# InnoDB 1.1 for MySQL 5.5 User's Guide

# InnoDB 1.1 for MySQL 5.5 User's Guide

**Abstract**

This is the User's Guide for the InnoDB storage engine 1.1 for MySQL 5.5.

Beginning with MySQL version 5.1, it is possible to swap out one version of the InnoDB storage engine and use another (the "plugin"). This manual documents the latest InnoDB plugin, version 1.1, which works with MySQL 5.5 and features cutting-edge improvements in performance and scalability.

This User's Guide documents the procedures and features that are specific to the InnoDB storage engine 1.1 for MySQL 5.5. It supplements the general InnoDB information in the MySQL Reference Manual.

Because InnoDB 1.1 is integrated with MySQL 5.5, it is generally available (GA) and production-ready.

**WARNING: Because the InnoDB storage engine 1.0 and above introduces a new file format, restrictions apply to the use of a database created with the InnoDB storage engine 1.0 and above, with earlier versions of InnoDB, when using `mysqldump` or MySQL replication and if you use the older InnoDB Hot Backup product rather than the newer MySQL Enterprise Backup product. See Section 1.4, "Compatibility Considerations for Downgrade and Backup".**

For legal information, see the Legal Notices.

Document generated on: 2014-01-30 (revision: 37565)

# Table of Contents

# Preface and Legal Notices

This is the User's Guide for the InnoDB storage engine 1.1 for MySQL 5.5.

## Legal Notices

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 InnoDB 1.1

## Table of Contents

InnoDB 1.1 combines the familiar reliability and performance of the InnoDB storage engine, with new performance and usability enhancements. InnoDB 1.1 includes all the features that were part of the InnoDB Plugin for MySQL 5.1, plus new features specific to MySQL 5.5 and higher.

Beginning with MySQL version 5.5, InnoDB is the default storage engine, rather than MyISAM, to promote greater data reliability and reducing the chance of corruption.

## 1.1 Features of the InnoDB Storage Engine

InnoDB in MySQL 5.5 contains several important new features:

- 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 Change 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 InnoDB I/O Threads"

    - Section 7.9, "Asynchronous I/O on Linux"

    - Section 7.10, "Group Commit"

    - Section 7.11, "Controlling the InnoDB Master Thread I/O Rate"

    - Section 7.12, "Controlling the Flushing Rate of Dirty Pages from the InnoDB Buffer Pool"

    - Section 7.13, "Using the PAUSE Instruction in InnoDB Spin Loops"

## Upward and Downward Compatibility

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 5.1 and earlier, is now called Antelope and does not support these features, but does support the other features introduced with the InnoDB storage engine.

The InnoDB storage engine is upward compatible from standard InnoDB as built in to, and distributed with, MySQL. Existing databases can be used with the InnoDB Storage Engine for MySQL. 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.

InnoDB since version 5.0.21 has a safety feature that prevents it from opening tables that are in an unknown format. However, the system tablespace may contain references to new-format tables that confuse the built-in InnoDB in MySQL 5.1 and earlier. These references are cleared in a slow shutdown.

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. To provide early feedback, InnoDB 1.1 checks the system tablespace

before startup to ensure that the file format used in the database is supported by the storage engine. See Section 4.2.1, "Compatibility Check When InnoDB Is Started" for the details.

## 1.2 Obtaining and Installing the InnoDB Storage Engine

Starting with MySQL 5.4.2, you do not need to do anything special to get or install the most up-to-date InnoDB storage engine. From that version forward, the InnoDB storage engine in the server is what was formerly known as the InnoDB Plugin. Earlier versions of MySQL required some extra build and configuration steps to get the Plugin-specific features such as fast index creation and table compression.

Report any bugs in the InnoDB storage engine using the My Oracle Support site. For general discussions about InnoDB Storage Engine for MySQL, see http://forums.mysql.com/list.php?22.

## 1.3 Viewing the InnoDB Storage Engine Version Number

InnoDB storage engine releases are numbered with version numbers independent of MySQL release numbers. The initial release of the InnoDB storage engine is version 1.0, and it is designed to work with MySQL 5.1. Version 1.1 of the InnoDB storage engine is for MySQL 5.5 and up.

- The first component of the InnoDB storage engine version number designates a major release level.

- The second component corresponds to the MySQL release. The digit 0 corresponds to MySQL 5.1. The digit 1 corresponds to MySQL 5.5.

- The third component indicates the specific release of the InnoDB storage engine (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 storage engine, 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 storage engine 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.

## 1.4 Compatibility Considerations for Downgrade and Backup

**Because InnoDB 1.1 supports the "Barracuda" file format, with new on-disk data structures within both the database and log files, pay special attention to file format compatibility with respect to the following scenarios:**

- Downgrading from MySQL 5.5 to the MySQL 5.1 or earlier (without the InnoDB Plugin enabled), or otherwise using earlier versions of MySQL with database files created by MySQL 5.5 and higher.

- Using `mysqldump`.

- Using MySQL replication.

- Using MySQL Enterprise Backup or InnoDB Hot Backup.

**WARNING:** Once you create any tables with the Barracuda file format, take care to avoid crashes and corruptions when using those files with an earlier version of MySQL. It is **strongly** recommended that you use a "slow shutdown" (`SET GLOBAL innodb_fast_shutdown=0`) when stopping the MySQL server before downgrading to MySQL 5.1 or earlier. This ensures that the log files and other system information do not cause consistency issues or startup problems when using a prior version of MySQL.

**WARNING:** If you dump a database containing compressed tables with `mysqldump`, the dump file may contain `CREATE TABLE` statements that attempt to create compressed tables, or those using `ROW_FORMAT=DYNAMIC` in the new database. Therefore, be sure the new database is running the InnoDB storage engine, 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, when the `mysqldump` file is loaded, MySQL and InnoDB ignore `CREATE TABLE` options they do not recognize, and the table(s) are created in a format used by the running server.

**WARNING:** If you use MySQL replication, ensure all slaves are configured with the InnoDB storage engine, 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 an older version of MySQL, it ignores the `CREATE TABLE` options to create a compressed table or one with `ROW_FORMAT=DYNAMIC`, and creates the table uncompressed, with `ROW_FORMAT=COMPACT`.

**WARNING:** Version 3.0 of InnoDB Hot Backup does not support the new Barracuda file format. Using InnoDB Hot Backup Version 3 to backup databases in this format causes unpredictable behavior. MySQL Enterprise Backup, the successor product to InnoDB Hot Backup, does support tables with the Barracuda file format. You can also back up such databases with `mysqldump`.

# Chapter 2 Fast Index Creation in the InnoDB Storage Engine

## Table of Contents

In MySQL 5.5 and higher, or in MySQL 5.1 with the InnoDB Plugin, creating and dropping secondary indexes does not copy the contents of the entire table, making this operation much more efficient than with prior releases.

## 2.1 Overview of Fast Index Creation

With MySQL 5.5 and higher, or MySQL 5.1 with the InnoDB Plugin, creating and dropping secondary indexes for InnoDB tables is much faster than before. Historically, adding or dropping an index on a table with existing data could be very slow. The `CREATE INDEX` and `DROP INDEX` statements worked by creating a new, empty table defined with the requested set of indexes, then copying the existing rows to the new table one-by-one, updating the indexes as the rows are inserted. After all rows from the original table were copied, the old table was dropped and the copy was renamed with the name of the original table.

The performance speedup for fast index creation applies to secondary indexes, not to the primary key index. The rows of an InnoDB table are stored in a clustered index organized based on the primary key, forming what some database systems call an "index-organized table". Because the table structure is so closely tied to the primary key, redefining the primary key still requires 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 of Fast Index Creation").

## 2.2 Examples of Fast Index Creation

It is possible to create multiple indexes on a table with one `ALTER TABLE` statement. 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 statements 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` statement, this approach is much more efficient than creating all the secondary indexes before loading the data.

You can 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` statement. Thus, the following statements are not as efficient as the `ALTER TABLE` statement 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 InnoDB secondary indexes also does not require any copying of table data. You can equally quickly drop multiple indexes with a single `ALTER TABLE` statement or multiple `DROP INDEX` statements:

```
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);
```

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 Details of Fast Index Creation

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 can be created or dropped without copying the data in the clustered index. Because each secondary index contains copies of the primary key values (used to access the clustered index when needed), when you change the definition of the primary key, 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 are 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 can 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 values of the secondary index key columns. The B-tree is then built in key-value order, which is more efficient than inserting rows into an index in random order. 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 for Fast Index Creation

While an InnoDB secondary index is being created or dropped, the table is locked in shared mode. Any writes to the table are blocked, but the data in the table can be read. When you alter the clustered index of a table, 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.

A `CREATE INDEX` or `ALTER TABLE` statement for an InnoDB table always waits for currently executing transactions that are accessing the table to commit or roll back. `ALTER TABLE` statements that redefine an InnoDB primary key wait for all `SELECT` statements that access the table to complete, or their containing transactions to commit. No transactions whose execution spans the creation of the index can be accessing the table, because the original table is dropped when the clustered index is restructured.

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

A newly-created InnoDB secondary index contains only the committed data in the table at the time the `CREATE INDEX` or `ALTER TABLE` statement begins to execute. It does not contain any uncommitted values, old versions of values, or values marked for deletion but not yet removed from the old index.

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 How Crash Recovery Works with Fast Index Creation

Although no data is lost if the server crashes while an `ALTER TABLE` statement is executing, the crash recovery process is different for clustered indexes and secondary indexes.

If the server crashes while creating an InnoDB secondary index, upon recovery, MySQL drops any partially created indexes. You must re-run the `ALTER TABLE` or `CREATE INDEX` statement.

When a crash occurs during the creation of an InnoDB clustered index, recovery is 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.

MySQL creates the new clustered index by copying the existing data from the original InnoDB 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 renamed with the name of the original table, and the original table is dropped from the database.

If a system crash occurs while creating a new clustered index, no data is lost, but you must complete the recovery process using the temporary tables that exist during the process. Since it is rare to re-create a clustered index or re-define primary keys on large tables, or to encounter a system crash during this operation, this manual does not provide information on recovering from this scenario. Contact MySQL support.

## 2.6 Limitations of Fast Index Creation

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

- During index creation, files are written to the temporary directory (`$TMPDIR` on Unix, `%TEMP%` on Windows, or the value of the `--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.

- An `ALTER TABLE` statement that contains `DROP INDEX` and `ADD INDEX` clauses that both name the same index uses a table copy, not Fast Index Creation.

- 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 statement `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 for Fast Index Creation".

  Prior to InnoDB storage engine 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 storage engine 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 storage engine 1.0.2 introduces some improvements in error handling when users attempt to drop indexes. See section Section 8.6, "Better Error Handling when Dropping Indexes" for details.

- MySQL 5.5 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 is copied, rather than using Fast Index Creation.

- `OPTIMIZE TABLE` for an `InnoDB` table is mapped to an `ALTER TABLE` operation to rebuild the table and update index statistics and free unused space in the clustered index. This operation does not use fast index creation. Secondary indexes are not created as efficiently because keys are inserted in the order they appeared in the primary key.

# Chapter 3 Working with `InnoDB` Compressed Tables

## Table of Contents

By using the SQL syntax and InnoDB configuration options for compression, you can create tables where the data is stored in compressed form. Compression can help to improve both raw performance and scalability. The compression means less data is transferred between disk and memory, and takes up less space on disk and in memory. The benefits are amplified for tables with secondary indexes, because index data is compressed also. Compression can be especially important for SSD storage devices, because they tend to have lower capacity than HDD devices.

## 3.1 Overview of Table Compression

Because processors and cache memories have increased in speed more than disk storage devices, many workloads are disk-bound. Data compression enables smaller database size, reduced I/O, and improved throughput, at the small cost of increased CPU utilization. Compression is especially valuable for read-intensive applications, on systems with enough RAM to keep frequently used data in memory.

An InnoDB table created with `ROW_FORMAT=COMPRESSED` can use a smaller page size on disk than the usual 16KB default. Smaller pages require less I/O to read from and write to disk, which is especially valuable for SSD devices.

The page size is specified through the `KEY_BLOCK_SIZE` parameter. The different page size means the table must be in its own `.ibd` file rather than in the system tablespace, which requires enabling the `innodb_file_per_table` option. The level of compression is the same regardless of the `KEY_BLOCK_SIZE` value. As you specify smaller values for `KEY_BLOCK_SIZE`, you get the I/O benefits of increasingly smaller pages. But if you specify a value that is too small, there is additional overhead to reorganize the pages when data values cannot be compressed enough to fit multiple rows in each page. There is a hard limit on how small `KEY_BLOCK_SIZE` can be for a table, based on the lengths of the key columns for each of its indexes. Specify a value that is too small, and the `CREATE TABLE` or `ALTER TABLE` statement fails.

In the buffer pool, the compressed data is held in small pages, with a page size based on the `KEY_BLOCK_SIZE` value. For extracting or updating the column values, MySQL also creates a 16KB page in the buffer pool with the uncompressed data. Within the buffer pool, any updates to the uncompressed page are also re-writtten back to the equivalent compressed page. You might need to size your buffer pool to accommodate the additional data of both compressed and uncompressed pages, although the uncompressed pages are evicted from the buffer pool when space is needed, and then uncompressed again on the next access.

## 3.2 Enabling Compression for a Table

The default uncompressed size of InnoDB data pages is 16KB. 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 uses a page size of 1KB, 2KB, 4KB, 8KB, or 16KB for the `.ibd` file of the table. (The actual compression algorithm is not affected by the `KEY_BLOCK_SIZE` value.)

> **Note**
>
> Compression is applicable to tables, not to individual rows, despite the option name `ROW_FORMAT`.

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` typically does not result in much compression, since the normal InnoDB page size is 16KB. This setting may still 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 in Section 3.4, "How Compression Works for InnoDB Tables".

All indexes of a table (including the clustered index) are compressed using the same page size, as specified in 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`. If necessary, you can set these parameters in the MySQL option file `my.cnf` or `my.ini`, or with the `SET` statement without shutting down the MySQL server.

Specifying `ROW_FORMAT=COMPRESSED` or `KEY_BLOCK_SIZE` in `CREATE TABLE` or `ALTER TABLE` statements produces these warnings if the Barracuda file format is not enabled. You can view them 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. |

| Level | Code | Message |
|-------|------|---------|
| 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. When InnoDB "strict mode" (see Section 8.4, "`InnoDB Strict Mode`") is enabled, InnoDB generates 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, MySQL 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 storage engine with the proper settings for the configuration parameters `innodb_file_format` and `innodb_file_per_table`,

# 3.3 Tuning Compression for InnoDB Tables

Most often, the internal optimizations described in InnoDB Data Storage and Compression ensure that the system runs well with compressed data. However, because the efficiency of compression depends on the nature of your data, you can make decisions that affect the performance of compressed tables:

- Which tables to compress.

- What compressed page size to use.

- Whether 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. Whether the workload is more like a data warehouse (primarily queries) or an OLTP system (mix of queries and DML).

- If the system performs DML operations on compressed tables, and the way the data is distributed leads to expensive compression failures at runtime, you might adjust additional advanced configuration options.

Use the guidelines in this section to help make those architectural and configuration choices. When you are ready to conduct long-term testing and put compressed tables into production, see Monitoring Compression at Runtime for ways to verify the effectiveness of those choices under real-world conditions.

## 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.

## 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.

You choose whether to turn on compression for each InnoDB table. 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 might result in long column values being stored "off-page", as discussed in Section 5.3, "`DYNAMIC` and `COMPRESSED` Row Formats". Those overflow pages may compress well. Given these considerations, for many applications, some tables compress more effectively than others, and you might 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.

## Compression and Application and Schema Design

Decide whether to compress data in your application or in the table; do not use both types of compression for the same data. When you compress the data in the application and store the results in a compressed table, extra space savings are extremely unlikely, and the double compression 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.

## 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 `BLOB`s 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 `INSERT`s 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 might find that compression can improve overall performance. When you test your application performance with different compression configurations, test on a platform similar to the planned configuration of the production system.

## Configuration Characteristics and Compression

Reading and writing database pages from and to disk is the slowest aspect of system performance. Compression attempts to reduce I/O by using CPU time to compress and uncompress data, and 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 storage engine 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. Still, 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.

## 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 Compression of B-Tree Pages.

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, you set the compressed page size to 8K or 4K bytes. Given that the maximum row size for an InnoDB table is around 8K, `KEY_BLOCK_SIZE=8` is usually a safe choice.

## Monitoring Compression at Runtime

Overall application performance, CPU and I/O utilization and the size of disk files are good indicators of how effective compression is for your application.

To dig deeper into performance considerations for compressed tables, you can monitor compression performance at runtime. using the Information Schema tables described in Example 6.1, "Using the Compression Information Schema Tables". These tables 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, 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 might 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. Increasing the size of the InnoDB buffer pool might also help performance, so that more uncompressed pages can stay in memory, reducing the need to uncompress pages that 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. If so, 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. If the ratio is low, then InnoDB is reorganizing, recompressing, and splitting B-tree nodes more often than is desirable. In this case, avoid compressing some tables, or increase `KEY_BLOCK_SIZE` for some of the compressed tables. You might turn off compression for tables that cause the number of "compression failures" in your application to be more than 1% or 2% of the total. (Such a failure ratio might be acceptable during a temporary operation such as a data load).

# 3.4 How Compression Works for InnoDB Tables

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

## 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.

InnoDB implements 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 `.idb` 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. You can adjust the balance between compression level and CPU overhead by modifying the `innodb_compression_level` configuration option.

# InnoDB Data Storage and Compression

All user data in InnoDB tables is stored in pages comprising a B-tree index (the 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 tables 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.

# 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 (a situation known as a compression failure), 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 storage engine version 1.0.2, and if

`innodb_strict_mode` is `ON`, the InnoDB storage engine 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`.

## Compressing BLOB, VARCHAR, and TEXT Columns

In an InnoDB table, `BLOB`, `VARCHAR`, and `TEXT` columns that are not part of the primary key may be stored on separately allocated overflow pages. We refer to these columns as off-page columns. Their 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, MySQL 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 older versions of InnoDB use the Antelope file format, which supports only `ROW_FORMAT=REDUNDANT` and `ROW_FORMAT=COMPACT`. In these formats, MySQL 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. Image data, such as `JPEG`, is typically already compressed and so does not benefit much from being stored in a compressed table; the double compression can waste CPU cycles for little or no space savings.

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.

## 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 (or a smaller size if `innodb_page_size` is set). To access the

data in a page, InnoDB reads the compressed page from disk if it is not already in the buffer pool, then uncompresses the page to its original 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. 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 might be written to disk, to free space for other data. Thus, at any given time, the buffer pool might contain both the compressed and uncompressed forms of the page, or only the compressed form of the page, or neither.

InnoDB keeps track of which pages to keep in memory and which to evict using a least-recently-used (LRU) list, so that hot (frequently accessed) data tends to stay in memory. When compressed tables are accessed, MySQL 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 goal 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, MySQL 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.

## Compression and the InnoDB Redo Log Files

Before a compressed page is written to a data file, MySQL 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 are usable for crash recovery, even in the unlikely case that the `zlib` library is upgraded and that change introduces a compatibility problem with the compressed data. 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 compressed tables use a different file format for the redo log and the per-table tablespaces than in MySQL 5.1 and earlier. The MySQL Enterprise Backup product supports this latest Barracuda file format for compressed InnoDB tables. The older InnoDB Hot Backup product can only back up tables using the file format Antelope, and thus does not support compressed InnoDB tables.

# 3.5 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 with `innodb_strict_mode` enabled, 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 |
|---|---|---|
| `ROW_FORMAT=REDUNDANT` | Storage format used prior to MySQL 5.0.3 | Less efficient than `ROW_FORMAT=COMPACT`; for backward compatibility |
| `ROW_FORMAT=COMPACT` | 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 |
| `ROW_FORMAT=DYNAMIC` | Available only with `innodb_file_format=Barracuda` | Store values within the clustered index page if they fit; if not, stores only a 20-byte pointer to an overflow page (no prefix) |
| `ROW_FORMAT=COMPRESSED` | Available only with `innodb_file_format=Barracuda` | Compresses the table and indexes using zlib to default compressed page size of 8K bytes; implies `ROW_FORMAT=DYNAMIC` |
| `KEY_BLOCK_SIZE=n` | Available only with `innodb_file_format=Barracuda` | Specifies compressed page size of 1, 2, 4, 8 or 16 kilobytes; implies `ROW_FORMAT=DYNAMIC` and `ROW_FORMAT=COMPRESSED` |

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 ignores 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. To see the full description of the error condition, issue the `SHOW ERRORS` statement: 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` | Ignored unless both `innodb_file_format=Barracuda` and `innodb_file_per_table` are enabled | `COMPACT` |

| Syntax | Warning or Error Condition | Resulting `ROW_FORMAT`, as shown in `SHOW TABLE STATUS` |
|---|---|---|
| or `KEY_BLOCK_SIZE` is specified | | |
| 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`, the InnoDB storage engine rejects invalid `ROW_FORMAT` or `KEY_BLOCK_SIZE` parameters. For compatibility with earlier versions of MySQL, strict mode is not enabled by default; instead, MySQL 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 when creating the table). The real compressed page size of the table cannot be displayed by MySQL.

# Chapter 4 `InnoDB` File-Format Management

## Table of Contents

As InnoDB evolves, new on-disk data structures are sometimes required to support new features. Features such as compressed tables (see Chapter 3, *Working with `InnoDB` Compressed Tables*), and long variable-length columns stored off-page (see Chapter 5, *`InnoDB` Row Storage and Row Formats*) require data file formats that are not compatible with prior versions of InnoDB. These features both require use of the new Barracuda file format.

> **Note**
>
> All other new features are compatible with the original Antelope file format and do not require the Barracuda file format.

This section discusses enabling file formats for new InnoDB tables, verifying compatibility of different file formats between MySQL releases, identifying the file format in use, downgrading the file format, and file format names that may be used in the future.

**Named File Formats.**     InnoDB 1.1 has 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 original InnoDB file format is called 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.

## 4.1 Enabling File Formats

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

Although Oracle recommends using the Barracuda format for new tables where practical, in MySQL 5.5 the default file format is still Antelope, for maximum compatibility with replication configurations containing different MySQL releases.

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.

## 4.2 Verifying File Format Compatibility

InnoDB 1.1 incorporates several checks to guard against the possible crashes and data corruptions that might occur if you run an older release of the MySQL server on InnoDB data files using a newer file format. 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 might arise.

# Backward Compatibility

Considerations of backward compatibility only apply when using a recent version of InnoDB (the InnoDB Plugin, or MySQL 5.5 and higher with InnoDB 1.1) alongside an older one (MySQL 5.1 or earlier, with the built-in InnoDB rather than the InnoDB Plugin). To minimize the chance of compatibility issues, you can standardize on the InnoDB Plugin for all your MySQL 5.1 and earlier database servers.

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. InnoDB 1.1 includes a 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 (such as performance improvements and bug fixes), and still preserve the option of using your database with a prior version of InnoDB, by preventing accidental use of new features that create downward-incompatible disk files.

If a version of InnoDB supports a particular file format (whether or not that format is the default), you can query and 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` statement, as described in Section 4.4, "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.

## Internal Details

Every InnoDB per-table tablespace (represented by a `*.ibd` file) file is labeled with a file format identifier. The system tablespace (represented by the `ibdata` files) 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.

Creating a compressed table, or a table with `ROW_FORMAT=DYNAMIC`, updates the file header for the corresponding `.ibd` file and the table type in the InnoDB data dictionary 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. (In many cases, the `ALTER TABLE` statement recreates a table and thus changes 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*.)

## Definition of ib-file set

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).

- InnoDB log files; usually two, `ib_logfile0` and `ib_logfile1`. Used for crash recovery and in backups.

An "ib-file set" does not include the corresponding `.frm` files that contain metadata about InnoDB tables. The `.frm` files are created and managed by MySQL, and can sometimes get out of sync with the internal metadata in InnoDB.

Multiple tables, even from more than one database, can be stored in a single "ib-file set". (In MySQL, a "database" is a logical collection of tables, what other systems refer to as a "schema" or "catalog".)

## 4.2.1 Compatibility Check When `InnoDB` Is Started

To prevent possible crashes or data corruptions when InnoDB 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 consistency issues with the new tables or startup problems for the MySQL server.

Beginning with version InnoDB 1.0.1, 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 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 InnoDB 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 the InnoDB Storage Engine*). 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`, InnoDB 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".)

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 4.2.2, "Compatibility Check When a Table Is Opened". 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 in Antelope or Barracuda file format; tables in Cheetah format cannot be accessed |
| OFF | `Cheetah` | `Barracuda` | `Barracuda` | Database cannot be opened; `innodb_file_format` cannot be set to Cheetah |
| ON | `Antelope` or `Barracuda` | `Barracuda` | `Barracuda` | Database can be opened; tables can be created in 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; `innodb_file_format` cannot be set to Cheetah |

## 4.2.2 Compatibility Check When a Table Is Opened

When a table is first accessed, InnoDB (including some releases prior to InnoDB 1.0) checks 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.

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 Named File Formats, versions of MySQL 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 (for example, 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
```

InnoDB also 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 4.3, "Identifying the File Format in Use".

Versions of InnoDB prior to MySQL 4.1 did not include table format identifiers in the database files, and versions prior to MySQL 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 file format management capability in InnoDB 1.0 and higher (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 MySQL 5.1 with the InnoDB Plugin. It is possible to "downgrade" such tables to the Antelope format with the procedure described in Section 4.4, "Downgrading the File Format".

# 4.3 Identifying the File Format in Use

After you enable a given `innodb_file_format`, this change applies only to newly created tables rather than existing ones. 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'`. (The `Row_format` is a separate column; 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'`.

## Internal Details

InnoDB has two different file formats (Antelope and Barracuda) and four different row formats (Redundant, Compact, Dynamic, and Compressed). The Antelope file format contains Redundant and Compact row

formats. A tablespace that uses the Barracuda file format uses either the Dynamic or Compressed row format.

File and row format information is written in 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 InnoDB storage engine up to MySQL 5.1.

The first 6 bits of the tablespace flags can be described this way:

- Bit 0: Zero for Antelope and no other bits will be set. One for Barracuda, and other bits may be set.

- Bits 1 to 4: A 4 bit number representing the compressed page size. zero = not compressed, 1 = 1k, 2 = 2k, 3 = 4k, 4 = 8k.

- Bit 5: Same value as Bit 0, zero for Antelope, and one for Barracuda. If Bit 0 and Bit 5 are set and Bits 1 to 4 are not, the row format is "Dynamic".

Until MySQL 5.6, no other bits are set in the tablespace flags. If bits 6 to 31 are not zero, the tablespace is corrupt or is not an InnoDB tablespace file.

## 4.4 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.

## 4.5 Future `InnoDB` File Formats

The file format used by the standard built-in InnoDB in MySQL 5.1 is the Antelope format. The file format introduced with InnoDB Plugin 1.0 is the Barracuda format. Although no new features have been announced that would require additional new file formats, the InnoDB file format mechanism allows for future 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

This section discusses how certain InnoDB features, such as table compression and off-page storage of long columns, are controlled by the `ROW_FORMAT` clause of the `CREATE TABLE` statement. It discusses considerations for choosing the right row format and compatibility of row formats between MySQL releases.

## 5.1 Overview of `InnoDB` Row Storage

The storage for rows and associated columns affects performance for queries and DML operations. As more rows fit into a single disk page, queries and index lookups can work faster, less cache memory is required in the InnoDB buffer pool, and less I/O is required to write out updated values for the numeric and short string columns.

The data in each InnoDB table is divided into pages. The pages that make up each table are arranged in a tree data structure called a B-tree index. Table data and secondary indexes both use this type of structure. The B-tree index that represents an entire table is known as the clustered index, which is organized according to the primary key columns. The nodes of the index data structure contain the values of all the columns in that row (for the clustered index) or the index columns and the primary key columns (for secondary indexes).

Variable-length columns are an exception to this rule. Columns such as `BLOB` and `VARCHAR` that are too long to fit on a B-tree page are stored on separately allocated disk pages called overflow pages. We call such columns off-page columns. The values of these columns are stored in 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 value is stored in the B-tree, to avoid wasting storage and eliminating the need to read a separate page.

The 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.

The next section describes the clauses you can use with the `CREATE TABLE` and `ALTER TABLE` statements to control how these variable-length columns are represented: `ROW_FORMAT` and `KEY_BLOCK_SIZE`. To use these clauses, you might also need to change the settings for the `innodb_file_per_table` and `innodb_file_format` configuration options.

## 5.2 Specifying the Row Format for a Table

You specify the row format for a table with the `ROW_FORMAT` clause of the `CREATE TABLE` and `ALTER TABLE` statements.

## 5.3 `DYNAMIC` and `COMPRESSED` Row Formats

This section discusses the `DYNAMIC` and `COMPRESSED` row formats for InnoDB tables. You can only create these kinds of tables when the `innodb_file_format` configuration option is set to `Barracuda`. (The Barracuda file format also allows the `COMPACT` and `REDUNDANT` row formats.)

When 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 based 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.

The `COMPRESSED` row format uses similar internal details for off-page storage as the `DYNAMIC` row format, with additional storage and performance considerations from the table and index data being compressed and using smaller page sizes. With the `COMPRESSED` row format, the option `KEY_BLOCK_SIZE` controls how much column data is stored in the clustered index, and how much is placed on overflow pages. For full details about the `COMPRESSED` row format, see Chapter 3, *Working with `InnoDB` Compressed Tables*.

## 5.4 COMPACT and REDUNDANT Row Formats

Early versions of InnoDB used an unnamed file format (now called Antelope) for database files. With that file format, tables are defined with `ROW_FORMAT=COMPACT` or `ROW_FORMAT=REDUNDANT`. InnoDB stores 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 pages.

To preserve compatibility with those prior versions, tables created with the newest InnoDB default to the `COMPACT` row format. See Section 5.3, "DYNAMIC and COMPRESSED Row Formats" for information about the newer `DYNAMIC` and `COMPRESSED` row formats.

With the Antelope file format, if the value of a column is 768 bytes or less, 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, 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.

# Chapter 6 `InnoDB INFORMATION_SCHEMA` tables

## Table of Contents

The `INFORMATION_SCHEMA` is a MySQL feature that helps you monitor server activity to diagnose capacity and performance issues. Several InnoDB-related `INFORMATION_SCHEMA` tables (`INNODB_CMP`, `INNODB_CMP_RESET`, `INNODB_CMPMEM`, `INNODB_CMPMEM_RESET`, `INNODB_TRX`, `INNODB_LOCKS` and `INNODB_LOCK_WAITS`) contain live information about 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).

The Information Schema tables are themselves plugins to the MySQL server, and must be activated by `INSTALL` statements. If they are installed, but the InnoDB storage engine plugin is not installed, these tables appear to be empty.

This section describes the InnoDB-related Information Schema tables and shows some examples of their use.

## 6.1 Information Schema Tables about Compression

Two new pairs of Information Schema tables 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.

### 6.1.1 `INNODB_CMP` and `INNODB_CMP_RESET`

The `INNODB_CMP` and `INNODB_CMP_RESET` tables contain status information on the operations related to compressed tables, which are covered in Chapter 3, *Working with `InnoDB` Compressed Tables*. 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 archive the output of `INNODB_CMP_RESET` every 60 minutes, you see the statistics for each hourly period. If you monitor the output of `INNODB_CMP` (making sure never to read `INNODB_CMP_RESET`), you see the cumulated statistics since InnoDB was started.

For the table definition, see Columns of `INNODB_CMP` and `INNODB_CMP_RESET`.

### 6.1.2 `INNODB_CMPMEM` and `INNODB_CMPMEM_RESET`

The `INNODB_CMPMEM` and `INNODB_CMPMEM_RESET` tables contain status information on the compressed pages that reside in the buffer pool. Please consult Chapter 3, *Working with `InnoDB` Compressed Tables* for further information on compressed tables and the use of the buffer pool. The `INNODB_CMP` and `INNODB_CMP_RESET` tables should provide more useful statistics on compression.

## Internal Details

InnoDB uses a 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.

For the table definition, see Columns of INNODB_CMPMEM and INNODB_CMPMEM_RESET.

## 6.1.3 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, *Working with `InnoDB` Compressed Tables*, `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 a light workload. 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 |

| page size | pages used | pages free | relocation ops | relocation time |
|-----------|-----------|-----------|----------------|-----------------|
| 8192 | 6169 | 0 | 5 | 0 |
| 16384 | 0 | 0 | 0 | 0 |

# 6.2 Information Schema Tables about Transactions

Three InnoDB-related Information Schema tables make it easy to monitor transactions and diagnose possible locking problems. The three tables are `INNODB_TRX`, `INNODB_LOCKS`, and `INNODB_LOCK_WAITS`.

- `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.

  For the table definition, see `INNODB_TRX` Columns.

- `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.

  For the table definition, see `INNODB_LOCKS` Columns.

- `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`.

  For the table definition, see `INNODB_LOCK_WAITS` Columns.

## 6.2.1 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 | RUN- NING | 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 6.3.3, "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 | RUN-NING | 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 … |

| trx id | trx state | trx started | trx requested lock id | trx wait started | trx weight | trx mysql thread id | trx query |
|---|---|---|---|---|---|---|---|
| *19C* | RUN-NING | 2008-01-15 13:09:10 | NULL | NULL | 2900 | 130 | insert into t2 values … |
| *E15* | RUN-NING | 2008-01-15 13:08:59 | NULL | NULL | 5395 | 61 | insert into t2 values … |
| *51D* | RUN-NING | 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 system running a heavy [workload](#).

| requesting trx id | requested lock id | blocking trx id | blocking lock id |
|---|---|---|---|
| *77F* | *77F*:806 | *77E* | *77E*:806 |
| *77F* | *77F*:806 | *77D* | *77D*:806 |
| *77F* | *77F*:806 | *77B* | *77B*:806 |
| *77E* | *77E*:806 | *77D* | *77D*:806 |
| *77E* | *77E*:806 | *77B* | *77B*:806 |
| *77D* | *77D*:806 | *77B* | *77B*:806 |
| *77B* | *77B*:733:12:1 | *77A* | *77A*:733:12:1 |
| *E56* | *E56*:743:6:2 | *E55* | *E55*:743:6:2 |
| *E55* | *E55*:743:38:2 | *19C* | *19C*: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 |
|---|---|---|---|---|---|---|---|---|---|
| *77F*:806 | *77F* | AUTO_INC | TABLE | `test`.`t09` | NULL | NULL | NULL | NULL | NULL |
| *77E*:806 | *77E* | AUTO_INC | TABLE | `test`.`t09` | NULL | NULL | NULL | NULL | NULL |
| *77D*:806 | *77D* | AUTO_INC | TABLE | `test`.`t09` | NULL | NULL | NULL | NULL | NULL |
| *77B*:806 | *77B* | AUTO_INC | TABLE | `test`.`t09` | NULL | NULL | NULL | NULL | NULL |
| *77B*:733:12:1 | *77B* | X | RECORD | `test`.`t09` | `PRIMARY` | 733 | 12 | 1 | supremum pseudo-record |
| *77A*:733:12:1 | *77A* | X | RECORD | `test`.`t09` | `PRIMARY` | 733 | 12 | 1 | supremum pseudo-record |
| *E56*:743:6:2 | *E56* | S | RECORD | `test`.`t2` | `PRIMARY` | 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 |
|---------|-------------|-----------|-----------|------------|------------|------------|-----------|----------|-----------|
| *E55*:743:6:2 | *E55* | X | RECORD | `` `test` `` `` .`t2` `` | `` `PRIMARY` `` | 743 | 6 | 2 | 0, 0 |
| *E55*:743 :38:2 | *E55* | S | RECORD | `` `test` `` `` .`t2` `` | `` `PRIMARY` `` | 743 | 38 | 2 | 1922, 1922 |
| *19C*:743 :38:2 | *19C* | X | RECORD | `` `test` `` `` .`t2` `` | `` `PRIMARY` `` | 743 | 38 | 2 | 1922, 1922 |

# 6.3 Special Locking Considerations for InnoDB `INFORMATION_SCHEMA` Tables

## 6.3.1 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 roll back).

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 a lock, it is in the `'RUNNING'` state. If a transaction is waiting for a lock, it is in the `'LOCK WAIT'` state.

The `INNODB_LOCKS` table 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 `INNODB_LOCK_WAITS` table shows which locks already held by a transaction are blocking locks requested by other transactions.

## 6.3.2 Granularity of `INFORMATION_SCHEMA` 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 changes only 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 `JOIN`s 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.

### 6.3.3 Possible Inconsistency with `PROCESSLIST`

As just described, while the transaction and locking data is correct and consistent when these `INFORMATION_SCHEMA` tables are populated. For example, 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. However, the underlying data changes so fast that similar glimpses at other, similarly fast-changing data, may not be in synchrony. Thus, you should be careful in comparing the data in the InnoDB transaction and locking tables with that in the `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 differ from the one in `PROCESSLIST.INFO`.

# Chapter 7 `InnoDB` Performance and Scalability Enhancements

## Table of Contents

This section discusses recent InnoDB enhancements to performance and scalability, covering the performance features in InnoDB 1.1 with MySQL 5.5, and the features in the InnoDB Plugin for MySQL 5.1. This information is useful to any DBA or developer who is concerned with performance and scalability. Although some of the enhancements do not require any action on your part, knowing this information can still help you diagnose performance issues more quickly and modernize systems and applications that rely on older, inefficient behavior.

## 7.1 Overview of InnoDB Performance

InnoDB has always been highly efficient, and includes several unique architectural elements to assure high performance and scalability. The latest InnoDB storage engine includes new features that take advantage of advances in operating systems and hardware platforms, such as multi-core processors and improved memory allocation systems. In addition, new configuration options let you better control some InnoDB internal subsystems to achieve the best performance with your workload.

Starting with MySQL 5.5 and InnoDB 1.1, the built-in InnoDB storage engine within MySQL is upgraded to the full feature set and performance of the former InnoDB Plugin. This change makes these performance and scalability enhancements available to a much wider audience than before, and eliminates the separate installation step of the InnoDB Plugin. After learning about the InnoDB performance features in this section, continue with Optimization to learn the best practices for overall MySQL performance, and Optimizing for `InnoDB` Tables in particular for InnoDB tips and guidelines.

## 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.

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 that is compiled with GCC version 4.1.2 or later uses 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, functions for Interlocked Variable Access are available that are similar to the built-in functions provided by GCC. On Windows 2000 and higher, InnoDB makes use of the Interlocked functions. 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, and InnoDB uses these functions by default. When MySQL 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, InnoDB uses the library functions.

This change improves the scalability of InnoDB on multi-core systems. This feature is enabled out-of-the-box on the platforms where it is supported. You do not have to set any parameter or option to take advantage of the improved performance. On platforms where the GCC, Windows, or Solaris functions for atomic memory access are not available, InnoDB uses the traditional Pthreads method of implementing mutexes and read/write locks.

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

Note: If you are building from source, ensure that the build process properly takes advantage of your platform capabilities.

For more information about the performance implications of locking, see Optimizing Locking Operations.

# 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), benefit from using a more highly tuned memory allocator as opposed to the internal, InnoDB-specific memory allocator.

You can control whether InnoDB uses its own memory allocator or an allocator of the operating system, by setting the value of the 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 uses 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, set `innodb_use_sys_malloc` to `0`.

**Note**

When the InnoDB memory allocator is disabled, InnoDB ignores 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 are fulfilled by the system memory allocator.

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.

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 only includes the buffer pool statistics in the "Total memory allocated". Any memory allocated using the `mem` subsystem or using `ut_malloc` is excluded.

For more information about the performance implications of InnoDB memory usage, see Buffering and Caching.

# 7.4 Controlling InnoDB Change Buffering

When `INSERT`, `UPDATE`, and `DELETE` operations are done to a table, often the values of indexed columns (particularly the values of secondary keys) are not in sorted order, requiring substantial I/O to bring secondary indexes up to date. 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 merged when the page is loaded to the buffer pool, and the updated page is later flushed to disk using the normal mechanism. The InnoDB main thread merges buffered changes when the server is nearly idle, and during a slow shutdown.

Because it can result in fewer disk reads and writes, this feature is most valuable for workloads that are I/O-bound, for example applications with a high volume of DML operations such as bulk inserts.

However, the insert buffer occupies a part of the buffer pool, reducing the memory available to cache data pages. If the working set almost fits in the buffer pool, or if your tables have relatively few secondary indexes, it may be useful to disable insert buffering. If the working set entirely fits in the buffer pool, insert buffering does not impose any extra overhead, because it only applies to pages that are not in the buffer pool.

You can control the extent to which InnoDB performs insert buffering with the system configuration parameter `innodb_change_buffering`. You can turn on and off buffering for inserts, delete operations

(when index records are initially marked for deletion) and purge operations (when index records are physically deleted). An update operation is represented as a combination of an insert and a delete. In MySQL 5.5 and higher, the default value is changed from `inserts` to `all`.

The allowed values of `innodb_change_buffering` are:

- **all**

  The default value: buffer inserts, delete-marking operations, and purges.

- **none**

  Do not buffer any operations.

- **inserts**

  Buffer insert operations.

- **deletes**

  Buffer delete-marking operations.

- **changes**

  Buffer both inserts and delete-marking.

- **purges**

  Buffer the physical deletion operations that happen in the background.

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.

For more information about speeding up `INSERT`, `UPDATE`, and `DELETE` statements, see Optimizing DML Statements.

# 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. MySQL monitors searches on each index defined for an InnoDB table. If it notices that certain index values are being accessed frequently, it automatically builds an in-memory hash table for that index. See Adaptive Hash Indexes for background information and usage guidelines for the adaptive hash index feature and the `innodb_adaptive_hash_index` configuration option.

# 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 modern operating systems and servers with multi-core processors, where context switching is efficient, most workloads run well without any limit on the number of concurrent threads. Scalability improvements in MySQL 5.5 and up reduce the need to limit the number of concurrently executing threads inside InnoDB.

In situations where it is 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 sleeps 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 is processed. Threads waiting for locks are not counted in the number of concurrently executing threads.

You can limit the number of concurrent threads by setting the configuration parameter `innodb_thread_concurrency`. Once the number of executing threads reaches this limit, additional threads sleep for a number of microseconds, set by the 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 InnoDB. Currently, 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 before 1.0.3 | (corresponding to Plugin) | 8 | 8 | 0 |
| InnoDB 1.0.3 and newer | (corresponding to Plugin) | 0 | No limit | 0 |

Note that InnoDB causes threads to sleep only when the number of concurrent threads is limited. When there is no limit on the number of threads, all 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 a new SQL statement starts, a thread has no tickets, and it must observe `innodb_thread_concurrency`. Once the thread is entitled to enter InnoDB, it is assigned a number of tickets that it can use for subsequently entering InnoDB. If the tickets run out, `innodb_thread_concurrency` is observed again and further tickets are 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 is given one ticket once the lock becomes available.

The correct values of these variables depend on your environment and workload. Try a range of different values to determine what value works for your applications. Before limiting the number of concurrently executing threads, 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.

For general performance information about MySQL thread handling, see How MySQL Uses Threads for Client Connections.

# 7.7 Changes in the Read-Ahead Algorithm

A read-ahead request is an I/O request to prefetch multiple pages in the buffer pool asynchronously, in anticipation that these pages will be needed soon. InnoDB uses or has used two read-ahead algorithms to improve I/O performance:

**Linear** read-ahead is a technique that predicts what pages might be needed soon based on pages in the buffer pool being accessed sequentially. You control when InnoDB performs a read-ahead operation by adjusting the number of sequential page accesses required to trigger an asynchronous read request, using the configuration parameter `innodb_read_ahead_threshold`. Before this parameter was added, InnoDB would only calculate whether to issue an asynchronous prefetch request for the entire next extent when it read in the last page of the current extent.

**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 turned off in 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`.

If the number of pages read from an extent of 64 pages is greater or equal to `innodb_read_ahead_threshold`, InnoDB initiates an asynchronous read-ahead operation 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 the access pattern check. For example, if you set the value to 48, InnoDB triggers 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` can 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.

The `SHOW ENGINE INNODB STATUS` command displays statistics to help you evaluate the effectiveness of the read-ahead algorithm. See Section 8.8, "More Read-Ahead Statistics" for more information.

For more information about I/O performance, see Optimizing `InnoDB` Disk I/O and Optimizing Disk I/O.

# 7.8 Multiple Background InnoDB I/O Threads

InnoDB uses background threads to service various types of I/O requests. You can configure the number of background threads that service read and write I/O on data pages, using the configuration parameters `innodb_read_io_threads` and `innodb_write_io_threads`. These parameters signify the number of background threads used for read and write requests respectively. They are effective on all supported platforms. You can set the value of these parameters in the MySQL option file (`my.cnf` or `my.ini`); you cannot change them dynamically. The default value for these parameters is `4` and the permissible values range from `1-64`.

These parameters replace `innodb_file_io_threads` from earlier versions of MySQL. If you try to set a value for this obsolete parameter, a warning is written to the log file and the value is ignored. This parameter only applied to Windows platforms. (On non-Windows platforms, there was only one thread each for read and write.)

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 × `innodb_read_io_threads` pending read requests in `SHOW ENGINE INNODB STATUS`, you might gain by increasing the value of `innodb_read_io_threads`.

For more information about InnoDB I/O performance, see Optimizing `InnoDB` Disk I/O.

# 7.9 Asynchronous I/O on Linux

Starting in InnoDB 1.1 with MySQL 5.5, the asynchronous I/O capability that InnoDB has had on Windows systems is now available on Linux systems. (Other Unix-like systems continue to use synchronous I/O calls.) This feature improves the scalability of heavily I/O-bound systems, which typically show many pending reads/writes in the output of the command `SHOW ENGINE INNODB STATUS\G`.

Running with a large number of `InnoDB` I/O threads, and especially running multiple such instances on the same server machine, can exceed capacity limits on Linux systems. In this case, you can fix the error:

```
EAGAIN: The specified maxevents exceeds the user's limit of available events.
```

by writing a higher limit to `/proc/sys/fs/aio-max-nr`.

In general, if a problem with the asynchronous I/O subsystem in the OS prevents InnoDB from starting, set the option `innodb_use_native_aio=0` in the configuration file. This new configuration option applies to Linux systems only, and cannot be changed once the server is running.

For more information about InnoDB I/O performance, see Optimizing `InnoDB` Disk I/O.

# 7.10 Group Commit

InnoDB, like any other ACID-compliant database engine, flushes 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 issues a single write to the log file to perform 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, and works once again with MySQL 5.1 with the InnoDB Plugin, and MySQL 5.5 and higher. The introduction of support for the distributed transactions and Two Phase Commit (2PC) in MySQL 5.0 interfered with the InnoDB group commit functionality. This issue is now resolved.

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 the MySQL Enterprise Backup product with InnoDB 1.0.4** (that is, the InnoDB Plugin with MySQL 5.1) and above. 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.

Group commit is transparent; you do not need to do anything to take advantage of this significant performance improvement.

For more information about performance of `COMMIT` and other transactional operations, see Optimizing `InnoDB` Transaction Management.

# 7.11 Controlling the InnoDB Master Thread I/O Rate

The master thread in InnoDB is a thread that performs various tasks in the background. Most of these tasks are I/O related, such as flushing dirty pages from the buffer pool or writing changes from the insert buffer 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.

The parameter `innodb_io_capacity` indicates the overall I/O capacity available to InnoDB. This parameter should be set to approximately the number of I/O operations that the system can perform per second. The value depends 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. The default value is `200`, reflecting that the performance of typical modern I/O devices is higher than in the early days of MySQL. 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.

The `innodb_io_capacity` setting is a total limit for all buffer pool instances. When dirty pages are flushed, the `innodb_io_capacity` limit is divided equally among buffer pool instances. For more information, see the `innodb_io_capacity` system variable description.

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.

For more information about InnoDB I/O performance, see Optimizing `InnoDB` Disk I/O.

# 7.12 Controlling the Flushing Rate of Dirty Pages from the InnoDB Buffer Pool

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 pool, a task performed by the master thread. Currently, InnoDB aggressively flushes buffer pool pages if the percentage of dirty pages in the buffer pool exceeds `innodb_max_dirty_pages_pct`.

InnoDB 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 is to smooth overall performance by ensuring that buffer flush activity keeps up with the need to keep the buffer pool "clean". Automatically adjusting the rate of flushing can help to avoid sudden dips in throughput, when excessive buffer pool flushing limits the I/O capacity available for ordinary read and write activity.

InnoDB uses its log files in a circular fashion. Before reusing a portion of a log file, InnoDB flushes to disk all dirty buffer pool pages whose redo entries are contained in that portion of the log file, a process known as a sharp checkpoint. If a workload is write-intensive, it generates a lot of redo information, all written to the log file. If all available space in the log files is used up, a sharp checkpoint occurs, causing a temporary reduction in throughput. This situation can happen even though `innodb_max_dirty_pages_pct` is not reached.

InnoDB uses a heuristic-based algorithm to avoid such a scenario, by measuring the number of dirty pages in the buffer pool and the rate at which redo is being generated. Based on these numbers, InnoDB decides how many dirty pages to flush from the buffer pool each second. This self-adapting algorithm is able to deal with sudden changes in the workload.

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, a new configuration parameter lets you turn off this feature. The default value of the 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.

For more information about InnoDB I/O performance, see Optimizing `InnoDB` Disk I/O.

## 7.13 Using the PAUSE Instruction in InnoDB Spin Loops

Synchronization inside InnoDB frequently involves the use of spin loops: 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 performance penalty on transaction throughput. Most modern processors implement the `PAUSE` instruction for use in spin loops, so the processor can be more efficient.

InnoDB 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.

You do not have to do anything to take advantage of this performance improvement.

For performance considerations for InnoDB locking operations, see Optimizing Locking Operations.

## 7.14 Control of Spin Lock Polling

Many InnoDB mutexes and rw-locks are reserved for a short time. On a multi-core system, it can be more efficient for a thread to continuously check if it can acquire a mutex or rw-lock for a while before sleeping. If the mutex or rw-lock becomes available during this polling period, the thread can continue immediately, in the same time slice. However, too-frequent polling by multiple threads of a shared object can cause "cache ping pong", different processors invalidating portions of each others' cache. InnoDB minimizes this issue by waiting a random time between subsequent polls. The delay is implemented as a busy loop.

You can control the maximum delay between testing a mutex or rw-lock using the parameter `innodb_spin_wait_delay`. The duration of the delay loop depends on the C compiler and the target processor. (In the 100MHz Pentium era, the unit of delay was one microsecond.) On a system where all processor cores share a fast cache memory, you might reduce the maximum delay or disable the busy loop altogether by setting `innodb_spin_wait_delay=0`. On a system with multiple processor chips, the effect of cache invalidation can be more significant and you might increase the maximum delay.

The default value of `innodb_spin_wait_delay` is `6`. The spin wait delay is a dynamic global parameter that you can specify in the MySQL option file (`my.cnf` or `my.ini`) or change 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.

For performance considerations for InnoDB locking operations, see Optimizing Locking Operations.

## 7.15 Making the Buffer Pool Scan Resistant

Rather than using a strictly LRU algorithm, InnoDB uses a technique to minimize the amount of data that is brought into the buffer pool and never accessed again. The goal is to make sure that frequently accessed

("hot") pages remain in the buffer pool, even as read-ahead and full table scans bring in new blocks that might or might not be accessed afterward.

Newly read blocks are inserted into the middle of the list representing the buffer pool. of the LRU list. All newly read pages are inserted at a location that by default is `3/8` from the tail of the LRU list. The pages are moved to the front of the list (the most-recently used end) when they are accessed in the buffer pool for the first time. Thus pages that are never accessed never make it to the front portion of the LRU list, and "age out" sooner than with a strict LRU approach. This arrangement divides the LRU list into two segments, where the pages downstream of the insertion point are considered "old" and are desirable victims for LRU eviction.

For an explanation of the inner workings of the InnoDB buffer pool and the specifics of its LRU replacement algorithm, see The `InnoDB` Buffer Pool.

You can control the insertion point in the LRU list, and choose whether InnoDB applies the same optimization to blocks brought into the buffer pool by table or index scans. The configuration parameter `innodb_old_blocks_pct` 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 value range is `5` (new pages in the buffer pool age out very quickly) to `95` (only 5% of the buffer pool is reserved for hot pages, making the algorithm close to the familiar LRU strategy).

The optimization that keeps the buffer pool from being churned by read-ahead can avoid similar problems due to table or index scans. In these scans, a data page is typically accessed a few times in quick succession and is never touched again. The configuration parameter `innodb_old_blocks_time` 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 most-recently used end of the buffer pool list when it is first accessed in the buffer pool. Increasing this value makes more and more blocks likely to age out faster from the buffer pool.

Both `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, the `SHOW ENGINE INNODB STATUS` command reports additional statistics. The `BUFFER POOL AND MEMORY` section 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 pool 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.

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 pool.

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.

For more information about the InnoDB buffer pool, see The `InnoDB` Buffer Pool.

# 7.16 Improvements to Crash Recovery Performance

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 than in MySQL 5.1 and earlier, due to improved algorithms for memory management. You do not need to take any actions to take advantage of this performance enhancement. If you kept the size of your redo log files artificially low because recovery took a long time, you can consider increasing the file size.

For more information about InnoDB recovery, see The `InnoDB` Recovery Process.

# 7.17 Integration with the MySQL Performance Schema

Starting with InnoDB 1.1 with MySQL 5.5, you can profile certain internal InnoDB operations using the MySQL Performance Schema feature. This type of tuning is primarily for expert users, those who push the limits of MySQL performance, read the MySQL source code, and evaluate optimization strategies to overcome performance bottlenecks. DBAs can also use this feature for capacity planning, to see whether their typical workload encounters any performance bottlenecks with a particular combination of CPU, RAM, and disk storage; and if so, to judge whether performance can be improved by increasing the capacity of some part of the system.

To use this feature to examine InnoDB performance:

- You must be running MySQL 5.5 or higher. You must build the database server from source, enabling the Performance Schema feature by building with the `--with-perfschema` option. Since the Performance Schema feature introduces some performance overhead, you should use it on a test or development system rather than on a production system.

- You must be running InnoDB 1.1 or higher.

- You must be generally familiar with how to use the Performance Schema feature, for example to query tables in the `performance_schema` database.

- Examine the following kinds of InnoDB objects by querying the appropriate `performance_schema` tables. The items associated with InnoDB all contain the substring `innodb` in the `EVENT_NAME` column.

  For the definitions of the `*_instances` tables, see Performance Schema Instance Tables. For the definitions of the `*_summary_*` tables, see Performance Schema Summary Tables. For the definition of the `thread` table, see Performance Schema Miscellaneous Tables. For the definition of the `*_current_*` and `*_history_*` tables, see Performance Schema Wait Event Tables.

  - Mutexes in the `mutex_instances` table. (Mutexes and RW-locks related to the `InnoDB` buffer pool are not included in this coverage; the same applies to the output of the `SHOW ENGINE INNODB MUTEX` command.)

  - RW-locks in the `rwlock_instances` table.

  - RW-locks in the `rwlock_instances` table.

  - File I/O operations in the `file_instances`, `file_summary_by_event_name`, and `file_summary_by_instance` tables.

  - Threads in the `PROCESSLIST` table.

- During performance testing, examine the performance data in the `events_waits_current` and `events_waits_history_long` tables. If you are interested especially in InnoDB-related objects, use the clause `WHERE EVENT_NAME LIKE '%innodb%'` to see just those entries; otherwise, examine the performance statistics for the overall MySQL server.

- You must be running MySQL 5.5, with the Performance Schema enabled by building with the `--with-perfschema` build option.

  For more information about the MySQL Performance Schema, see MySQL Performance Schema.

# 7.18 Improvements to Performance from Multiple Buffer Pools

This performance enhancement is primarily useful for people with a large buffer pool size, typically in the multi-gigabyte range. To take advantage of this speedup, you must set the new `innodb_buffer_pool_instances` configuration option, and you might also adjust the `innodb_buffer_pool_size` value.

When the InnoDB buffer pool is large, many data requests can be satisfied by retrieving from memory. You might encounter bottlenecks from multiple threads trying to access the buffer pool at once. Starting in InnoDB 1.1 and MySQL 5.5, you can enable multiple buffer pools to minimize this contention. Each page that is stored in or read from the buffer pool is assigned to one of the buffer pools randomly, using a hashing function. Each buffer pool manages its own free lists, flush lists, LRUs, and all other data structures connected to a buffer pool, and is protected by its own buffer pool mutex.

To enable this feature, set the `innodb_buffer_pool_instances` configuration option to a value greater than 1 (the default) up to 64 (the maximum). This option takes effect only when you set the

`innodb_buffer_pool_size` to a size of 1 gigabyte or more. The total size you specify is divided among all the buffer pools. For best efficiency, specify a combination of `innodb_buffer_pool_instances` and `innodb_buffer_pool_size` so that each buffer pool instance is at least 1 gigabyte.

For more information about the InnoDB buffer pool, see The `InnoDB` Buffer Pool.

## 7.19 Better Scalability with Multiple Rollback Segments

Starting in InnoDB 1.1 with MySQL 5.5, the limit on concurrent transactions is greatly expanded, removing a bottleneck with the InnoDB rollback segment that affected high-capacity systems. The limit applies to concurrent transactions that change any data; read-only transactions do not count against that maximum.

The single rollback segment is now divided into 128 segments, each of which can support up to 1023 transactions that perform writes, for a total of approximately 128K concurrent transactions. The original transaction limit was 1023.

Each transaction is assigned to one of the rollback segments, and remains tied to that rollback segment for the duration. This enhancement improves both scalability (higher number of concurrent transactions) and performance (less contention when different transactions access the rollback segments).

To take advantage of this feature, you do not need to create any new database or tables, or reconfigure anything. You must do a slow shutdown before upgrading from MySQL 5.1 or earlier, or some time afterward. InnoDB makes the required changes inside the system tablespace automatically, the first time you restart after performing a slow shutdown.

For more information about performance of InnoDB under high transactional load, see Optimizing `InnoDB` Transaction Management.

## 7.20 Better Scalability with Improved Purge Scheduling

Starting in InnoDB 1.1 with MySQL 5.5, the purge operations (a type of garbage collection) that InnoDB performs automatically can be done in a separate thread, rather than as part of the master thread. This change improves scalability, because the main database operations run independently from maintenance work happening in the background.

To enable this feature, set the configuration option `innodb_purge_threads=1`, as opposed to the default of 0, which combines the purge operation into the master thread.

You might not notice a significant speedup, because the purge thread might encounter new types of contention; the single purge thread really lays the groundwork for further tuning and possibly multiple purge threads in the future. There is another new configuration option, `innodb_purge_batch_size` with a default of 20 and maximum of 5000. This option is mainly intended for experimentation and tuning of purge operations, and should not be interesting to typical users.

For more information about InnoDB I/O performance, see Optimizing `InnoDB` Disk I/O.

## 7.21 Improved Log Sys Mutex

This is another performance improvement that comes for free, with no user action or configuration needed. The details here are intended for performance experts who delve into the InnoDB source code, or interpret reports with keywords such as "mutex" and "log_sys".

The mutex known as the log sys mutex has historically done double duty, controlling access to internal data structures related to log records and the LSN, as well as pages in the buffer pool that are changed when a mini-transaction is committed. Starting in InnoDB 1.1 with MySQL 5.5, these two kinds of

operations are protected by separate mutexes, with a new `log_buf` mutex controlling writes to buffer pool pages due to mini-transactions.

For performance considerations for InnoDB locking operations, see Optimizing Locking Operations.

# 7.22 Separate Flush List Mutex

Starting with InnoDB 1.1 with MySQL 5.5, concurrent access to the buffer pool is faster. Operations involving the flush list, a data structure related to the buffer pool, are now controlled by a separate mutex and do not block access to the buffer pool. You do not need to configure anything to take advantage of this speedup; it is fully automatic.

For more information about the InnoDB buffer pool, see The `InnoDB` Buffer Pool.

# Chapter 8 `InnoDB` Features for Flexibility, Ease of Use and Reliability

## Table of Contents

This section describes several recently added InnoDB features that offer new flexibility and improve ease of use, reliability and performance. The Barracuda file format improves efficiency for storing large variable-length columns, and enables table compression. Configuration options that once were unchangeable after startup, are now flexible and can be changed dynamically. Some improvements are automatic, such as faster and more efficient `TRUNCATE TABLE`. Others allow you the flexibility to control InnoDB behavior; for example, you can control whether certain problems cause errors or just warnings. And informational messages and error reporting continue to be made more user-friendly.

## 8.1 The Barracuda File Format

InnoDB has started using named file formats to improve compatibility in upgrade and downgrade situations, or heterogeneous systems running different levels of MySQL. Many important InnoDB features, such as table compression and the `DYNAMIC` row format for more efficient BLOB storage, require creating tables in the Barracuda file format. The original file format, which previously didn't have a name, is known now as Antelope.

To create new tables that take advantage of the Barracuda features, enable that file format using the configuration parameter `innodb_file_format`. The value of this parameter determines whether a newly created table or index can use compression or the new `DYNAMIC` row format.

To preclude the use of new features that would make your database inaccessible to the built-in InnoDB in MySQL 5.1 and prior releases, omit `innodb_file_format` or set it to Antelope.

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` statement.

For more information about managing file formats, see Chapter 4, *InnoDB File-Format Management*.

## 8.2 Dynamic Control of System Configuration Parameters

In MySQL 5.5 and higher, you can change certain system configuration parameters without shutting down and restarting the server, as was necessary in MySQL 5.1 and lower. This increases uptime, and

makes it easier to test and prototype new SQL and application code. The following sections explain these parameters.

## 8.2.1 Dynamically Changing `innodb_file_per_table`

Since MySQL version 4.1, InnoDB has provided two alternatives for how tables are stored on disk. You can create a new table and its indexes in the shared system tablespace, physically stored in the ibdata files. Or, you can store a new table and its indexes in a separate tablespace (a .ibd file). The storage layout for each InnoDB table is determined by the the configuration parameter `innodb_file_per_table` at the time the table is created.

In MySQL 5.5 and higher, the configuration parameter `innodb_file_per_table` is dynamic, and can be set `ON` or `OFF` using the `SET GLOBAL`. Previously, the only way to set this parameter was in the MySQL configuration file (`my.cnf` or `my.ini`), and changing it required shutting down and restarting the server.

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 enabled can use the Barracuda file format, and `TRUNCATE` returns the disk space for those tables to the operating system. The Barracuda file format in turn enables features such as table compression and the `DYNAMIC` row format. Tables created when `innodb_file_per_table` is off cannot use these features. To take advantage of those features for an existing table, you can turn on the file-per-table setting and run `ALTER TABLE t ENGINE=INNODB` for that table.

When you redefine the primary key for an InnoDB table, the table is re-created using the current settings for `innodb_file_per_table` and `innodb_file_format`. This behavior does not apply when adding or dropping InnoDB secondary indexes, as explained in Chapter 2, *Fast Index Creation in the InnoDB Storage Engine*. When a secondary index is created without rebuilding the table, the index is stored in the same file as the table data, regardless of the current `innodb_file_per_table` setting.

## 8.2.2 Dynamically Changing `innodb_stats_on_metadata`

In MySQL 5.5 and higher, you can change the setting of `innodb_stats_on_metadata` dynamically at runtime, to control whether or not InnoDB performs statistics gathering when metadata statements are executed. To change the setting, issue the statement `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.

This setting is related to the feature described in Section 8.5, "Controlling Optimizer Statistics Estimation".

## 8.2.3 Dynamically Changing `innodb_lock_wait_timeout`

The length of time a transaction waits for a resource, before giving up and rolling back the statement, is determined by the value of the configuration parameter `innodb_lock_wait_timeout`. (In MySQL 5.0.12 and earlier, the entire transaction was rolled back, not just the statement.) Your application can try the statement again (usually after waiting for a while), or roll back the entire transaction and restart.

The error returned when the timeout period is exceeded is:

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

In MySQL 5.5 and higher, the configuration parameter `innodb_lock_wait_timeout` can be set at runtime with the `SET GLOBAL` or `SET SESSION` statement. 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.

In MySQL 5.1 and earlier, the only way to set this parameter was in the MySQL configuration file (`my.cnf` or `my.ini`), and changing it required shutting down and restarting the server.

## 8.2.4 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.

The configuration option `innodb_adaptive_hash_index` lets you disable the adaptive hash index. It is enabled by default. You can modify this parameter through the `SET GLOBAL` statement, without restarting the server. Changing the setting requires the `SUPER` privilege.

Disabling the adaptive hash index empties the hash table immediately. Normal operations can continue while the hash table is emptied, and executing queries that were using the hash table access the index B-trees directly instead. When the adaptive hash index is re-enabled, the hash table is populated again during normal operation.

## 8.3 `TRUNCATE TABLE` Reclaims Space

When you truncate a table that is stored in a `.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 table is dropped and re-created in a new `.ibd` file. This operation is much faster than deleting the rows one by one. The operating system can reuse the disk space, in contrast to tables within the InnoDB system tablespace, where only InnoDB can use the space after they are truncated. Physical backups can also be smaller, without big blocks of unused space in the middle of the system tablespace.

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 TABLE` statement is an arbitrary number.

> **Note**
>
> If there is a foreign key constraint between two columns in the same table, that table can still be truncated using this fast technique.
>
> If there are foreign key constraints between the table being truncated and other tables, the truncate operation fails. This is a change to the previous behavior, which would transform the `TRUNCATE` operation to a `DELETE` operation that removed all the rows and triggered `ON DELETE` operations on child tables.

## 8.4 `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 statements, InnoDB provides a strict mode of operations. In this mode, InnoDB raises error conditions in certain cases, rather than issuing a warning and processing the specified statement (perhaps with unintended behavior). This is analogous to `sql_mode` in MySQL, which controls what SQL syntax MySQL accepts, and determines whether it silently ignores errors, or validates input syntax and data values. Since InnoDB strict mode is relatively new, some statements that execute without errors with earlier versions of MySQL might generate errors unless you disable strict mode.

The setting of InnoDB strict mode affects the handling of syntax errors on the `CREATE TABLE`, `ALTER TABLE` and `CREATE INDEX` statements. The strict mode also enables a record size check, so that an `INSERT` or `UPDATE` never fails due to the record being too large for the selected page size.

Oracle recommends enabling `innodb_strict_mode` when using the `ROW_FORMAT` and `KEY_BLOCK_SIZE` clauses on `CREATE TABLE`, `ALTER TABLE`, and `CREATE INDEX` statements. Without strict mode, InnoDB ignores conflicting clauses and creates the table or index, with only a warning in the message log. The resulting table might have different behavior than you intended, such as having no compression when you tried to create a compressed table. When InnoDB strict mode is on, such problems generate an immediate error and the table or index is not created, avoiding a troubleshooting session later.

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` or `my.ini`. You can also enable or disable InnoDB strict mode at runtime with the statement `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`, and the setting affects only that client.

# 8.5 Controlling Optimizer Statistics Estimation

The MySQL query optimizer uses estimated statistics about key distributions to choose the indexes for an execution plan, based on the relative selectivity of the index. Certain operations cause InnoDB to sample random pages from each index on a table to estimate the cardinality of the index. (This technique is known as random dives.) These operations include the `ANALYZE TABLE` statement, the `SHOW TABLE STATUS` statement, and accessing the table for the first time after a restart.

To give you control over the quality of the statistics estimate (and thus better information for the query optimizer), you can now change the number of sampled pages using the parameter `innodb_stats_sample_pages`. Previously, the number of sampled pages was always 8, which could be insufficient to produce an accurate estimate, leading to poor index choices by the query optimizer. This technique is especially important for large tables and tables used in joins. Unnecessary full table scans for such tables can be a substantial performance issue.

You can set the global parameter `innodb_stats_sample_pages`, at runtime. The default value for this parameter is 8, preserving the same behavior as in past releases.

**Note**

The value of `innodb_stats_sample_pages` affects the index sampling for *all* tables and indexes. 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.

- Increasing the `innodb_stats_sample_pages` value might require more disk reads. 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 might choose very different query plans based on different estimates of index selectivity.

To disable the cardinality estimation for metadata statements such as `SHOW TABLE STATUS`, execute the statement `SET GLOBAL innodb_stats_on_metadata=OFF` (or `0`). The ability to set this option dynamically is also relatively new.

All InnoDB tables are opened, and the statistics are re-estimated for all associated indexes, when the `mysql` client starts if the auto-rehash setting is set on (the default). To improve the start up time of the `mysql` client, you can turn auto-rehash off. The auto-rehash feature enables automatic name completion of database, table, and column names for interactive users.

Whatever value of `innodb_stats_sample_pages` works best for a system, set the option and leave it at that value. Choose a value that results in reasonably accurate estimates for all tables in your database without requiring excessive I/O. Because the statistics are automatically recalculated at various times other than on execution of `ANALYZE TABLE`, it does not make sense to increase the index sample size, run `ANALYZE TABLE`, then decrease sample size again. The more accurate statistics calculated by `ANALYZE` running with a high value of `innodb_stats_sample_pages` can be wiped away later.

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

# 8.6 Better Error Handling when Dropping Indexes

For optimal performance with DML statements, InnoDB requires an index to exist on foreign key columns, so that `UPDATE` and `DELETE` operations on a parent table can easily check whether corresponding rows exist in the child table. MySQL creates or drops such indexes automatically when needed, as a side-effect of `CREATE TABLE`, `CREATE INDEX`, and `ALTER TABLE` statements.

When you drop an index, InnoDB checks whether the index is not used for checking a foreign key constraint. It is still OK to drop the index if there is another index that can be used to enforce the same constraint. InnoDB prevents you from dropping the last index that can enforce a particular referential constraint.

The message that reports this error condition is:

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

This message is friendlier than the earlier message it replaces:

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

A similar change in error reporting applies to an attempt to drop the primary key index. For tables without an explicit `PRIMARY KEY`, InnoDB creates an implicit clustered index using the first columns of the table that are declared `UNIQUE` and `NOT NULL`. When you drop such an index, InnoDB automatically copies the table and rebuilds the index using a different `UNIQUE NOT NULL` group of columns or a system-generated key. Since this operation changes the primary key, it uses the slow method of copying the table and re-creating the index, rather than the Fast Index Creation technique from Section 2.3, "Implementation Details of Fast Index Creation".

Previously, an attempt to drop an implicit clustered index (the first `UNIQUE NOT NULL` index) failed if the table did not contain a `PRIMARY KEY`:

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

# 8.7 More Compact Output of `SHOW ENGINE INNODB MUTEX`

The statement `SHOW ENGINE INNODB MUTEX` displays information about InnoDB mutexes and rw-locks. Although this information is useful for tuning on multi-core systems, the amount of output can be

overwhelming on systems with a big buffer pool. There is one mutex and one rw-lock in each 16K buffer pool block, and there are 65,536 blocks per gigabyte. It is unlikely that a single block mutex or rw-lock from the buffer pool could become a performance bottleneck.

`SHOW ENGINE INNODB MUTEX` now skips the mutexes and rw-locks of buffer pool blocks. It also does not list any mutexes or rw-locks that have never been waited on (`os_waits=0`). Thus, `SHOW ENGINE INNODB MUTEX` only displays information about mutexes and rw-locks outside of the buffer pool that have caused at least one OS-level wait.

# 8.8 More Read-Ahead Statistics

As described in Section 7.7, "Changes in the Read-Ahead Algorithm", a read-ahead request is an asynchronous I/O request issued in anticipation that a page will be used in the near future. Knowing how many pages are read through this read-ahead mechanism, and how many of them are evicted from the buffer pool without ever being accessed, can be useful to help fine-tune the parameter `innodb_read_ahead_threshold`.

`SHOW ENGINE INNODB STATUS` output displays the global status variables `Innodb_buffer_pool_read_ahead` and `Innodb_buffer_pool_read_ahead_evicted`. These variables indicate the number of pages brought into the buffer pool by read-ahead requests, and the number of such pages evicted from the buffer pool without ever being accessed respectively. These counters provide global values since the last server restart.

`SHOW ENGINE INNODB STATUS` also shows the rate at which the read-ahead pages are read in and the rate at which such pages are evicted without being accessed. The per-second averages are based on the statistics collected since the last invocation of `SHOW ENGINE INNODB STATUS` and are displayed in the `BUFFER POOL AND MEMORY` section of the output.

Since the InnoDB read-ahead mechanism has been simplified to remove random read-ahead, the status variables `Innodb_buffer_pool_read_ahead_rnd` and `Innodb_buffer_pool_read_ahead_seq` are no longer part of the `SHOW ENGINE INNODB STATUS` output.

# Chapter 9 Installing the InnoDB Storage Engine

When you use the InnoDB storage engine 1.1 and above, with MySQL 5.5 and above, you do not need to do anything special to install: everything comes configured as part of the MySQL source and binary distributions. This is a change from earlier releases of the InnoDB Plugin, where you were required to match up MySQL and InnoDB version numbers and update your build and configuration processes.

The InnoDB storage engine 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 storage engine; it is not available from the InnoDB web site.

If you used any scripts or configuration files with the earlier InnoDB storage engine 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.

Because the InnoDB storage engine has now replaced the built-in InnoDB, you no longer need to specify options like `--ignore-builtin-innodb` and `--plugin-load` during startup.

To take best advantage of current InnoDB features, we recommend specifying the following options in your configuration file:

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

For information about these new features, see Section 8.2.1, "Dynamically Changing `innodb_file_per_table`", Chapter 4, *InnoDB File-Format Management*, and Section 8.4, "`InnoDB` Strict Mode". You might need to continue to use the previous values for these parameters in some replication and similar configurations involving both new and older versions of MySQL.

# Chapter 10 Upgrading the InnoDB Storage Engine

Prior to MySQL 5.5, some upgrade scenarios involved upgrading the separate instance of InnoDB known as the InnoDB Plugin. In MySQL 5.5 and higher, the features of the InnoDB Plugin have been folded back into built-in InnoDB, so the upgrade procedure for InnoDB is the same as the one for the MySQL server. For details, see Upgrading MySQL.

# Chapter 11 Downgrading the InnoDB Storage Engine

## Table of Contents

## 11.1 Overview

Prior to MySQL 5.5, some downgrade scenarios involved switching the separate instance of InnoDB known as the InnoDB Plugin back to the built-in InnoDB storage engine. In MySQL 5.5 and higher, the features of the InnoDB Plugin have been folded back into built-in InnoDB, so the downgrade procedure for InnoDB is the same as the one for the MySQL server. For details, see Downgrading MySQL.

# Chapter 12 InnoDB Storage Engine Change History

## Table of Contents

## 12.1 Changes in InnoDB Storage Engine 1.x

Since InnoDB 1.1 is tightly integrated with MySQL 5.5, for changes after the initial InnoDB 1.1 release, see the MySQL 5.5 Release Notes.

## 12.2 Changes in InnoDB Storage Engine 1.1 (April 13, 2010)

For an overview of the changes, see this introduction article for MySQL 5.5 with InnoDB 1.1. The following is a condensed version of the change log.

Fix for Bug #52580: Crash in ha_innobase::open on executing INSERT with concurrent ALTER TABLE.

Change in MySQL Bug #51557 releases the mutex LOCK_open before ha_innobase::open(), causing racing condition for index translation table creation. Fix it by adding dict_sys mutex for the operation.

Add support for multiple buffer pools.

Fix Bug #26590: MySQL does not allow more than 1023 open transactions. Create additional rollback segments on startup. Reduce the upper limit of total rollback segments from 256 to 128. This is because we can't use the sign bit. It has not caused problems in the past because we only created one segment. InnoDB has always had the capability to use the additional rollback segments, therefore this patch is backward compatible. The only requirement to maintain backward compatibility has been to ensure that the additional segments are created after the double write buffer. This is to avoid breaking assumptions in the existing code.

Implement Performance Schema in InnoDB. Objects in four different modules in InnoDB have been performance instrumented, these modules are: mutexes, rwlocks, file I/O, and threads. We mostly preserved the existing APIs, but APIs would point to instrumented function wrappers if performance schema is defined. There are 4 different defines that control the instrumentation of each module. The feature is off by default, and will be compiled in with a special build option, and require a configure option to turn it on when the server boots.

Implement the buf_pool_watch for DeleteBuffering in the page hash table. This serves two purposes. It allows multiple watches to be set at the same time (by multiple purge threads) and it removes a race condition when the read of a block completes about the time the buffer pool watch is being set.

Introduce a new mutex to protect flush_list. Redesign mtr_commit() in a way that log_sys mutex is not held while all mtr_memos are popped and is released just after the modified blocks are inserted into the flush_list. This should reduce contention on log_sys mutex.

Implement the global variable `innodb_change_buffering`, with the following values:

- `none`: buffer nothing

- `inserts`: buffer inserts (like InnoDB so far)

- `deletes`: buffer delete-marks

- `changes`: buffer inserts and delete-marks

- `purges`: buffer delete-marks and deletes

- `all`: buffer all operations (insert, delete-mark, delete)

The default is `all`. All values except `none` and `inserts` will make InnoDB write new-format records to the insert buffer, even for inserts.

Provide support for native AIO on Linux.

## 12.3 Changes in InnoDB Plugin 1.0.x

The InnoDB 1.0.x releases that accompany MySQL 5.1 have their own change history. Changes up to InnoDB 1.0.8 are listed at http://dev.mysql.com/doc/innodb-plugin/1.0/en/innodb-changes.html. Changes from InnoDB 1.0.9 and up are incorporated in the MySQL 5.1 Release Notes.

# Appendix A Third-Party Software

## Table of Contents

Oracle acknowledges that certain Third Party and Open Source software has been used to develop or is incorporated in the InnoDB storage engine. This appendix includes required third-party license information.

## A.1 Performance Patches from Google

Oracle 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 7.11, "Controlling the InnoDB 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

Oracle gratefully acknowledges the contribution of Percona, Inc. to improve InnoDB performance by implementing configurable background threads, as discussed in Section 7.8, "Multiple Background InnoDB 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

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

- Introducing the PAUSE instruction inside spin loops, as discussed in Section 7.13, "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`, `ut0ut.c`, and `ut0ut.h`.

This contribution is incorporated subject to the conditions contained in the file `COPYING.Sun_Microsystems`, which are reproduced here.

# Appendix B `InnoDB` Parameter Changes

## Table of Contents

## B.1 New Parameters

### New InnoDB Plugin and InnoDB 1.1 Parameters

Throughout the course of development of InnoDB 1.1 and its predecessor, the InnoDB Plugin, the following configuration parameters were introduced. The `InnoDB` storage engine version in which each parameter was introduced is provided (where applicable). Starting in MySQL 5.5.30, separate version numbering for the `InnoDB` storage engine was discontinued in favor of MySQL server `version` numbering. See `InnoDB` Startup Options and System Variables for parameter descriptions and the MySQL release in which parameters were added.

- `innodb_adaptive_flushing`: introduced in InnoDB storage engine 1.0.4

- `innodb_buffer_pool_instances`

- `innodb_change_buffering`: introduced in InnoDB storage engine 1.0.3

- `innodb_file_format`: introduced in InnoDB storage engine 1.0.1

- `innodb_file_format_check`: introduced in InnoDB storage engine 1.0.1

- `innodb_file_format_max`

- `innodb_io_capacity`: introduced in InnoDB storage engine 1.0.4

- `innodb_old_blocks_pct`: introduced in InnoDB storage engine 1.0.5

- `innodb_old_blocks_time`: introduced in InnoDB storage engine 1.0.5

- `innodb_purge_batch_size`

- `innodb_purge_threads`

- `innodb_read_ahead_threshold`: introduced in InnoDB storage engine 1.0.4

- `innodb_read_io_threads`: introduced in InnoDB storage engine 1.0.4

- `innodb_spin_wait_delay`: introduced in InnoDB storage engine 1.0.4

- `innodb_stats_sample_pages`: introduced in InnoDB storage engine 1.0.2

- `innodb_strict_mode`: introduced in InnoDB storage engine 1.0.2

- `innodb_use_native_aio`

- `innodb_use_sys_malloc`: introduced in InnoDB storage engine 1.0.3

- `innodb_write_io_threads`: introduced in InnoDB storage engine 1.0.4

# New `InnoDB` Parameters in MySQL 5.5

The following `InnoDB` configuration parameters were added in MySQL 5.5. Some of the parameters listed below may have also been added to other MySQL versions. See `InnoDB` Startup Options and System Variables for parameter descriptions and the specific release in which parameters were added.

- `innodb_buffer_pool_instances`

- `innodb_file_format_max`

- `innodb_force_load_corrupted`

- `innodb_large_prefix`

- `innodb_print_all_deadlocks`

- `innodb_purge_batch_size`

- `innodb_purge_threads`

- `innodb_random_read_ahead`

- `innodb_rollback_segments`

- `innodb_stats_method`

- `innodb_stats_on_metadata`

- `innodb_use_native_aio`

# B.2 Deprecated Parameters

- `innodb_file_io_threads`

  This parameter was deprecated in InnoDB storage engine 1.0.4 and replaced by `innodb_read_io_threads` and `innodb_write_io_threads`. See Section 7.8, "Multiple Background InnoDB I/O Threads" for additional information.

- `ignore_builtin_innodb`: removed in MySQL 5.5.22.

# B.3 Parameters with New Defaults

For better out-of-the-box performance, default values for the following `InnoDB` configuration parameters were changed in MySQL 5.5. See `InnoDB` Startup Options and System Variables for information about new parameter defaults and the specific release that changes took effect.

- `innodb_change_buffering`

- `innodb_file_format`

- `innodb_file_format_check`

- `innodb_file_per_table`

# 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 5.4, "`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, *Working with `InnoDB` Compressed Tables* for usage details.

For additional information about `InnoDB COMPRESSED` row format, see Section 5.3, "`DYNAMIC` and `COMPRESSED` Row Formats".
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, *Working with `InnoDB` Compressed Tables* 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` and `COMPRESSED` Row Formats".
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 and 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 create 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 are 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 InnoDB File-Per-Table Mode.

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, ilist.

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 My SQL 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 5.4, "`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 if 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 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 **SSD**s.
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 InnoDB File-
    Per-Table Mode 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, 73
.ARZ file, 73
.cfg file, 78
.frm file, 89
.ibd file, 94
.ibz file, 94
.isl file, 94
.MRG file, 102
.MYD file, 103
.MYI file, 103
.OPT file, 107
.PAR file, 109
.TRG file, 123
.TRN file, 123

## A

ACID, 73
adaptive flushing, 73
adaptive hash index, 42, 55, 73
AHI, 74
AIO, 74
ALTER TABLE
    ROW_FORMAT, 27
Antelope, 74
Antelope file format, 21, 28
application programming interface (API), 74
apply, 74
asynchronous I/O, 45, 74
atomic, 75
atomic instruction, 75
auto-increment, 75
auto-increment locking, 75
autocommit, 75
availability, 75

## B

B-tree, 75
background threads
    master, 46, 46
    read, 44
    write, 44
backticks, 76
backup, 76
Barracuda, 76
Barracuda file format, 9, 21, 27, 53
beta, 76
binary log, 76
binlog, 77
blind query expansion, 77

bottleneck, 77
bounce, 77
buddy allocator, 29, 77
buffer, 77
buffer pool, 47, 77
    and compressed tables, 16
buffer pool instance, 78
built-in, 78
business rules, 78

## C

cache, 79
cardinality, 79
change buffer, 79
change buffering, 79
    disabling, 41
checkpoint, 80
checksum, 80
child table, 80
clean page, 80
clean shutdown, 80
client, 80
clustered index, 80
cold backup, 81
column, 81
column index, 81
column prefix, 81
commit, 81
compact row format, 28, 81
composite index, 81
compressed backup, 82
compressed row format, 27, 82
compression, 9, 82
    algorithms, 14
    application and schema design, 12
    BLOBs, VARCHAR and TEXT, 16
    buffer pool considerations, 16
    compressed page size, 13
    configuration characteristics, 13
    data and indexes, 15
    data characteristics, 11
    enabling for a table, 9
    implementation, 14
    information schema, 29, 29
    innodb_strict_mode, 55
    KEY_BLOCK_SIZE, 13
    log file format, 17
    modification log, 15
    monitoring, 14
    overflow pages, 16
    overview, 9
    tuning, 11
    workload characteristics, 13

# J

## K

KEY_BLOCK_SIZE, 9, 13, 100

## L

latch, 100
list, 100
lock, 100
lock escalation, 100
lock mode, 100
lock wait timeout, 54
locking, 101
    information schema, 29, 31, 36
locking read, 101
log, 101
log buf mutex, 51
log buffer, 101
log file, 101
log group, 101
log sys mutex, 51
logical, 101
logical backup, 102
loose_, 102
low-water mark, 102
LRU, 102
LRU page replacement, 47
LSN, 102

## M

master server, 103
master thread, 103
MDL, 103
memcached, 103
memory allocator
    innodb_use_sys_malloc, 40
merge, 103
metadata lock, 103
metrics counter, 104
midpoint insertion, 47
midpoint insertion strategy, 104
mini-transaction, 104
mixed-mode insert, 104
multi-core, 104
multiple buffer pools, 50
multiple rollback segments, 51
mutex, 104
MVCC, 104
my.cnf, 105
my.ini, 105
mysql, 105
MySQL Enterprise Backup, 105
mysqlbackup command, 105
mysqld, 105
mysqldump, 105

## N

natural key, 105
neighbor page, 106
next-key lock, 106
non-blocking I/O, 106
non-locking read, 106
non-repeatable read, 106
normalized, 106
NoSQL, 107
NOT NULL constraint, 107
NULL, 107

## O

off-page column, 107
OLTP, 108
online, 108
online DDL, 108
optimistic, 108
optimizer, 108
Optimizer Statistics Estimation, 54, 56
option, 109
option file, 109
overflow page, 109

## P

page, 109
page cleaner, 109
page size, 109
parent table, 110
partial backup, 110
partial index, 110
Performance Schema, 49, 110
persistent statistics, 110
pessimistic, 110
phantom, 110
physical, 110
physical backup, 111
PITR, 111
plan stability, 111
plugin, 111
point-in-time recovery, 111
prepared backup, 111
primary key, 111
process, 112
PROCESSLIST
    possible inconsistency with INFORMATION_SCHEMA
    tables, 37
pseudo-record, 112
Pthreads, 112
purge, 112
purge buffering, 112
purge lag, 112
purge scheduling, 51

two-phase commit, 126

# U

undo, 126
undo log, 126
undo tablespace, 126
unique constraint, 126
unique index, 126
unique key, 127

# V

victim, 127

# W

wait, 127
warm backup, 127
warm up, 127
Windows, 127
workload, 128
write combining, 128

# X

XA, 128

# Y

young, 128