of a master servers that fails. They are also commonly used for testing software upgrades and new settings, to ensure that database configuration changes do not cause problems with performance or reliability.

Slave servers typically have high workloads, because they process all the **DML** (write) operations relayed from the master, as well as user queries. To ensure that slave servers can apply changes from the master fast enough, they frequently have fast I/O devices and sufficient CPU and memory to run multiple database instances on the same slave server. For example, the master server might use hard drive storage while the slave servers use **SSDs**.

See Also [DML](#), [replication](#), [server](#), [SSD](#).

slow query log

A type of **log** used for performance tuning of SQL statements processed by the MySQL server. The log information is stored in a file. You must enable this feature to use it. You control which categories of "slow" SQL statements are logged. For more information, see [The Slow Query Log](#).

See Also [general query log](#), [log](#).

slow shutdown

A type of shutdown that does additional [InnoDB](#) flushing operations before completing. Also known as a **clean shutdown**. Specified by the configuration parameter `innodb_fast_shutdown=0` or the command `SET GLOBAL innodb_fast_shutdown=0;`. Although the shutdown itself can take longer, that time will be saved on the subsequent startup.

See Also [clean shutdown](#), [fast shutdown](#), [shutdown](#).

snapshot

A representation of data at a particular time, which remains the same even as changes are **committed** by other **transactions**. Used by certain **isolation levels** to allow **consistent reads**.

See Also [commit](#), [consistent read](#), [isolation level](#), [transaction](#).

space ID

An identifier used to uniquely identify an [InnoDB](#) **tablespace** within a MySQL instance. The space ID for the **system tablespace** is always zero; this same ID applies to all tables within the system tablespace. Each tablespace file created in **file-per-table** mode also has its own space ID.

Prior to MySQL 5.6, this hardcoded value presented difficulties in moving [InnoDB](#) tablespace files between MySQL instances. Starting in MySQL 5.6, you can copy tablespace files between instances by using the **transportable tablespace** feature involving the statements `FLUSH TABLES ... FOR EXPORT`, `ALTER TABLE ... DISCARD TABLESPACE`, and `ALTER TABLE ... IMPORT TABLESPACE`. The information needed to adjust the space ID is conveyed in the **.cfg file** which you copy along with the tablespace. See [Improved Tablespace Management](#) for details.

See Also [.cfg file](#), [file-per-table](#), [.ibd file](#), [system tablespace](#), [tablespace](#), [transportable tablespace](#).

spin

A type of **wait** operation that continuously tests whether a resource becomes available. This technique is used for resources that are typically held only for brief periods, where it is more efficient to wait in a "busy loop" than to put the thread to sleep and perform a context switch. If the resource does not become available within a short time, the spin loop ceases and another wait technique is used.

See Also [latch](#), [lock](#), [mutex](#), [wait](#).

SQL

The Structured Query Language that is standard for performing database operations. Often divided into the categories **DDL**, **DML**, and **queries**. MySQL includes some additional statement categories such as **replication**. See [Language Structure](#) for the building blocks of SQL syntax, [Data Types](#) for the data types to use for MySQL table columns, [SQL Statement Syntax](#) for details about SQL statements and their associated categories, and [Functions and Operators](#) for standard and MySQL-specific functions to use in queries.

See Also [DDL](#), [DML](#), [query](#), [replication](#).

SSD

Acronym for "solid-state drive". A type of storage device with different performance characteristics than a traditional hard disk drive (**HDD**): smaller storage capacity, faster for random reads, no moving parts, and with

a number of considerations affecting write performance. Its performance characteristics can influence the throughput of a **disk-bound** workload.
See Also [disk-bound](#), [SSD](#).

startup

The process of starting the MySQL server. Typically done by one of the programs listed in [MySQL Server and Server-Startup Programs](#). The opposite of **shutdown**.
See Also [shutdown](#).

statement-based replication

A form of **replication** where SQL statements are sent from the **master** server and replayed on the **slave** server. It requires some care with the setting for the `innodb_autoinc_lock_mode` option, to avoid potential timing problems with **auto-increment locking**.
See Also [auto-increment locking](#), [innodb_autoinc_lock_mode](#), [master server](#), [replication](#), [row-based replication](#), [slave server](#).

statistics

Estimated values relating to each **InnoDB table** and **index**, used to construct an efficient **query execution plan**. The main values are the **cardinality** (number of distinct values) and the total number of table rows or index entries. The statistics for the table represent the data in its **primary key** index. The statistics for a **secondary index** represent the rows covered by that index.

The values are estimated rather than counted precisely because at any moment, different **transactions** can be inserting and deleting rows from the same table. To keep the values from being recalculated frequently, you can enable **persistent statistics**, where the values are stored in **InnoDB** system tables, and refreshed only when you issue an `ANALYZE TABLE` statement.

You can control how **NULL** values are treated when calculating statistics through the `innodb_stats_method` configuration option.

Other types of statistics are available for database objects and database activity through the **INFORMATION_SCHEMA** and **PERFORMANCE_SCHEMA** tables.

See Also [cardinality](#), [index](#), [INFORMATION_SCHEMA](#), [NULL](#), [Performance Schema](#), [persistent statistics](#), [primary key](#), [query execution plan](#), [secondary index](#), [table](#), [transaction](#).

stemming

The ability to search for different variations of a word based on a common root word, such as singular and plural, or past, present, and future verb tense. This feature is currently supported in MyISAM **full-text search** feature but not in **FULLTEXT indexes** for InnoDB tables.
See Also [full-text search](#), [FULLTEXT index](#).

stopword

In a **FULLTEXT index**, a word that is considered common or trivial enough that it is omitted from the **search index** and ignored in search queries. Different configuration settings control stopword processing for **InnoDB** and **MyISAM** tables. See [Full-Text Stopwords](#) for details.
See Also [FULLTEXT index](#), [search index](#).

storage engine

A component of the MySQL database that performs the low-level work of storing, updating, and querying data. In MySQL 5.5 and higher, **InnoDB** is the default storage engine for new tables, superceding MyISAM. Different storage engines are designed with different tradeoffs between factors such as memory usage versus disk usage, read speed versus write speed, and speed versus robustness. Each storage engine manages specific tables, so we refer to **InnoDB** tables, **MyISAM** tables, and so on.

The **MySQL Enterprise Backup** product is optimized for backing up InnoDB tables. It can also back up tables handled by MyISAM and other storage engines.

See Also [InnoDB](#), [MySQL Enterprise Backup](#), [table type](#).

strict mode

The general name for the setting controlled by the `innodb_strict_mode` option. Turning on this setting causes certain conditions that are normally treated as warnings, to be considered errors. For example, certain invalid combinations of options related to **file format** and **row format**, that normally produce a warning and continue with default values, now cause the `CREATE TABLE` operation to fail.

MySQL also has something called strict mode.

See Also [file format](#), [innodb_strict_mode](#), [row format](#).

sublist

Within the list structure that represents the buffer pool, pages that are relatively old and relatively new are represented by different portions of the list. A set of parameters control the size of these portions and the dividing point between the new and old pages.

See Also [buffer pool](#), [eviction](#), [list](#), [LRU](#).

supremum record

A **pseudo-record** in an index, representing the **gap** above the largest value in that index. If a transaction has a statement such as `SELECT ... FOR UPDATE ... WHERE col > 10;`, and the largest value in the column is 20, it is a lock on the supremum record that prevents other transactions from inserting even larger values such as 50, 100, and so on.

See Also [gap](#), [infimum record](#), [pseudo-record](#).

surrogate key

Synonym name for **synthetic key**.

See Also [synthetic key](#).

synthetic key

A indexed column, typically a **primary key**, where the values are assigned arbitrarily. Often done using an **auto-increment** column. By treating the value as completely arbitrary, you can avoid overly restrictive rules and faulty application assumptions. For example, a numeric sequence representing employee numbers might have a gap if an employee was approved for hiring but never actually joined. Or employee number 100 might have a later hiring date than employee number 500, if they left the company and later rejoined. Numeric values also produce shorter values of predictable length. For example, storing numeric codes meaning "Road", "Boulevard", "Expressway", and so on is more space-efficient than repeating those strings over and over.

Also known as a **surrogate key**. Contrast with **natural key**.

See Also [auto-increment](#), [natural key](#), [primary key](#), [surrogate key](#).

system tablespace

A small set of data files (the **ibdata** files) containing the metadata for InnoDB-related objects (the **data dictionary**), and the storage areas for the **undo log**, the **change buffer**, and the **doublewrite buffer**. Depending on the setting of the `innodb_file_per_table`, when tables are created, it might also contain table and index data for some or all InnoDB tables. The data and metadata in the system tablespace apply to all the **databases** in a MySQL **instance**.

Prior to MySQL 5.6.7, the default was to keep all InnoDB tables and indexes inside the system tablespace, often causing this file to become very large. Because the system tablespace never shrinks, storage problems could arise if large amounts of temporary data were loaded and then deleted. In MySQL 5.6.7 and higher, the default is **file-per-table** mode, where each table and its associated indexes are stored in a separate **.ibd file**. This new default makes it easier to use InnoDB features that rely on the **Barracuda** file format, such as table **compression** and the **DYNAMIC** row format.

In MySQL 5.6 and higher, setting a value for the `innodb_undo_tablespaces` option splits the **undo log** into one or more separate tablespace files. These files are still considered part of the system tablespace.

Keeping all table data in the system tablespace or in separate `.ibd` files has implications for storage management in general. The **MySQL Enterprise Backup** product might back up a small set of large files, or many smaller files. On systems with thousands of tables, the filesystem operations to process thousands of `.ibd` files can cause bottlenecks.

See Also [Barracuda](#), [change buffer](#), [compression](#), [data dictionary](#), [database](#), [doublewrite buffer](#), [dynamic row format](#), [file-per-table](#), [.ibd file](#), [ibdata file](#), [innodb_file_per_table](#), [instance](#), [MySQL Enterprise Backup](#), [tablespace](#), [undo log](#).

T

.TRG file

A file containing **trigger** parameters. Files with this extension are always included in backups produced by the `mysqlbackup` command of the **MySQL Enterprise Backup** product.

See Also [MySQL Enterprise Backup](#), [mysqlbackup command](#), [.TRN file](#).

.TRN file

A file containing trigger namespace information. Files with this extension are always included in backups produced by the `mysqlbackup` command of the **MySQL Enterprise Backup** product.

See Also [MySQL Enterprise Backup](#), [mysqlbackup command](#), [.TRG file](#).

table

Each MySQL table is associated with a particular **storage engine**. **InnoDB** tables have particular **physical** and **logical** characteristics that affect performance, **scalability**, **backup**, administration, and application development.

In terms of file storage, each InnoDB table is either part of the single big InnoDB **system tablespace**, or in a separate `.ibd` file if the table is created in **file-per-table** mode. The `.ibd` file holds data for the table and all its **indexes**, and is known as a **tablespace**.

InnoDB tables created in file-per-table mode can use the **Barracuda** file format. Barracuda tables can use the **DYNAMIC row format** or the **COMPRESSED row format**. These relatively new settings enable a number of InnoDB features, such as **compression**, **fast index creation**, and **off-page columns**

For backward compatibility with MySQL 5.1 and earlier, InnoDB tables inside the system tablespace must use the **Antelope** file format, which supports the **compact row format** and the **redundant row format**.

The **rows** of an InnoDB table are organized into an index structure known as the **clustered index**, with entries sorted based on the **primary key** columns of the table. Data access is optimized for queries that filter and sort on the primary key columns, and each index contains a copy of the associated primary key columns for each entry. Modifying values for any of the primary key columns is an expensive operation. Thus an important aspect of InnoDB table design is choosing a primary key with columns that are used in the most important queries, and keeping the primary key short, with rarely changing values.

See Also [Antelope](#), [backup](#), [Barracuda](#), [clustered index](#), [compact row format](#), [compressed row format](#), [compression](#), [dynamic row format](#), [Fast Index Creation](#), [file-per-table](#), [.ibd file](#), [index](#), [off-page column](#), [primary key](#), [redundant row format](#), [row](#), [system tablespace](#), [tablespace](#).

table lock

A lock that prevents any other **transaction** from accessing a table. InnoDB makes considerable effort to make such locks unnecessary, by using techniques such as **online DDL**, **row locks** and **consistent reads** for processing **DML** statements and **queries**. You can create such a lock through SQL using the `LOCK TABLE` statement; one of the steps in migrating from other database systems or MySQL storage engines is to remove such statements wherever practical.

See Also [consistent read](#), [DML](#), [lock](#), [locking](#), [online DDL](#), [query](#), [row lock](#), [table](#), [transaction](#).

table scan

See [full table scan](#).

table statistics

See [statistics](#).

table type

Obsolete synonym for **storage engine**. We refer to [InnoDB](#) tables, [MyISAM](#) tables, and so on.
See Also [InnoDB](#), [storage engine](#).

tablespace

A data file that can hold data for one or more InnoDB **tables** and associated **indexes**. The **system tablespace** contains the tables that make up the **data dictionary**, and prior to MySQL 5.6 holds all the other InnoDB tables by default. Turning on the [innodb_file_per_table](#) option, the default in MySQL 5.6 and higher, allows newly created tables to each have their own tablespace, with a separate **data file** for each table.

Using multiple tablespaces, by turning on the [innodb_file_per_table](#) option, is vital to using many MySQL features such as table compression and transportable tablespaces, and managing disk usage. See [Using Per-Table Tablespaces](#) and [Improved Tablespace Management](#) for details.

Tablespaces created by the built-in InnoDB storage engine are upward compatible with the InnoDB Plugin. Tablespaces created by the InnoDB Plugin are downward compatible with the built-in InnoDB storage engine, if they use the **Antelope** file format.

MySQL Cluster also groups its tables into tablespaces. See [MySQL Cluster Disk Data Objects](#) for details.
See Also [Antelope](#), [Barracuda](#), [compressed row format](#), [data dictionary](#), [data files](#), [file-per-table](#), [index](#), [innodb_file_per_table](#), [system tablespace](#), [table](#).

tablespace dictionary

A representation of the **data dictionary** metadata for a table, within the InnoDB **tablespace**. This metadata can be checked against the **.frm file** for consistency when the table is opened, to diagnose errors resulting from out-of-date **.frm** files. This information is present for InnoDB tables that are part of the **system tablespace**, as well as for tables that have their own **.ibd file** because of the **file-per-table** option.

See Also [data dictionary](#), [file-per-table](#), [.frm file](#), [.ibd file](#), [system tablespace](#), [tablespace](#).

temporary table

A table whose data does not need to be truly permanent. For example, temporary tables might be used as storage areas for intermediate results in complicated calculations or transformations; this intermediate data would not need to be recovered after a crash. Database products can take various shortcuts to improve the performance of operations on temporary tables, by being less scrupulous about writing data to disk and other measures to protect the data across restarts.

Sometimes, the data itself is removed automatically at a set time, such as when the transaction ends or when the session ends. With some database products, the table itself is removed automatically too.

See Also [table](#).

temporary tablespace

The tablespace for non-compressed [InnoDB](#) temporary tables and related objects. This tablespace was introduced in MySQL 5.7.1. The configuration file option, [innodb_temp_data_file_path](#), allows users to define a relative path for the temporary data file. If [innodb_temp_data_file_path](#) is not specified, the default behavior is to create a single auto-extending 12MB data file named [ibtmp1](#) in the data directory, alongside [ibdata1](#). The temporary tablespace is recreated on each server start and receives a dynamically generated space-id, which helps avoid conflicts with existing space-ids. The temporary tablespace cannot reside on a raw device. Inability or error creating the temporary table is treated as fatal and server startup will be refused.

The tablespace is removed on normal shutdown or on init abort, which may occur when a user specifies the wrong startup options, for example. The temporary tablespace is not removed when a crash occurs. In this case, the database administrator can remove the tablespace manually or restart the server with the same configuration, which will remove and recreate the temporary tablespace.

See Also [ibtmp file](#).

text collection

The set of columns included in a **FULLTEXT index**.
See Also [FULLTEXT index](#).

thread

A unit of processing that is typically more lightweight than a **process**, allowing for greater **concurrency**.
See Also [concurrency](#), [master thread](#), [process](#), [Pthreads](#).

torn page

An error condition that can occur due to a combination of I/O device configuration and hardware failure. If data is written out in chunks smaller than the InnoDB **page size** (by default, 16KB), a hardware failure while writing could result in only part of a page being stored to disk. The InnoDB **doublewrite buffer** guards against this possibility.
See Also [doublewrite buffer](#).

TPS

Acronym for "**transactions** per second", a unit of measurement sometimes used in benchmarks. Its value depends on the **workload** represented by a particular benchmark test, combined with factors that you control such as the hardware capacity and database configuration.
See Also [transaction](#), [workload](#).

transaction

Transactions are atomic units of work that can be committed or rolled back. When a transaction makes multiple changes to the database, either all the changes succeed when the transaction is committed, or all the changes are undone when the transaction is rolled back.

Database transactions, as implemented by InnoDB, have properties that are collectively known by the acronym **ACID**, for atomicity, consistency, isolation, and durability.
See Also [ACID](#), [commit](#), [isolation level](#), [lock](#), [rollback](#).

transaction ID

An internal field associated with each row. This field is physically changed by INSERT, UPDATE, and DELETE operations to record which transaction has locked the row.
See Also [implicit row lock](#).

transportable tablespace

A feature that allows a **tablespace** to be moved from one instance to another. Traditionally, this has not been possible for InnoDB tablespaces because all table data was part of the **system tablespace**. In MySQL 5.6 and higher, the `FLUSH TABLES ... FOR EXPORT` syntax prepares an InnoDB table for copying to another server; running `ALTER TABLE ... DISCARD TABLESPACE` and `ALTER TABLE ... IMPORT TABLESPACE` on the other server brings the copied data file into the other instance. A separate `.cfg` file, copied along with the `.ibd` file, is used to update the table metadata (for example the **space ID**) as the tablespace is imported. See [Improved Tablespace Management](#) for usage information.
See Also [.ibd file](#), [space ID](#), [system tablespace](#), [tablespace](#).

troubleshooting

Resources for troubleshooting InnoDB reliability and performance issues include: the Information Schema tables.

truncate

A **DDL** operation that removes the entire contents of a table, while leaving the table and related indexes intact. Contrast with **drop**. Although conceptually it has the same result as a `DELETE` statement with no `WHERE` clause, it operates differently behind the scenes: InnoDB creates a new empty table, drops the old table, then renames the new table to take the place of the old one. Because this is a DDL operation, it cannot be **rolled back**.

If the table being truncated contains foreign keys that reference another table, the truncation operation uses a slower method of operation, deleting one row at a time so that corresponding rows in the referenced table can be deleted as needed by any `ON DELETE CASCADE` clause. (MySQL 5.5 and higher do not allow this slower form of truncate, and return an error instead if foreign keys are involved. In this case, use a `DELETE` statement instead.

See Also [DDL](#), [drop](#), [foreign key](#), [rollback](#).

tuple

A technical term designating an ordered set of elements. It is an abstract notion, used in formal discussions of database theory. In the database field, tuples are usually represented by the columns of a table row. They could also be represented by the result sets of queries, for example, queries that retrieved only some columns of a table, or columns from joined tables.

See Also [cursor](#).

two-phase commit

An operation that is part of a distributed **transaction**, under the **XA** specification. (Sometimes abbreviated as 2PC.) When multiple databases participate in the transaction, either all databases **commit** the changes, or all databases **roll back** the changes.

See Also [commit](#), [rollback](#), [transaction](#), [XA](#).

U

undo

Data that is maintained throughout the life of a **transaction**, recording all changes so that they can be undone in case of a **rollback** operation. It is stored in the **undo log**, also known as the **rollback segment**, either within the **system tablespace** or in separate **undo tablespaces**.

See Also [rollback](#), [rollback segment](#), [system tablespace](#), [transaction](#), [undo log](#), [undo tablespace](#).

undo buffer

See [undo log](#).

undo log

A storage area that holds copies of data modified by active **transactions**. If another transaction needs to see the original data (as part of a **consistent read** operation), the unmodified data is retrieved from this storage area.

By default, this area is physically part of the **system tablespace**. In MySQL 5.6 and higher, you can use the [innodb_undo_tablespaces](#) and [innodb_undo_directory](#) configuration options to split it into one or more separate **tablespace** files, the **undo tablespaces**, optionally stored on another storage device such as an **SSD**.

The undo log is split into separate portions, the **insert undo buffer** and the **update undo buffer**. Collectively, these parts are also known as the **rollback segment**, a familiar term for Oracle DBAs.

See Also [consistent read](#), [rollback segment](#), [SSD](#), [system tablespace](#), [transaction](#), [undo tablespace](#).

undo tablespace

One of a set of files containing the **undo log**, when the undo log is separated from the **system tablespace** by the [innodb_undo_tablespaces](#) and [innodb_undo_directory](#) configuration options. Only applies to MySQL 5.6 and higher.

See Also [system tablespace](#), [undo log](#).

unique constraint

A kind of **constraint** that asserts that a column cannot contain any duplicate values. In terms of **relational algebra**, it is used to specify 1-to-1 relationships. For efficiency in checking whether a value can be inserted (that is, the value does not already exist in the column), a unique constraint is supported by an underlying **unique index**.

See Also [constraint](#), [relational](#), [unique index](#).

unique index

An index on a column or set of columns that have a **unique constraint**. Because the index is known not to contain any duplicate values, certain kinds of lookups and count operations are more efficient than in the normal kind of index. Most of the lookups against this type of index are simply to determine if a certain value exists or not.

The number of values in the index is the same as the number of rows in the table, or at least the number of rows with non-null values for the associated columns.

The **insert buffering** optimization does not apply to unique indexes. As a workaround, you can temporarily set `unique_checks=0` while doing a bulk data load into an InnoDB table.

See Also [cardinality](#), [insert buffering](#), [unique constraint](#), [unique key](#).

unique key

The set of columns (one or more) comprising a **unique index**. When you can define a `WHERE` condition that matches exactly one row, and the query can use an associated unique index, the lookup and error handling can be performed very efficiently.

See Also [cardinality](#), [unique constraint](#), [unique index](#).

V

victim

The transaction that is automatically chosen to be rolled back when a **deadlock** is detected. InnoDB rolls back the transaction that has updated the fewest rows.

See Also [deadlock](#), [deadlock detection](#), [innodb_lock_wait_timeout](#).

W

wait

When an operation, such as acquiring a **lock**, **mutex**, or **latch**, cannot be completed immediately, InnoDB pauses and tries again. The mechanism for pausing is elaborate enough that this operation has its own name, the **wait**. Individual threads are paused using a combination of internal InnoDB scheduling, operating system `wait()` calls, and short-duration **spin** loops.

On systems with heavy load and many transactions, you might use the output from the `SHOW INNODB STATUS` command to determine whether threads are spending too much time waiting, and if so, how you can improve **concurrency**.

See Also [concurrency](#), [latch](#), [lock](#), [mutex](#), [spin](#).

warm backup

A **backup** taken while the database is running, but that restricts some database operations during the backup process. For example, tables might become read-only. For busy applications and web sites, you might prefer a **hot backup**.

See Also [backup](#), [cold backup](#), [hot backup](#).

warm up

To run a system under a typical **workload** for some time after startup, so that the **buffer pool** and other memory regions are filled as they would be under normal conditions.

This process happens naturally over time when a MySQL server is restarted or subjected to a new workload. Starting in MySQL 5.6, you can speed up the warmup process by setting the configuration variables `innodb_buffer_pool_dump_at_shutdown=ON` and `innodb_buffer_pool_load_at_startup=ON`, to bring the contents of the buffer pool back into memory after a restart. Typically, you run a workload for some time to warm up the buffer pool before running performance tests, to ensure consistent results across multiple runs; otherwise, performance might be artificially low during the first run.

See Also [buffer pool](#), [workload](#).

Windows

The built-in **InnoDB** storage engine and the InnoDB **Plugin** are supported on all the same Microsoft Windows versions as the MySQL server. The **MySQL Enterprise Backup** product has more comprehensive support for Windows systems than the **InnoDB Hot Backup** product that it supersedes.

See Also [InnoDB](#), [MySQL Enterprise Backup](#), [plugin](#).

workload

The combination and volume of **SQL** and other database operations, performed by a database application during typical or peak usage. You can subject the database to a particular workload during performance testing to identify **bottlenecks**, or during capacity planning.

See Also [bottleneck](#), [disk-bound](#), [disk-bound](#), [SQL](#).

write combining

An optimization technique that reduces write operations when **dirty pages** are **flushed** from the InnoDB **buffer pool**. If a row in a page is updated multiple times, or multiple rows on the same page are updated, all of those changes are stored to the data files in a single write operation rather than one write for each change.

See Also [buffer pool](#), [dirty page](#), [flush](#).

X

XA

A standard interface for coordinating distributed **transactions**, allowing multiple databases to participate in a transaction while maintaining **ACID** compliance. For full details, see [XA Transactions](#).

XA Distributed Transaction support is turned on by default. If you are not using this feature, you can disable the `innodb_support_xa` configuration option, avoiding the performance overhead of an extra fsync for each transaction.

See Also [commit](#), [transaction](#), [two-phase commit](#).

Y

young

A characteristic of a **page** in the [InnoDB buffer pool](#) meaning it has been accessed recently, and so is moved within the buffer pool data structure, so that it will not be **flushed** soon by the **LRU** algorithm. This term is used in some **information schema** column names of tables related to the buffer pool.

See Also [buffer pool](#), [flush](#), [INFORMATION_SCHEMA](#), [LRU](#), [page](#).

Index

Symbols

.ARM file, 95
.ARZ file, 95
.cfg file, 100
.frm file, 111
.ibd file, 116
.ibz file, 116
.isl file, 116
.MRG file, 124
.MYD file, 125
.MYI file, 125
.OPT file, 129
.PAR file, 131
.TRG file, 145
.TRN file, 145

A

ACID, 95
adaptive flushing, 95
adaptive hash index, 42, 51, 95
AHI, 96
AIO, 96
ALTER TABLE
 ROW_FORMAT, 26
Antelope, 96
Antelope file format, 19
application programming interface (API), 96
apply, 96
asynchronous I/O, 96
atomic, 97
atomic instruction, 97
auto-increment, 97
auto-increment locking, 97
autocommit, 97
availability, 97

B

B-tree, 97
background threads
 master, 45, 45
 read, 44
 write, 44
backticks, 98
backup, 98
Barracuda, 98
Barracuda file format, 19
beta, 98
binary log, 98
binlog, 99
blind query expansion, 99

bottleneck, 99
bounce, 99
buddy allocator, 28, 99
buffer, 99
buffer cache, 46
buffer pool, 99
buffer pool instance, 100
built-in, 100
business rules, 100

C

cache, 101
cardinality, 101
change buffer, 101
change buffering, 77, 101
 disabling, 41
checkpoint, 102
checksum, 102
child table, 102
clean page, 102
clean shutdown, 102
client, 102
clustered index, 102
cold backup, 103
column, 103
column index, 103
column prefix, 103
commit, 103
compact row format, 103
compiling, 65
composite index, 103
compressed backup, 104
compressed row format, 104
compression, 104
 algorithms, 16
 application and schema design, 13
 BLOBs, VARCHAR and TEXT, 17
 buffer pool, 18
 compressed page size, 15
 configuration characteristics, 14
 data and indexes, 16
 data characteristics, 13
 enabling for a table, 9
 implementation, 16
 information schema, 27, 27
 innodb_strict_mode, 52
 KEY_BLOCK_SIZE, 15
 log files, 18
 modification log, 16
 monitoring, 15
 overflow pages, 17
 overview, 9
 tuning, 12

- workload characteristics, 14
- compression failure, 104
- concurrency, 105
- configuration file, 105
- configuring, 68
- consistent read, 105
- constraint, 105
- counter, 106
- covering index, 106
- crash, 106
- crash recovery, 106
- CREATE INDEX, 5
- CREATE TABLE
 - KEY_BLOCK_SIZE, 15
 - options for table compression, 9
 - ROW_FORMAT, 26
- CRUD, 106
- cursor, 106

D

- data dictionary, 107
- data directory, 107
- data files, 107
- data warehouse, 107
- database, 107
- DCL, 107
- DDL, 107
- deadlock, 108
- deadlock detection, 108
- delete, 108
- delete buffering, 108
- denormalized, 109
- descending index, 109
- dirty page, 109
- dirty read, 109
- disk-based, 109
- disk-bound, 109, 109
- DML, 109
- document id, 110
- doublewrite buffer, 110
- downgrading, 73
- drop, 110
- DROP INDEX, 5
- dynamic row format, 110

E

- early adopter, 110
- error log, 110
- eviction, 111
- exclusive lock, 111
- extent, 111

F

- Fast Index Creation, 111
 - concurrency, 7
 - crash recovery, 7
 - examples, 5
 - implementation, 6
 - limitations, 8
 - overview, 5
- fast shutdown, 112
- file format, 19, 112
 - Antelope, 17
 - Barracuda, 9
 - downgrading, 24
 - identifying, 23
- file format management
 - downgrading, 73
 - enabling new file formats, 49
- file per table, 50
- file-per-table, 112
- fill factor, 112
- fixed row format, 112
- flush, 112
- flush list, 113
- foreign key, 113
- FOREIGN KEY constraint, 113
- FOREIGN KEY constraints
 - and fast index creation, 8
 - and TRUNCATE_TABLE, 51
- FTS, 113
- full backup, 113
- full table scan, 113
- full-text search, 114
- FULLTEXT index, 114
- fuzzy checkpointing, 114

G

- GA, 114
- gap, 114
- gap lock, 114
- general query log, 114
- global_transaction, 115
- group commit, 44, 115

H

- hash index, 115
- HDD, 115
- heartbeat, 115
- high-water mark, 115
- history list, 115
- hot, 115
- hot backup, 116

I

- ib-file set, 20, 117
- ibbackup_logfile, 117
- ibdata file, 117
- ibtmp file, 117
- ib_logfile, 117
- ignore_builtin_innodb, 58
 - and skip_innodb, 89
- ilist, 117
- implicit row lock, 117
- in-memory database, 117
- incremental backup, 118
- index, 118
- index cache, 118
- index dives (for statistics estimation), 52
- index hint, 118
- index prefix, 118
- indexes
 - creating and dropping, 6
 - primary (clustered) and secondary, 6
- infimum record, 119
- INFORMATION_SCHEMA, 27, 119
 - INNODB_CMP table, 27
 - INNODB_CMPMEM table, 28
 - INNODB_CMPMEM_RESET table, 28
 - INNODB_CMP_RESET table, 27
 - INNODB_LOCKS table, 30
 - INNODB_LOCK_WAITS table, 31
 - INNODB_TRX table, 29
- InnoDB, 119
 - troubleshooting
 - fast index creation, 8
- InnoDB Plugin
 - compatibility, 1
 - downloading, 3
 - features, 1
 - installing, 3
 - restrictions, 4
- innodb_adaptive_flushing, 45
- innodb_adaptive_hash_index, 42
 - and innodb_thread_concurrency, 42
 - dynamically changing, 51
- innodb_additional_mem_pool_size
 - and innodb_use_sys_malloc, 40
- innodb_autoinc_lock_mode, 119
- innodb_change_buffering, 41
- innodb_concurrency_tickets, 42
- innodb_file_format, 19, 119
 - Antelope, 17
 - Barracuda, 9
 - downgrading, 73
 - enabling new file formats, 49
 - identifying, 23
 - innodb_file_format_check, 21
 - innodb_file_io_threads, 44
 - innodb_file_per_table, 9, 120
 - dynamically changing, 50
 - innodb_io_capacity, 45
 - innodb_lock_wait_timeout, 120
 - dynamically changing, 51
 - innodb_max_dirty_pages_pct, 45
 - innodb_old_blocks_pct, 46
 - innodb_old_blocks_time, 46
 - innodb_read_ahead_threshold, 43
 - innodb_read_io_threads, 44
 - innodb_spin_wait_delay, 46
 - innodb_stats_on_metadata
 - dynamically changing, 50
 - innodb_stats_sample_pages, 52
 - innodb_strict_mode, 52, 120
 - innodb_thread_concurrency, 42
 - innodb_thread_sleep_delay, 42
 - innodb_use_sys_malloc, 40
 - and innodb_thread_concurrency, 42
 - innodb_write_io_threads, 44
- insert, 120
- insert buffer, 120
- insert buffering, 77, 120
 - disabling, 41
- installing
 - binary InnoDB Plugin, 58, 89
- instance, 121
- instrumentation, 121
- intention lock, 121
- internal memory allocator
 - disabling, 40
- inverted index, 121
- IOPS, 121
- isolation level, 121

J

- join, 122

K

- KEY_BLOCK_SIZE, 9, 15, 122

L

- latch, 122
- list, 122
- lock, 122
- lock escalation, 122
- lock mode, 122
- lock wait timeout, 51
- locking, 123
 - information schema, 27, 29, 36
- locking read, 123

- log, 123
- log buffer, 123
- log file, 123
- log group, 123
- logical, 123
- logical backup, 124
- loose_, 124
- low-water mark, 124
- LRU, 124
- LRU page replacement, 46
- LSN, 124

M

- master server, 125
- master thread, 125
- MDL, 125
- memcached, 125
- memory allocator
 - innodb_use_sys_malloc, 40
- merge, 125
- metadata lock, 125
- metrics counter, 126
- midpoint insertion, 46
- midpoint insertion strategy, 126
- mini-transaction, 126
- mixed-mode insert, 126
- multi-core, 126
- mutex, 126
- MVCC, 126
- my.cnf, 127
- my.ini, 127
- mysql, 127
- MySQL Enterprise Backup, 127
- mysqlbackup command, 127
- mysqld, 127
- mysqldump, 127

N

- natural key, 127
- neighbor page, 128
- next-key lock, 128
- non-blocking I/O, 128
- non-locking read, 128
- non-repeatable read, 128
- normalized, 128
- NoSQL, 129
- NOT NULL constraint, 129
- NULL, 129

O

- off-page column, 129
- OLTP, 130
- online, 130

- online DDL, 130
- optimistic, 130
- optimizer, 130
- optimizer statistics estimation, 50, 52
- option, 131
- option file, 131
- overflow page, 131

P

- page, 131
- page cleaner, 131
- page size, 131
- parameters, deprecated, 93
 - innodb_file_io_threads, 44
- parameters, new, 91
 - innodb_adaptive_flushing, 45
 - innodb_change_buffering, 41
 - innodb_file_format, 49
 - innodb_file_format_check, 21
 - innodb_io_capacity, 45
 - innodb_read_ahead_threshold, 43
 - innodb_read_io_threads, 44
 - innodb_spin_wait_delay, 46
 - innodb_stats_sample_pages, 52
 - innodb_strict_mode, 52
 - innodb_use_sys_malloc, 40
 - innodb_write_io_threads, 44
- parameters, with new defaults, 93
 - innodb_additional_mem_pool_size, 93
 - innodb_buffer_pool_size, 93
 - innodb_log_buffer_size, 93
 - innodb_max_dirty_pages_pct, 45, 93
 - innodb_sync_spin_loops, 93
 - innodb_thread_concurrency, 93
- parent table, 132
- partial backup, 132
- partial index, 132
- Performance Schema, 132
- persistent statistics, 132
- pessimistic, 132
- phantom, 132
- physical, 132
- physical backup, 133
- PITR, 133
- plan stability, 133
- plugin, 133
- point-in-time recovery, 133
- prepared backup, 133
- primary key, 133
- process, 134
- PROCESSLIST
 - possible inconsistency with INFORMATION_SCHEMA tables, 37

pseudo-record, 134
Pthreads, 134
purge, 134
purge buffering, 134
purge lag, 134
purge thread, 134

Q

query, 134
query execution plan, 134
quiesce, 135

R

RAID, 135
random dive, 135
raw backup, 135
read ahead, 54
 linear, 43
 random, 43
READ COMMITTED, 135
READ UNCOMMITTED, 135
read view, 136
read-ahead, 136
read-only transaction, 136
record lock, 136
redo, 136
redo log, 136
redundant row format, 137
referential integrity, 137
relational, 137
relevance, 137
REPEATABLE READ, 137
replication, 138
restore, 138
rollback, 138
rollback segment, 138
row, 138
row format, 138
row lock, 138
row-based replication, 139
row-level locking, 139
ROW_FORMAT
 COMPACT, 25
 COMPRESSED, 9, 25
 DYNAMIC, 25
 REDUNDANT, 25
rw-lock, 139

S

savepoint, 139
scalability, 139
scale out, 139
scale up, 139

schema, 139
search index, 140
secondary index, 140
segment, 140
selectivity, 140
semi-consistent read, 140
SERIALIZABLE, 141
server, 141
shared lock, 141
shared tablespace, 141
sharp checkpoint, 141
SHOW ENGINE INNODB MUTEX, 54
SHOW ENGINE INNODB STATUS
 and innodb_adaptive_hash_index, 42
 and innodb_use_sys_malloc, 40
shutdown, 141
skip_innodb
 and ignore_builtin_innodb, 89
slave server, 141
slow query log, 142
slow shutdown, 142
snapshot, 142
source code, 65
space ID, 142
spin, 142
SQL, 142
SSD, 9, 142
startup, 143
statement-based replication, 143
statistics, 143
stemming, 143
stopword, 143
storage engine, 143
strict mode, 52, 144
sublist, 144
supremum record, 144
surrogate key, 144
synthetic key, 144
system tablespace, 144

T

table, 145
table lock, 145
table scan, 46
table type, 146
tablespace, 146
tablespace dictionary, 146
temporary table, 146
temporary tablespace, 146
text collection, 147
thread, 147
torn page, 147
TPS, 147

transaction, 147
transaction ID, 147
transportable tablespace, 147
troubleshooting, 147
truncate, 147
TRUNCATE TABLE, 51
tuple, 148
two-phase commit, 148

U

undo, 148
undo log, 148
undo tablespace, 148
unique constraint, 148
unique index, 148
unique key, 149
upgrading
 converting compressed tables, 72
 dynamic plugin, 71
 static plugin, 71

V

victim, 149

W

wait, 149
warm backup, 149
warm up, 149
Windows, 149
workload, 150
write combining, 150

X

XA, 150

Y

young, 150