

...

Privileges belonging to a specific user are displayed by the `SHOW GRANTS` statement. See [Section 13.7.7.21, “SHOW GRANTS Statement”](#), for more information.

13.7.7.27 SHOW PROCEDURE CODE Statement

```
SHOW PROCEDURE CODE proc_name
```

This statement is a MySQL extension that is available only for servers that have been built with debugging support. It displays a representation of the internal implementation of the named stored procedure. A similar statement, `SHOW FUNCTION CODE`, displays information about stored functions (see [Section 13.7.7.19, “SHOW FUNCTION CODE Statement”](#)).

To use either statement, you must be the user named as the routine `DEFINER`, have the `SHOW_ROUTINE` privilege, or have the `SELECT` privilege at the global level.

If the named routine is available, each statement produces a result set. Each row in the result set corresponds to one “instruction” in the routine. The first column is `Pos`, which is an ordinal number beginning with 0. The second column is `Instruction`, which contains an SQL statement (usually changed from the original source), or a directive which has meaning only to the stored-routine handler.

```
mysql> DELIMITER //
mysql> CREATE PROCEDURE p1 ()
  BEGIN
    DECLARE fanta INT DEFAULT 55;
    DROP TABLE t2;
    LOOP
      INSERT INTO t3 VALUES (fanta);
    END LOOP;
  END//
```

Query OK, 0 rows affected (0.01 sec)

```
mysql> SHOW PROCEDURE CODE p1//
```

Pos	Instruction
0	set fanta@0 55
1	stmt 9 "DROP TABLE t2"
2	stmt 5 "INSERT INTO t3 VALUES (fanta)"
3	jump 2

4 rows in set (0.00 sec)

```
mysql> CREATE FUNCTION test.hello (s CHAR(20))
  RETURNS CHAR(50) DETERMINISTIC
  RETURN CONCAT('Hello, ',s,'!');

Query OK, 0 rows affected (0.00 sec)
```

```
mysql> SHOW FUNCTION CODE test.hello;
```

Pos	Instruction
0	freturn 254 concat('Hello, ',s@0,'!')

1 row in set (0.00 sec)

In this example, the nonexecutable `BEGIN` and `END` statements have disappeared, and for the `DECLARE variable_name` statement, only the executable part appears (the part where the default is assigned). For each statement that is taken from source, there is a code word `stmt` followed by a type (9 means `DROP`, 5 means `INSERT`, and so on). The final row contains an instruction `jump 2`, meaning `GOTO instruction #2`.

13.7.7.28 SHOW PROCEDURE STATUS Statement

```
SHOW PROCEDURE STATUS
  [LIKE 'pattern' | WHERE expr]
```

This statement is a MySQL extension. It returns characteristics of a stored procedure, such as the database, name, type, creator, creation and modification dates, and character set information. A similar statement, [SHOW FUNCTION STATUS](#), displays information about stored functions (see [Section 13.7.7.20, “SHOW FUNCTION STATUS Statement”](#)).

To use either statement, you must be the user named as the routine [DEFINER](#), have the [SHOW_ROUTINE](#) privilege, have the [SELECT](#) privilege at the global level, or have the [CREATE ROUTINE](#), [ALTER ROUTINE](#), or [EXECUTE](#) privilege granted at a scope that includes the routine.

The [LIKE](#) clause, if present, indicates which procedure or function names to match. The [WHERE](#) clause can be given to select rows using more general conditions, as discussed in [Section 26.8, “Extensions to SHOW Statements”](#).

```
mysql> SHOW PROCEDURE STATUS LIKE 'spl'\G
***** 1. row *****
      Db: test
      Name: spl
      Type: PROCEDURE
    Definer: testuser@localhost
  Modified: 2018-08-08 13:54:11
   Created: 2018-08-08 13:54:11
  Security_type: DEFINER
      Comment:
character_set_client: utf8mb4
collation_connection: utf8mb4_0900_ai_ci
Database Collation: utf8mb4_0900_ai_ci

mysql> SHOW FUNCTION STATUS LIKE 'hello'\G
***** 1. row *****
      Db: test
      Name: hello
      Type: FUNCTION
    Definer: testuser@localhost
  Modified: 2020-03-10 11:10:03
   Created: 2020-03-10 11:10:03
  Security_type: DEFINER
      Comment:
character_set_client: utf8mb4
collation_connection: utf8mb4_0900_ai_ci
Database Collation: utf8mb4_0900_ai_ci
```

[character_set_client](#) is the session value of the [character_set_client](#) system variable when the routine was created. [collation_connection](#) is the session value of the [collation_connection](#) system variable when the routine was created. [Database Collation](#) is the collation of the database with which the routine is associated.

Stored routine information is also available from the [INFORMATION_SCHEMA PARAMETERS](#) and [ROUTINES](#) tables. See [Section 26.3.20, “The INFORMATION_SCHEMA PARAMETERS Table”](#), and [Section 26.3.30, “The INFORMATION_SCHEMA ROUTINES Table”](#).

13.7.7.29 SHOW PROCESSLIST Statement

```
SHOW [ FULL ] PROCESSLIST
```

The MySQL process list indicates the operations currently being performed by the set of threads executing within the server. The [SHOW PROCESSLIST](#) statement is one source of process information. For a comparison of this statement with other sources, see [Sources of Process Information](#).



Note

As of MySQL 8.0.22, an alternative implementation for [SHOW PROCESSLIST](#) is available based on the Performance Schema [processlist](#) table, which, unlike the default [SHOW PROCESSLIST](#) implementation, does not require a mutex and has better performance characteristics. For details, see [Section 27.12.21.6, “The processlist Table”](#).

If you have the `PROCESS` privilege, you can see all threads, even those belonging to other users. Otherwise (without the `PROCESS` privilege), nonanonymous users have access to information about their own threads but not threads for other users, and anonymous users have no access to thread information.

Without the `FULL` keyword, `SHOW PROCESSLIST` displays only the first 100 characters of each statement in the `Info` field.

The `SHOW PROCESSLIST` statement is very useful if you get the “too many connections” error message and want to find out what is going on. MySQL reserves one extra connection to be used by accounts that have the `CONNECTION_ADMIN` privilege (or the deprecated `SUPER` privilege), to ensure that administrators should always be able to connect and check the system (assuming that you are not giving this privilege to all your users).

Threads can be killed with the `KILL` statement. See Section 13.7.8.4, “`KILL Statement`”.

Example of `SHOW PROCESSLIST` output:

```
mysql> SHOW FULL PROCESSLIST\G
***** 1. row *****
    Id: 1
    User: system user
    Host:
      db: NULL
Command: Connect
    Time: 1030455
  State: Waiting for master to send event
  Info: NULL
***** 2. row *****
    Id: 2
    User: system user
    Host:
      db: NULL
Command: Connect
    Time: 1004
  State: Has read all relay log; waiting for the slave
         I/O thread to update it
  Info: NULL
***** 3. row *****
    Id: 3112
    User: replikator
    Host: artemis:2204
      db: NULL
Command: Binlog Dump
    Time: 2144
  State: Has sent all binlog to slave; waiting for binlog to be updated
  Info: NULL
***** 4. row *****
    Id: 3113
    User: replikator
    Host: iconnect2:45781
      db: NULL
Command: Binlog Dump
    Time: 2086
  State: Has sent all binlog to slave; waiting for binlog to be updated
  Info: NULL
***** 5. row *****
    Id: 3123
    User: stefan
    Host: localhost
      db: apollon
Command: Query
    Time: 0
  State: NULL
  Info: SHOW FULL PROCESSLIST
```

`SHOW PROCESSLIST` output has these columns:

- `Id`

The connection identifier. This is the same value displayed in the `ID` column of the `INFORMATION_SCHEMA PROCESSLIST` table, displayed in the `PROCESSLIST_ID` column of the Performance Schema `threads` table, and returned by the `CONNECTION_ID()` function within the thread.

- [User](#)

The MySQL user who issued the statement. A value of `system user` refers to a nonclient thread spawned by the server to handle tasks internally, for example, a delayed-row handler thread or an I/O (receiver) or SQL (applier) thread used on replica hosts. For `system user`, there is no host specified in the `Host` column. `unauthenticated user` refers to a thread that has become associated with a client connection but for which authentication of the client user has not yet occurred. `event_scheduler` refers to the thread that monitors scheduled events (see [Section 25.4, “Using the Event Scheduler”](#)).



Note

A `User` value of `system user` is distinct from the `SYSTEM_USER` privilege. The former designates internal threads. The latter distinguishes the system user and regular user account categories (see [Section 6.2.11, “Account Categories”](#)).

- [Host](#)

The host name of the client issuing the statement (except for `system user`, for which there is no host). The host name for TCP/IP connections is reported in `host_name:client_port` format to make it easier to determine which client is doing what.

- [db](#)

The default database for the thread, or `NULL` if none has been selected.

- [Command](#)

The type of command the thread is executing on behalf of the client, or `Sleep` if the session is idle. For descriptions of thread commands, see [Section 8.14, “Examining Server Thread \(Process\) Information”](#). The value of this column corresponds to the `COM_xxx` commands of the client/server protocol and `Com_xxx` status variables. See [Section 5.1.10, “Server Status Variables”](#).

- [Time](#)

The time in seconds that the thread has been in its current state. For a replica SQL thread, the value is the number of seconds between the timestamp of the last replicated event and the real time of the replica host. See [Section 17.2.3, “Replication Threads”](#).

- [State](#)

An action, event, or state that indicates what the thread is doing. For descriptions of `State` values, see [Section 8.14, “Examining Server Thread \(Process\) Information”](#).

Most states correspond to very quick operations. If a thread stays in a given state for many seconds, there might be a problem that needs to be investigated.

- [Info](#)

The statement the thread is executing, or `NULL` if it is executing no statement. The statement might be the one sent to the server, or an innermost statement if the statement executes other statements. For example, if a `CALL` statement executes a stored procedure that is executing a `SELECT` statement, the `Info` value shows the `SELECT` statement.

13.7.7.30 SHOW PROFILE Statement

```

SHOW PROFILE [type [, type] ... ]
  [FOR QUERY n]
  [LIMIT row_count [OFFSET offset]]

type: {
  ALL
  BLOCK IO
  CONTEXT SWITCHES
  CPU
  IPC
  MEMORY
  PAGE FAULTS
  SOURCE
  SWAPS
}

```

The `SHOW PROFILE` and `SHOW PROFILES` statements display profiling information that indicates resource usage for statements executed during the course of the current session.



Note

The `SHOW PROFILE` and `SHOW PROFILES` statements are deprecated; expect them to be removed in a future MySQL release. Use the [Performance Schema](#) instead; see [Section 27.19.1, “Query Profiling Using Performance Schema”](#).

To control profiling, use the `profiling` session variable, which has a default value of 0 (`OFF`). Enable profiling by setting `profiling` to 1 or `ON`:

```
mysql> SET profiling = 1;
```

`SHOW PROFILES` displays a list of the most recent statements sent to the server. The size of the list is controlled by the `profiling_history_size` session variable, which has a default value of 15. The maximum value is 100. Setting the value to 0 has the practical effect of disabling profiling.

All statements are profiled except `SHOW PROFILE` and `SHOW PROFILES`, so neither of those statements appears in the profile list. Malformed statements are profiled. For example, `SHOW PROFILING` is an illegal statement, and a syntax error occurs if you try to execute it, but it shows up in the profiling list.

`SHOW PROFILE` displays detailed information about a single statement. Without the `FOR QUERY n` clause, the output pertains to the most recently executed statement. If `FOR QUERY n` is included, `SHOW PROFILE` displays information for statement `n`. The values of `n` correspond to the `Query_ID` values displayed by `SHOW PROFILES`.

The `LIMIT row_count` clause may be given to limit the output to `row_count` rows. If `LIMIT` is given, `OFFSET offset` may be added to begin the output `offset` rows into the full set of rows.

By default, `SHOW PROFILE` displays `Status` and `Duration` columns. The `Status` values are like the `State` values displayed by `SHOW PROCESSLIST`, although there might be some minor differences in interpretation for the two statements for some status values (see [Section 8.14, “Examining Server Thread \(Process\) Information”](#)).

Optional `type` values may be specified to display specific additional types of information:

- `ALL` displays all information
- `BLOCK IO` displays counts for block input and output operations
- `CONTEXT SWITCHES` displays counts for voluntary and involuntary context switches
- `CPU` displays user and system CPU usage times
- `IPC` displays counts for messages sent and received
- `MEMORY` is not currently implemented

SHOW Statements

- **PAGE FAULTS** displays counts for major and minor page faults
- **SOURCE** displays the names of functions from the source code, together with the name and line number of the file in which the function occurs
- **SWAPS** displays swap counts

Profiling is enabled per session. When a session ends, its profiling information is lost.

```
mysql> SELECT @@profiling;
+-----+
| @@profiling |
+-----+
|          0 |
+-----+
1 row in set (0.00 sec)

mysql> SET profiling = 1;
Query OK, 0 rows affected (0.00 sec)

mysql> DROP TABLE IF EXISTS t1;
Query OK, 0 rows affected, 1 warning (0.00 sec)

mysql> CREATE TABLE t1 (id INT);
Query OK, 0 rows affected (0.01 sec)

mysql> SHOW PROFILES;
+-----+-----+-----+
| Query_ID | Duration | Query           |
+-----+-----+-----+
|      0 | 0.000088 | SET PROFILING = 1
|      1 | 0.000136 | DROP TABLE IF EXISTS t1
|      2 | 0.011947 | CREATE TABLE t1 (id INT)
+-----+-----+-----+
3 rows in set (0.00 sec)

mysql> SHOW PROFILE;
+-----+-----+
| Status          | Duration |
+-----+-----+
| checking permissions | 0.000040 |
| creating table   | 0.000056 |
| After create     | 0.011363 |
| query end        | 0.000375 |
| freeing items    | 0.000089 |
| logging slow query | 0.000019 |
| cleaning up       | 0.000005 |
+-----+-----+
7 rows in set (0.00 sec)

mysql> SHOW PROFILE FOR QUERY 1;
+-----+-----+
| Status          | Duration |
+-----+-----+
| query end        | 0.000107 |
| freeing items    | 0.000008 |
| logging slow query | 0.000015 |
| cleaning up       | 0.000006 |
+-----+-----+
4 rows in set (0.00 sec)

mysql> SHOW PROFILE CPU FOR QUERY 2;
+-----+-----+-----+-----+
| Status          | Duration | CPU_user | CPU_system |
+-----+-----+-----+-----+
| checking permissions | 0.000040 | 0.000038 | 0.000002 |
| creating table   | 0.000056 | 0.000028 | 0.000028 |
| After create     | 0.011363 | 0.000217 | 0.001571 |
| query end        | 0.000375 | 0.000013 | 0.000028 |
| freeing items    | 0.000089 | 0.000010 | 0.000014 |
| logging slow query | 0.000019 | 0.000009 | 0.000010 |
| cleaning up       | 0.000005 | 0.000003 | 0.000002 |
+-----+-----+-----+-----+
```

```
+-----+-----+-----+
7 rows in set (0.00 sec)
```

**Note**

Profiling is only partially functional on some architectures. For values that depend on the `getrusage()` system call, `NULL` is returned on systems such as Windows that do not support the call. In addition, profiling is per process and not per thread. This means that activity on threads within the server other than your own may affect the timing information that you see.

Profiling information is also available from the `INFORMATION_SCHEMA PROFILING` table. See [Section 26.3.24, “The INFORMATION_SCHEMA PROFILING Table”](#). For example, the following queries are equivalent:

```
SHOW PROFILE FOR QUERY 2;

SELECT STATE, FORMAT(DURATION, 6) AS DURATION
FROM INFORMATION_SCHEMA.PROFILING
WHERE QUERY_ID = 2 ORDER BY SEQ;
```

13.7.7.31 SHOW PROFILES Statement

```
SHOW PROFILES
```

The `SHOW PROFILES` statement, together with `SHOW PROFILE`, displays profiling information that indicates resource usage for statements executed during the course of the current session. For more information, see [Section 13.7.7.30, “SHOW PROFILE Statement”](#).

**Note**

The `SHOW PROFILE` and `SHOW PROFILES` statements are deprecated; expect it to be removed in a future MySQL release. Use the Performance Schema instead; see [Section 27.19.1, “Query Profiling Using Performance Schema”](#).

13.7.7.32 SHOW RELAYLOG EVENTS Statement

```
SHOW RELAYLOG EVENTS
[ IN 'log_name' ]
[ FROM pos ]
[ LIMIT [offset,] row_count ]
[ channel_option ]

channel_option:
  FOR CHANNEL channel
```

Shows the events in the relay log of a replica. If you do not specify '`log_name`', the first relay log is displayed. This statement has no effect on the source. `SHOW RELAYLOG EVENTS` requires the `REPLICATION SLAVE` privilege.

The `LIMIT` clause has the same syntax as for the `SELECT` statement. See [Section 13.2.13, “SELECT Statement”](#).

**Note**

Issuing a `SHOW RELAYLOG EVENTS` with no `LIMIT` clause could start a very time- and resource-consuming process because the server returns to the client the complete contents of the relay log (including all statements modifying data that have been received by the replica).

The optional `FOR CHANNEL channel` clause enables you to name which replication channel the statement applies to. Providing a `FOR CHANNEL channel` clause applies the statement to a specific replication channel. If no channel is named and no extra channels exist, the statement applies to the default channel.

When using multiple replication channels, if a `SHOW RELAYLOG EVENTS` statement does not have a channel defined using a `FOR CHANNEL channel` clause an error is generated. See [Section 17.2.2, “Replication Channels”](#) for more information.

`SHOW RELAYLOG EVENTS` displays the following fields for each event in the relay log:

- `Log_name`

The name of the file that is being listed.

- `Pos`

The position at which the event occurs.

- `Event_type`

An identifier that describes the event type.

- `Server_id`

The server ID of the server on which the event originated.

- `End_log_pos`

The value of `End_log_pos` for this event in the source's binary log.

- `Info`

More detailed information about the event type. The format of this information depends on the event type.

For compressed transaction payloads, the `Transaction_payload_event` is first printed as a single unit, then it is unpacked and each event inside it is printed.

Some events relating to the setting of user and system variables are not included in the output from `SHOW RELAYLOG EVENTS`. To get complete coverage of events within a relay log, use `mysqlbinlog`.

13.7.7.33 SHOW REPLICAS Statement

```
{SHOW REPLICAS}
```

Displays a list of replicas currently registered with the source. From MySQL 8.0.22, use `SHOW REPLICAS` in place of `SHOW SLAVE HOSTS`, which is deprecated from that release. In releases before MySQL 8.0.22, use `SHOW SLAVE HOSTS`. `SHOW REPLICAS` requires the `REPLICATION SLAVE` privilege.

`SHOW REPLICAS` should be executed on a server that acts as a replication source. The statement displays information about servers that are or have been connected as replicas, with each row of the result corresponding to one replica server, as shown here:

```
mysql> SHOW REPLICAS;
+-----+-----+-----+-----+
| Server_id | Host      | Port | Source_id | Replica_UUID           |
+-----+-----+-----+-----+
|       10 | iconnect2 | 3306 |         3 | 14cb6624-7f93-11e0-b2c0-c80aa9429562 |
|       21 | athena    | 3306 |         3 | 07af4990-f41f-11df-a566-7ac56fdaf645 |
+-----+-----+-----+-----+
```

- `Server_id`: The unique server ID of the replica server, as configured in the replica server's option file, or on the command line with `--server-id=value`.
- `Host`: The host name of the replica server, as specified on the replica with the `--report-host` option. This can differ from the machine name as configured in the operating system.

- **User**: The replica server user name, as specified on the replica with the `--report-user` option. Statement output includes this column only if the source server is started with the `--show-replica-auth-info` or `--show-slave-auth-info` option.
- **Password**: The replica server password, as specified on the replica with the `--report-password` option. Statement output includes this column only if the source server is started with the `--show-replica-auth-info` or `--show-slave-auth-info` option.
- **Port**: The port on the source to which the replica server is listening, as specified on the replica with the `--report-port` option.

A zero in this column means that the replica port (`--report-port`) was not set.

- **Source_id**: The unique server ID of the source server that the replica server is replicating from. This is the server ID of the server on which `SHOW REPLICAS` is executed, so this same value is listed for each row in the result.
- **Replica_UUID**: The globally unique ID of this replica, as generated on the replica and found in the replica's `auto.cnf` file.

13.7.7.34 SHOW SLAVE HOSTS | SHOW REPLICAS Statement

```
{SHOW SLAVE HOSTS | SHOW REPLICAS}
```

Displays a list of replicas currently registered with the source. From MySQL 8.0.22, `SHOW SLAVE HOSTS` is deprecated and the alias `SHOW REPLICAS` should be used instead. The statement works in the same way as before, only the terminology used for the statement and its output has changed. Both versions of the statement update the same status variables when used. Please see the documentation for `SHOW REPLICAS` for a description of the statement.

13.7.7.35 SHOW REPLICAS STATUS Statement

```
SHOW {REPLICAS | SLAVE} STATUS [FOR CHANNEL channel]
```

This statement provides status information on essential parameters of the replica threads. From MySQL 8.0.22, use `SHOW REPLICAS STATUS` in place of `SHOW SLAVE STATUS`, which is deprecated from that release. In releases before MySQL 8.0.22, use `SHOW SLAVE STATUS`. The statement requires the `REPLICATION CLIENT` privilege (or the deprecated `SUPER` privilege).

`SHOW REPLICAS STATUS` is nonblocking. When run concurrently with `STOP REPLICAS`, `SHOW REPLICAS STATUS` returns without waiting for `STOP REPLICAS` to finish shutting down the replication SQL (applier) thread or replication I/O (receiver) thread (or both). This permits use in monitoring and other applications where getting an immediate response from `SHOW REPLICAS STATUS` is more important than ensuring that it returned the latest data. The `SLAVE` keyword was replaced with `REPLICAS` in MySQL 8.0.22.

If you issue this statement using the `mysql` client, you can use a `\G` statement terminator rather than a semicolon to obtain a more readable vertical layout:

```
mysql> SHOW REPLICAS STATUS\G
***** 1. row *****
    Replica_IO_State: Waiting for source to send event
          Source_Host: localhost
          Source_User: repl
          Source_Port: 13000
      Connect_Retry: 60
        Source_Log_File: source-bin.000002
     Read_Source_Log_Pos: 1307
        Relay_Log_File: replica-relay-bin.000003
     Relay_Log_Pos: 1508
Relay_Source_Log_File: source-bin.000002
   Replica_IO_Running: Yes
  Replica_SQL_Running: Yes
```

```

        Replicate_Do_DB:
        Replicate_Ignore_DB:
        Replicate_Do_Table:
        Replicate_Ignore_Table:
        Replicate_Wild_Do_Table:
        Replicate_Wild_Ignore_Table:
        Last_Error:
        Skip_Counter: 0
        Exec_Source_Log_Pos: 1307
        Relay_Log_Space: 1858
        Until_Condition: None
        Until_Log_File:
        Until_Log_Pos: 0
        Source_SSL_Allowed: No
        Source_SSL_CA_File:
        Source_SSL_CA_Path:
        Source_SSL_Cert:
        Source_SSL_Cipher:
        Source_SSL_Key:
        Seconds_Behind_Source: 0
Source_SSL_Verify_Server_Cert: No
        Last_IO_Error:
        Last_SQL_Error:
Replicate_Ignore_Server_Ids:
        Source_Server_Id: 1
        Source_UUID: 3e11fa47-71ca-11e1-9e33-c80aa9429562
        Source_Info_File:
        SQL_Delay: 0
        SQL_Remaining_Delay: NULL
Replica_SQL_Running_State: Reading event from the relay log
        Source_Retry_Count: 10
        Source_Bind:
        Last_IO_Error_Timestamp:
Last_SQL_Error_Timestamp:
        Source_SSL_Crl:
        Source_SSL_Crlpath:
        Retrieved_Gtid_Set: 3e11fa47-71ca-11e1-9e33-c80aa9429562:1-5
        Executed_Gtid_Set: 3e11fa47-71ca-11e1-9e33-c80aa9429562:1-5
        Auto_Position: 1
Replicate_Rewrite_DB:
        Channel_name:
        Source_TLS_Version: TLSv1.2
Source_public_key_path: public_key.pem
        Get_source_public_key: 0
        Network_Namespace:

```

The Performance Schema provides tables that expose replication information. This is similar to the information available from the `SHOW REPLICAS STATUS` statement, but represented in table form. For details, see [Section 27.12.11, “Performance Schema Replication Tables”](#).

From MySQL 8.0.27, you can set the `GTID_ONLY` option on the `CHANGE REPLICATION SOURCE TO` statement to stop a replication channel from persisting file names and file positions in the replication metadata repositories. With this setting, file positions for the source binary log file and the relay log file are tracked in memory. The `SHOW REPLICAS STATUS` statement still displays file positions in normal use. However, because the file positions are not being regularly updated in the connection metadata repository and the applier metadata repository except in a few situations, they are likely to be out of date if the server is restarted.

For a replication channel with the `GTID_ONLY` setting after a server start, the read and applied file positions for the source binary log file (`Read_Source_Log_Pos` and `Exec_Source_Log_Pos`) are set to zero, and the file names (`Source_Log_File` and `Relay_Source_Log_File`) are set to `INVALID`. The relay log file name (`Relay_Log_File`) is set according to the `relay_log_recovery` setting, either a new file that was created at server start or the first relay log file present. The file position (`Relay_Log_Pos`) is set to position 4, and GTID auto-skip is used to skip any transactions in the file that were already applied.

When the receiver thread contacts the source and gets valid position information, the read position (`Read_Source_Log_Pos`) and file name (`Source_Log_File`) are updated with the correct data and become valid. When the applier thread applies a transaction from the source, or skips an already executed transaction, the executed position (`Exec_Source_Log_Pos`) and file name (`Relay_Source_Log_File`) are updated with the correct data and become valid. The relay log file position (`Relay_Log_Pos`) is also updated at that time.

The following list describes the fields returned by `SHOW REPLICAS STATUS`. For additional information about interpreting their meanings, see [Section 17.1.7.1, “Checking Replication Status”](#).

- `Replica_IO_State`

A copy of the `State` field of the `SHOW PROCESSLIST` output for the replica I/O (receiver) thread. This tells you what the thread is doing: trying to connect to the source, waiting for events from the source, reconnecting to the source, and so on. For a listing of possible states, see [Section 8.14.5, “Replication I/O \(Receiver\) Thread States”](#).

- `Source_Host`

The source host that the replica is connected to.

- `Source_User`

The user name of the account used to connect to the source.

- `Source_Port`

The port used to connect to the source.

- `Connect_Retry`

The number of seconds between connect retries (default 60). This can be set with a `CHANGE REPLICATION SOURCE TO` statement (from MySQL 8.0.23) or `CHANGE MASTER TO` statement (before MySQL 8.0.23).

- `Source_Log_File`

The name of the source binary log file from which the I/O (receiver) thread is currently reading. This is set to `INVALID` for a replication channel with the `GTID_ONLY` setting after a server start. It will be updated when the replica contacts the source.

- `Read_Source_Log_Pos`

The position in the current source binary log file up to which the I/O (receiver) thread has read. This is set to zero for a replication channel with the `GTID_ONLY` setting after a server start. It will be updated when the replica contacts the source.

- `Relay_Log_File`

The name of the relay log file from which the SQL (applier) thread is currently reading and executing.

- `Relay_Log_Pos`

The position in the current relay log file up to which the SQL (applier) thread has read and executed.

- `Relay_Source_Log_File`

The name of the source binary log file containing the most recent event executed by the SQL (applier) thread. This is set to `INVALID` for a replication channel with the `GTID_ONLY` setting after a server start. It will be updated when a transaction is executed or skipped.

- `Replica_IO_Running`

Whether the replication I/O (receiver) thread is started and has connected successfully to the source. Internally, the state of this thread is represented by one of the following three values:

- **MYSQL_REPLICA_NOT_RUN.** The replication I/O (receiver) thread is not running. For this state, `Replica_IO_Running` is `No`.
- **MYSQL_REPLICA_RUN_NOT_CONNECT.** The replication I/O (receiver) thread is running, but is not connected to a replication source. For this state, `Replica_IO_Running` is `Connecting`.
- **MYSQL_REPLICA_RUN_CONNECT.** The replication I/O (receiver) thread is running, and is connected to a replication source. For this state, `Replica_IO_Running` is `Yes`.
- `Replica_SQL_Running`

Whether the replication SQL (applier) thread is started.

- `Replicate_Do_DB`, `Replicate_Ignore_DB`

The names of any databases that were specified with the `--replicate-do-db` and `--replicate-ignore-db` options, or the `CHANGE REPLICATION FILTER` statement. If the `FOR CHANNEL` clause was used, the channel specific replication filters are shown. Otherwise, the replication filters for every replication channel are shown.

- `Replicate_Do_Table`, `Replicate_Ignore_Table`, `Replicate_Wild_Do_Table`, `Replicate_Wild_Ignore_Table`

The names of any tables that were specified with the `--replicate-do-table`, `--replicate-ignore-table`, `--replicate-wild-do-table`, and `--replicate-wild-ignore-table` options, or the `CHANGE REPLICATION FILTER` statement. If the `FOR CHANNEL` clause was used, the channel specific replication filters are shown. Otherwise, the replication filters for every replication channel are shown.

- `Last_Errorno`, `Last_Error`

These columns are aliases for `Last_SQL_Errorno` and `Last_SQL_Error`.

Issuing `RESET MASTER` or `RESET REPLICA` resets the values shown in these columns.



Note

When the replication SQL thread receives an error, it reports the error first, then stops the SQL thread. This means that there is a small window of time during which `SHOW REPLICAS STATUS` shows a nonzero value for `Last_SQL_Errorno` even though `Replica_SQL_Running` still displays `Yes`.

- `Skip_Counter`

The current value of the `sql_slave_skip_counter` system variable. See [SET GLOBAL sql_slave_skip_counter Syntax](#).

- `Exec_Source_Log_Pos`

The position in the current source binary log file to which the replication SQL thread has read and executed, marking the start of the next transaction or event to be processed. This is set to zero for a replication channel with the `GTID_ONLY` setting after a server start. It will be updated when a transaction is executed or skipped.

You can use this value with the `CHANGE REPLICATION SOURCE TO` statement's `SOURCE_LOG_POS` option (from MySQL 8.0.23) or the `CHANGE MASTER TO` statement's `MASTER_LOG_POS` option (before MySQL 8.0.23) when starting a new replica from an existing replica, so that the new replica reads from this point. The coordinates given by

(`Relay_Source_Log_File`, `Exec_Source_Log_Pos`) in the source's binary log correspond to the coordinates given by (`Relay_Log_File`, `Relay_Log_Pos`) in the relay log.

Inconsistencies in the sequence of transactions from the relay log which have been executed can cause this value to be a "low-water mark". In other words, transactions appearing before the position are guaranteed to have committed, but transactions after the position may have committed or not. If these gaps need to be corrected, use `START REPLICA UNTIL SQL_AFTER_MTS_GAPS`. See [Section 17.5.1.34, "Replication and Transaction Inconsistencies"](#) for more information.

- `Relay_Log_Space`

The total combined size of all existing relay log files.

- `Until_Condition`, `Until_Log_File`, `Until_Log_Pos`

The values specified in the `UNTIL` clause of the `START REPLICA` statement.

`Until_Condition` has these values:

- `None` if no `UNTIL` clause was specified.
- `Source` if the replica is reading until a given position in the source's binary log.
- `Relay` if the replica is reading until a given position in its relay log.
- `SQL_BEFORE_GTIDS` if the replication SQL thread is processing transactions until it has reached the first transaction whose GTID is listed in the `gtid_set`.
- `SQL_AFTER_GTIDS` if the replication threads are processing all transactions until the last transaction in the `gtid_set` has been processed by both threads.
- `SQL_AFTER_MTS_GAPS` if a multithreaded replica's SQL threads are running until no more gaps are found in the relay log.

`Until_Log_File` and `Until_Log_Pos` indicate the log file name and position that define the coordinates at which the replication SQL thread stops executing.

For more information on `UNTIL` clauses, see [Section 13.4.2.9, "START SLAVE Statement"](#).

- `Source_SSL_Allowed`, `Source_SSL_CA_File`, `Source_SSL_CA_Path`, `Source_SSL_Cert`, `Source_SSL_Cipher`, `Source_SSL_CRL_File`, `Source_SSL_CRL_Path`, `Source_SSL_Key`, `Source_SSL_Verify_Server_Cert`

These fields show the SSL parameters used by the replica to connect to the source, if any.

`Source_SSL_Allowed` has these values:

- `Yes` if an SSL connection to the source is permitted.
- `No` if an SSL connection to the source is not permitted.
- `Ignored` if an SSL connection is permitted but the replica server does not have SSL support enabled.

The values of the other SSL-related fields correspond to the values of the `SOURCE_SSL_*` options of the `CHANGE REPLICATION SOURCE TO` statement (from MySQL 8.0.23), or the `MASTER_SSL_*` options of the `CHANGE MASTER TO` statement (before MySQL 8.0.23). See [Section 13.4.2.1, "CHANGE MASTER TO Statement"](#).

- `Seconds_Behind_Source`

This field is an indication of how "late" the replica is:

- When the replica is actively processing updates, this field shows the difference between the current timestamp on the replica and the original timestamp logged on the source for the event currently being processed on the replica.
- When no event is currently being processed on the replica, this value is 0.

In essence, this field measures the time difference in seconds between the replication SQL (applier) thread and the replication I/O (receiver) thread. If the network connection between source and replica is fast, the replication receiver thread is very close to the source, so this field is a good approximation of how late the replication applier thread is compared to the source. If the network is slow, this is *not* a good approximation; the replication applier thread may quite often be caught up with the slow-reading replication receiver thread, so `Seconds_Behind_Source` often shows a value of 0, even if the replication receiver thread is late compared to the source. In other words, *this column is useful only for fast networks*.

This time difference computation works even if the source and replica do not have identical clock times, provided that the difference, computed when the replica receiver thread starts, remains constant from then on. Any changes, including NTP updates, can lead to clock skews that can make calculation of `Seconds_Behind_Source` less reliable.

In MySQL 8.0, this field is `NULL` (undefined or unknown) if the replication applier thread is not running, or if the applier thread has consumed all of the relay log and the replication receiver thread is not running. (In older versions of MySQL, this field was `NULL` if the replication applier thread or the replication receiver thread was not running or was not connected to the source.) If the replication receiver thread is running but the relay log is exhausted, `Seconds_Behind_Source` is set to 0.

The value of `Seconds_Behind_Source` is based on the timestamps stored in events, which are preserved through replication. This means that if a source M1 is itself a replica of M0, any event from M1's binary log that originates from M0's binary log has M0's timestamp for that event. This enables MySQL to replicate `TIMESTAMP` successfully. However, the problem for `Seconds_Behind_Source` is that if M1 also receives direct updates from clients, the `Seconds_Behind_Source` value randomly fluctuates because sometimes the last event from M1 originates from M0 and sometimes is the result of a direct update on M1.

When using a multithreaded replica, you should keep in mind that this value is based on `Exec_Source_Log_Pos`, and so may not reflect the position of the most recently committed transaction.

- `Last_IO_Errno`, `Last_IO_Error`

The error number and error message of the most recent error that caused the replication I/O (receiver) thread to stop. An error number of 0 and message of the empty string mean “no error.” If the `Last_IO_Error` value is not empty, the error values also appear in the replica's error log.

I/O error information includes a timestamp showing when the most recent I/O (receiver)thread error occurred. This timestamp uses the format `YYMMDD hh:mm:ss`, and appears in the `Last_IO_Error_Timestamp` column.

Issuing `RESET MASTER` or `RESET REPLICA` resets the values shown in these columns.

- `Last_SQL_Errno`, `Last_SQL_Error`

The error number and error message of the most recent error that caused the replication SQL (applier) thread to stop. An error number of 0 and message of the empty string mean “no error.” If the `Last_SQL_Error` value is not empty, the error values also appear in the replica's error log.

If the replica is multithreaded, the replication SQL thread is the coordinator for worker threads. In this case, the `Last_SQL_Error` field shows exactly what the `Last_Error_Message` column in the Performance Schema `replication_applier_status_by_coordinator` table shows.

The field value is modified to suggest that there may be more failures in the other worker threads which can be seen in the `replication_applier_status_by_worker` table that shows each worker thread's status. If that table is not available, the replica error log can be used. The log or the `replication_applier_status_by_worker` table should also be used to learn more about the failure shown by `SHOW REPLICA STATUS` or the coordinator table.

SQL error information includes a timestamp showing when the most recent SQL (applier) thread error occurred. This timestamp uses the format `YYMMDD hh:mm:ss`, and appears in the `Last_SQL_Error_Timestamp` column.

Issuing `RESET MASTER` or `RESET REPLICA` resets the values shown in these columns.

In MySQL 8.0, all error codes and messages displayed in the `Last_SQL_Errno` and `Last_SQL_Error` columns correspond to error values listed in [Server Error Message Reference](#). This was not always true in previous versions. (Bug #11760365, Bug #52768)

- [Replicate_Ignore_Server_Ids](#)

Any server IDs that have been specified using the `IGNORE_SERVER_IDS` option of the `CHANGE REPLICATION SOURCE TO | CHANGE MASTER TO` statement, so that the replica ignores events from these servers. This option is used in a circular or other multi-source replication setup when one of the servers is removed. If any server IDs have been set in this way, a comma-delimited list of one or more numbers is shown. If no server IDs have been set, the field is blank.



Note

The `Ignored_server_ids` value in the `slave_master_info` table also shows the server IDs to be ignored, but as a space-delimited list, preceded by the total number of server IDs to be ignored. For example, if a `CHANGE REPLICATION SOURCE TO | CHANGE MASTER TO` statement containing the `IGNORE_SERVER_IDS = (2,6,9)` option has been issued to tell a replica to ignore sources having the server ID 2, 6, or 9, that information appears as shown here:

```
Replicate_Ignore_Server_Ids: 2, 6, 9
```

```
Ignored_server_ids: 3, 2, 6, 9
```

`Replicate_Ignore_Server_Ids` filtering is performed by the I/O (receiver) thread, rather than by the SQL (applier) thread, which means that events which are filtered out are not written to the relay log. This differs from the filtering actions taken by server options such `--replicate-do-table`, which apply to the applier thread.



Note

From MySQL 8.0, a deprecation warning is issued if `SET GTID_MODE=ON` is issued when any channel has existing server IDs set with `IGNORE_SERVER_IDS`. Before starting GTID-based replication, use `SHOW REPLICA STATUS` to check for and clear all ignored server ID lists on the servers involved. You can clear a list by issuing a `CHANGE REPLICATION SOURCE TO | CHANGE MASTER TO` statement containing the `IGNORE_SERVER_IDS` option with an empty list.

- [Source_Server_Id](#)

The `server_id` value from the source.

- [Source_UUID](#)

The `server_uuid` value from the source.

- [Source_Info_File](#)

The location of the `master.info` file, the use of which is now deprecated. By default from MySQL 8.0, a table is used instead for the replica's connection metadata repository.

- [SQL_Delay](#)

The number of seconds that the replica must lag the source.

- [SQL_Remaining_Delay](#)

When `Replica_SQL_Running_State` is `Waiting until MASTER_DELAY seconds after source executed event`, this field contains the number of delay seconds remaining. At other times, this field is `NULL`.

- [Replica_SQL_Running_State](#)

The state of the SQL thread (analogous to `Replica_IO_State`). The value is identical to the `State` value of the SQL thread as displayed by `SHOW PROCESSLIST`. [Section 8.14.6, “Replication SQL Thread States”](#), provides a listing of possible states.

- [Source_Retry_Count](#)

The number of times the replica can attempt to reconnect to the source in the event of a lost connection. This value can be set using the `SOURCE_RETRY_COUNT | MASTER_RETRY_COUNT` option of the `CHANGE REPLICATION SOURCE TO` statement (from MySQL 8.0.23) or `CHANGE MASTER TO` statement (before MySQL 8.0.23), or the older `--master-retry-count` server option (still supported for backward compatibility).

- [Source_Bind](#)

The network interface that the replica is bound to, if any. This is set using the `SOURCE_BIND | MASTER_BIND` option for the `CHANGE REPLICATION SOURCE TO` statement (from MySQL 8.0.23) or `CHANGE MASTER TO` statement (before MySQL 8.0.23).

- [Last_IO_Error_Timestamp](#)

A timestamp in `YYMMDD hh:mm:ss` format that shows when the most recent I/O error took place.

- [Last_SQL_Error_Timestamp](#)

A timestamp in `YYMMDD hh:mm:ss` format that shows when the most recent SQL error occurred.

- [Retrieved_Gtid_Set](#)

The set of global transaction IDs corresponding to all transactions received by this replica. Empty if GTIDs are not in use. See [GTID Sets](#) for more information.

This is the set of all GTIDs that exist or have existed in the relay logs. Each GTID is added as soon as the `Gtid_log_event` is received. This can cause partially transmitted transactions to have their GTIDs included in the set.

When all relay logs are lost due to executing `RESET REPLICA` or `CHANGE REPLICATION SOURCE TO | CHANGE MASTER TO`, or due to the effects of the `--relay-log-recovery` option, the set is cleared. When `relay_log_purge = 1`, the newest relay log is always kept, and the set is not cleared.

- [Executed_Gtid_Set](#)

The set of global transaction IDs written in the binary log. This is the same as the value for the global `gtid_executed` system variable on this server, as well as the value for `Executed_Gtid_Set` in the output of `SHOW MASTER STATUS` on this server. Empty if GTIDs are not in use. See [GTID Sets](#) for more information.

- `Auto_Position`

1 if GTID auto-positioning is in use for the channel, otherwise 0.

- `Replicate_Rewrite_DB`

The `Replicate_Rewrite_DB` value displays any replication filtering rules that were specified. For example, if the following replication filter rule was set:

```
CHANGE REPLICATION FILTER REPLICATE_REWRITE_DB=( (db1,db2), (db3,db4));
```

the `Replicate_Rewrite_DB` value displays:

```
Replicate_Rewrite_DB: (db1,db2),(db3,db4)
```

For more information, see [Section 13.4.2.2, “CHANGE REPLICATION FILTER Statement”](#).

- `Channel_name`

The replication channel which is being displayed. There is always a default replication channel, and more replication channels can be added. See [Section 17.2.2, “Replication Channels”](#) for more information.

- `Master_TLS_Version`

The TLS version used on the source. For TLS version information, see [Section 6.3.2, “Encrypted Connection TLS Protocols and Ciphers”](#).

- `Source_public_key_path`

The path name to a file containing a replica-side copy of the public key required by the source for RSA key pair-based password exchange. The file must be in PEM format. This column applies to replicas that authenticate with the `sha256_password` or `caching_sha2_password` authentication plugin.

If `Source_public_key_path` is given and specifies a valid public key file, it takes precedence over `Get_source_public_key`.

- `Get_source_public_key`

Whether to request from the source the public key required for RSA key pair-based password exchange. This column applies to replicas that authenticate with the `caching_sha2_password` authentication plugin. For that plugin, the source does not send the public key unless requested.

If `Source_public_key_path` is given and specifies a valid public key file, it takes precedence over `Get_source_public_key`.

- `Network_Namespace`

The network namespace name; empty if the connection uses the default (global) namespace. For information about network namespaces, see [Section 5.1.14, “Network Namespace Support”](#). This column was added in MySQL 8.0.22.

13.7.7.36 SHOW SLAVE | REPLICA STATUS Statement

```
SHOW {SLAVE | REPLICA} STATUS [FOR CHANNEL channel]
```

This statement provides status information on essential parameters of the replica threads. From MySQL 8.0.22, `SHOW SLAVE STATUS` is deprecated and the alias `SHOW REPLICA STATUS` should be used instead. The statement works in the same way as before, only the terminology used for the statement and its output has changed. Both versions of the statement update the same status variables when used. Please see the documentation for `SHOW REPLICA STATUS` for a description of the statement.

13.7.7.37 SHOW STATUS Statement

```
SHOW [GLOBAL | SESSION] STATUS
[LIKE 'pattern' | WHERE expr]
```

`SHOW STATUS` provides server status information (see [Section 5.1.10, “Server Status Variables”](#)). This statement does not require any privilege. It requires only the ability to connect to the server.

Status variable information is also available from these sources:

- Performance Schema tables. See [Section 27.12.15, “Performance Schema Status Variable Tables”](#).
- The `mysqladmin extended-status` command. See [Section 4.5.2, “mysqladmin — A MySQL Server Administration Program”](#).

For `SHOW STATUS`, a `LIKE` clause, if present, indicates which variable names to match. A `WHERE` clause can be given to select rows using more general conditions, as discussed in [Section 26.8, “Extensions to SHOW Statements”](#).

`SHOW STATUS` accepts an optional `GLOBAL` or `SESSION` variable scope modifier:

- With a `GLOBAL` modifier, the statement displays the global status values. A global status variable may represent status for some aspect of the server itself (for example, `Aborted_connects`), or the aggregated status over all connections to MySQL (for example, `Bytes_received` and `Bytes_sent`). If a variable has no global value, the session value is displayed.
- With a `SESSION` modifier, the statement displays the status variable values for the current connection. If a variable has no session value, the global value is displayed. `LOCAL` is a synonym for `SESSION`.
- If no modifier is present, the default is `SESSION`.

The scope for each status variable is listed at [Section 5.1.10, “Server Status Variables”](#).

Each invocation of the `SHOW STATUS` statement uses an internal temporary table and increments the global `Created_tmp_tables` value.

Partial output is shown here. The list of names and values may differ for your server. The meaning of each variable is given in [Section 5.1.10, “Server Status Variables”](#).

```
mysql> SHOW STATUS;
+-----+-----+
| Variable_name      | Value   |
+-----+-----+
| Aborted_clients    | 0       |
| Aborted_connects   | 0       |
| Bytes_received     | 155372598 |
| Bytes_sent          | 1176560426 |
| Connections         | 30023   |
| Created_tmp_disk_tables | 0       |
| Created_tmp_tables  | 8340    |
| Created_tmp_files   | 60      |
...
| Open_tables          | 1       |
| Open_files            | 2       |
| Open_streams          | 0       |
| Opened_tables         | 44600   |
| Questions             | 2026873 |
...
| Table_locks_immediate | 1920382 |
| Table_locks_waited   | 0       |
| Threads_cached        | 0       |
| Threads_created       | 30022   |
| Threads_connected     | 1       |
| Threads_running       | 1       |
| Uptime                | 80380   |
+-----+-----+
```

With a `LIKE` clause, the statement displays only rows for those variables with names that match the pattern:

```
mysql> SHOW STATUS LIKE 'Key%';
+-----+-----+
| Variable_name      | Value   |
+-----+-----+
| Key_blocks_used    | 14955   |
| Key_read_requests  | 96854827|
| Key_reads          | 162040  |
| Key_write_requests | 7589728 |
| Key_writes         | 3813196 |
+-----+-----+
```

13.7.7.38 SHOW TABLE STATUS Statement

```
SHOW TABLE STATUS
[ {FROM | IN} db_name ]
[LIKE 'pattern' | WHERE expr]
```

`SHOW TABLE STATUS` works like `SHOW TABLES`, but provides a lot of information about each non-Temporary table. You can also get this list using the `mysqlshow --status db_name` command. The `LIKE` clause, if present, indicates which table names to match. The `WHERE` clause can be given to select rows using more general conditions, as discussed in [Section 26.8, “Extensions to SHOW Statements”](#).

This statement also displays information about views.

`SHOW TABLE STATUS` output has these columns:

- `Name`

The name of the table.

- `Engine`

The storage engine for the table. See [Chapter 15, The InnoDB Storage Engine](#), and [Chapter 16, Alternative Storage Engines](#).

For partitioned tables, `Engine` shows the name of the storage engine used by all partitions.

- `Version`

This column is unused. With the removal of `.frm` files in MySQL 8.0, this column now reports a hardcoded value of `10`, which is the last `.frm` file version used in MySQL 5.7.

- `Row_format`

The row-storage format (`Fixed`, `Dynamic`, `Compressed`, `Redundant`, `Compact`). For `MyISAM` tables, `Dynamic` corresponds to what `myisamchk -dvv` reports as `Packed`.

- `Rows`

The number of rows. Some storage engines, such as `MyISAM`, store the exact count. For other storage engines, such as `InnoDB`, this value is an approximation, and may vary from the actual value by as much as 40% to 50%. In such cases, use `SELECT COUNT(*)` to obtain an accurate count.

The `Rows` value is `NULL` for `INFORMATION_SCHEMA` tables.

For `InnoDB` tables, the row count is only a rough estimate used in SQL optimization. (This is also true if the `InnoDB` table is partitioned.)

- `Avg_row_length`

The average row length.

- [Data_length](#)

For [MyISAM](#), [Data_length](#) is the length of the data file, in bytes.

For [InnoDB](#), [Data_length](#) is the approximate amount of space allocated for the clustered index, in bytes. Specifically, it is the clustered index size, in pages, multiplied by the [InnoDB](#) page size.

Refer to the notes at the end of this section for information regarding other storage engines.

- [Max_data_length](#)

For [MyISAM](#), [Max_data_length](#) is maximum length of the data file. This is the total number of bytes of data that can be stored in the table, given the data pointer size used.

Unused for [InnoDB](#).

Refer to the notes at the end of this section for information regarding other storage engines.

- [Index_length](#)

For [MyISAM](#), [Index_length](#) is the length of the index file, in bytes.

For [InnoDB](#), [Index_length](#) is the approximate amount of space allocated for non-clustered indexes, in bytes. Specifically, it is the sum of non-clustered index sizes, in pages, multiplied by the [InnoDB](#) page size.

Refer to the notes at the end of this section for information regarding other storage engines.

- [Data_free](#)

The number of allocated but unused bytes.

[InnoDB](#) tables report the free space of the tablespace to which the table belongs. For a table located in the shared tablespace, this is the free space of the shared tablespace. If you are using multiple tablespaces and the table has its own tablespace, the free space is for only that table. Free space means the number of bytes in completely free extents minus a safety margin. Even if free space displays as 0, it may be possible to insert rows as long as new extents need not be allocated.

For NDB Cluster, [Data_free](#) shows the space allocated on disk for, but not used by, a Disk Data table or fragment on disk. (In-memory data resource usage is reported by the [Data_length](#) column.)

For partitioned tables, this value is only an estimate and may not be absolutely correct. A more accurate method of obtaining this information in such cases is to query the [INFORMATION_SCHEMA PARTITIONS](#) table, as shown in this example:

```
SELECT SUM(DATA_FREE)
      FROM INFORMATION_SCHEMA.PARTITIONS
     WHERE TABLE_SCHEMA = 'mydb'
       AND TABLE_NAME    = 'mytable';
```

For more information, see [Section 26.3.21, “The INFORMATION_SCHEMA PARTITIONS Table”](#).

- [Auto_increment](#)

The next [AUTO_INCREMENT](#) value.

- [Create_time](#)

When the table was created.

- [Update_time](#)

When the data file was last updated. For some storage engines, this value is `NULL`. For example, [InnoDB](#) stores multiple tables in its [system tablespace](#) and the data file timestamp does not apply. Even with [file-per-table](#) mode with each [InnoDB](#) table in a separate `.ibd` file, [change buffering](#) can delay the write to the data file, so the file modification time is different from the time of the last insert, update, or delete. For [MyISAM](#), the data file timestamp is used; however, on Windows the timestamp is not updated by updates, so the value is inaccurate.

[Update_time](#) displays a timestamp value for the last [UPDATE](#), [INSERT](#), or [DELETE](#) performed on [InnoDB](#) tables that are not partitioned. For MVCC, the timestamp value reflects the [COMMIT](#) time, which is considered the last update time. Timestamps are not persisted when the server is restarted or when the table is evicted from the [InnoDB](#) data dictionary cache.

- [Check_time](#)

When the table was last checked. Not all storage engines update this time, in which case, the value is always `NULL`.

For partitioned [InnoDB](#) tables, [Check_time](#) is always `NULL`.

- [Collation](#)

The table default collation. The output does not explicitly list the table default character set, but the collation name begins with the character set name.

- [Checksum](#)

The live checksum value, if any.

- [Create_options](#)

Extra options used with [CREATE TABLE](#).

[Create_options](#) shows [partitioned](#) for a partitioned table.

Prior to MySQL 8.0.16, [Create_options](#) shows the [ENCRYPTION](#) clause specified for tables created in file-per-table tablespaces. As of MySQL 8.0.16, it shows the encryption clause for file-per-table tablespaces if the table is encrypted or if the specified encryption differs from the schema encryption. The encryption clause is not shown for tables created in general tablespaces. To identify encrypted file-per-table and general tablespaces, query the [INNODB_TABLESPACES_ENCRYPTION](#) column.

When creating a table with [strict mode](#) disabled, the storage engine's default row format is used if the specified row format is not supported. The actual row format of the table is reported in the [Row_format](#) column. [Create_options](#) shows the row format that was specified in the [CREATE TABLE](#) statement.

When altering the storage engine of a table, table options that are not applicable to the new storage engine are retained in the table definition to enable reverting the table with its previously defined options to the original storage engine, if necessary. [Create_options](#) may show retained options.

- [Comment](#)

The comment used when creating the table (or information as to why MySQL could not access the table information).

Notes

- For [InnoDB](#) tables, [SHOW TABLE STATUS](#) does not give accurate statistics except for the physical size reserved by the table. The row count is only a rough estimate used in SQL optimization.

- For `NDB` tables, the output of this statement shows appropriate values for the `Avg_row_length` and `Data_length` columns, with the exception that `BLOB` columns are not taken into account.
- For `NDB` tables, `Data_length` includes data stored in main memory only; the `Max_data_length` and `Data_free` columns apply to Disk Data.
- For NDB Cluster Disk Data tables, `Max_data_length` shows the space allocated for the disk part of a Disk Data table or fragment. (In-memory data resource usage is reported by the `Data_length` column.)
- For `MEMORY` tables, the `Data_length`, `Max_data_length`, and `Index_length` values approximate the actual amount of allocated memory. The allocation algorithm reserves memory in large amounts to reduce the number of allocation operations.
- For views, most columns displayed by `SHOW TABLE STATUS` are 0 or `NULL` except that `Name` indicates the view name, `Create_time` indicates the creation time, and `Comment` says `VIEW`.

Table information is also available from the `INFORMATION_SCHEMA TABLES` table. See [Section 26.3.38, “The INFORMATION_SCHEMA TABLES Table”](#).

13.7.7.39 SHOW TABLES Statement

```
SHOW [EXTENDED] [FULL] TABLES
[ {FROM | IN} db_name ]
[LIKE 'pattern' | WHERE expr]
```

`SHOW TABLES` lists the non-`TEMPORARY` tables in a given database. You can also get this list using the `mysqlshow db_name` command. The `LIKE` clause, if present, indicates which table names to match. The `WHERE` clause can be given to select rows using more general conditions, as discussed in [Section 26.8, “Extensions to SHOW Statements”](#).

Matching performed by the `LIKE` clause is dependent on the setting of the `lower_case_table_names` system variable.

The optional `EXTENDED` modifier causes `SHOW TABLES` to list hidden tables created by failed `ALTER TABLE` statements. These temporary tables have names beginning with `#sql` and can be dropped using `DROP TABLE`.

This statement also lists any views in the database. The optional `FULL` modifier causes `SHOW TABLES` to display a second output column with values of `BASE TABLE` for a table, `VIEW` for a view, or `SYSTEM VIEW` for an `INFORMATION_SCHEMA` table.

If you have no privileges for a base table or view, it does not show up in the output from `SHOW TABLES` or `mysqlshow db_name`.

Table information is also available from the `INFORMATION_SCHEMA TABLES` table. See [Section 26.3.38, “The INFORMATION_SCHEMA TABLES Table”](#).

13.7.7.40 SHOW TRIGGERS Statement

```
SHOW TRIGGERS
[ {FROM | IN} db_name ]
[LIKE 'pattern' | WHERE expr]
```

`SHOW TRIGGERS` lists the triggers currently defined for tables in a database (the default database unless a `FROM` clause is given). This statement returns results only for databases and tables for which you have the `TRIGGER` privilege. The `LIKE` clause, if present, indicates which table names (not trigger names) to match and causes the statement to display triggers for those tables. The `WHERE` clause can be given to select rows using more general conditions, as discussed in [Section 26.8, “Extensions to SHOW Statements”](#).

For the `ins_sum` trigger defined in [Section 25.3, “Using Triggers”](#), the output of `SHOW TRIGGERS` is as shown here:

```
mysql> SHOW TRIGGERS LIKE 'acc%\G
***** 1. row *****
Trigger: ins_sum
Event: INSERT
Table: account
Statement: SET @sum = @sum + NEW.amount
Timing: BEFORE
Created: 2018-08-08 10:10:12.61
sql_mode: ONLY_FULL_GROUP_BY,STRICT_TRANS_TABLES,
NO_ZERO_IN_DATE,NO_ZERO_DATE,
ERROR_FOR_DIVISION_BY_ZERO,
NO_ENGINE_SUBSTITUTION
Definer: me@localhost
character_set_client: utf8mb4
collation_connection: utf8mb4_0900_ai_ci
Database Collation: utf8mb4_0900_ai_ci
```

`SHOW TRIGGERS` output has these columns:

- [Trigger](#)

The name of the trigger.

- [Event](#)

The trigger event. This is the type of operation on the associated table for which the trigger activates. The value is [INSERT](#) (a row was inserted), [DELETE](#) (a row was deleted), or [UPDATE](#) (a row was modified).

- [Table](#)

The table for which the trigger is defined.

- [Statement](#)

The trigger body; that is, the statement executed when the trigger activates.

- [Timing](#)

Whether the trigger activates before or after the triggering event. The value is [BEFORE](#) or [AFTER](#).

- [Created](#)

The date and time when the trigger was created. This is a [TIMESTAMP\(2\)](#) value (with a fractional part in hundredths of seconds) for triggers.

- [sql_mode](#)

The SQL mode in effect when the trigger was created, and under which the trigger executes. For the permitted values, see [Section 5.1.11, “Server SQL Modes”](#).

- [Definer](#)

The account of the user who created the trigger, in '`user_name`'@'`host_name`' format.

- [character_set_client](#)

The session value of the `character_set_client` system variable when the trigger was created.

- [collation_connection](#)

The session value of the `collation_connection` system variable when the trigger was created.

- [Database Collation](#)

The collation of the database with which the trigger is associated.

Trigger information is also available from the [INFORMATION_SCHEMA TRIGGERS](#) table. See [Section 26.3.45, “The INFORMATION_SCHEMA TRIGGERS Table”](#).

13.7.7.41 SHOW VARIABLES Statement

```
SHOW [GLOBAL | SESSION] VARIABLES
    [LIKE 'pattern' | WHERE expr]
```

`SHOW VARIABLES` shows the values of MySQL system variables (see [Section 5.1.8, “Server System Variables”](#)). This statement does not require any privilege. It requires only the ability to connect to the server.

System variable information is also available from these sources:

- Performance Schema tables. See [Section 27.12.14, “Performance Schema System Variable Tables”](#).
- The `mysqladmin variables` command. See [Section 4.5.2, “mysqladmin — A MySQL Server Administration Program”](#).

For `SHOW VARIABLES`, a `LIKE` clause, if present, indicates which variable names to match. A `WHERE` clause can be given to select rows using more general conditions, as discussed in [Section 26.8, “Extensions to SHOW Statements”](#).

`SHOW VARIABLES` accepts an optional `GLOBAL` or `SESSION` variable scope modifier:

- With a `GLOBAL` modifier, the statement displays global system variable values. These are the values used to initialize the corresponding session variables for new connections to MySQL. If a variable has no global value, no value is displayed.
- With a `SESSION` modifier, the statement displays the system variable values that are in effect for the current connection. If a variable has no session value, the global value is displayed. `LOCAL` is a synonym for `SESSION`.
- If no modifier is present, the default is `SESSION`.

The scope for each system variable is listed at [Section 5.1.8, “Server System Variables”](#).

`SHOW VARIABLES` is subject to a version-dependent display-width limit. For variables with very long values that are not completely displayed, use `SELECT` as a workaround. For example:

```
SELECT @@GLOBAL.innodb_data_file_path;
```

Most system variables can be set at server startup (read-only variables such as `version_comment` are exceptions). Many can be changed at runtime with the `SET` statement. See [Section 5.1.9, “Using System Variables”](#), and [Section 13.7.6.1, “SET Syntax for Variable Assignment”](#).

Partial output is shown here. The list of names and values may differ for your server. [Section 5.1.8, “Server System Variables”](#), describes the meaning of each variable, and [Section 5.1.1, “Configuring the Server”](#), provides information about tuning them.

```
mysql> SHOW VARIABLES;
+-----+-----+
| Variable_name          | Value   |
+-----+-----+
| activate_all_roles_on_login | OFF
| auto_generate_certs      | ON
| auto_increment_increment | 1
| auto_increment_offset    | 1
| autocommit               | ON
| automatic_sp_privileges | ON
| avoid_temporal_upgrade  | OFF
| back_log                 | 151
```

SHOW Statements

basedir	/usr/
big_tables	OFF
bind_address	*
binlog_cache_size	32768
binlog_checksum	CRC32
binlog_direct_non_transactional_updates	OFF
binlog_error_action	ABORT_SERVER
binlog_expire_logs_seconds	2592000
binlog_format	ROW
binlog_group_commit_sync_delay	0
binlog_group_commit_sync_no_delay_count	0
binlog_gtid_simple_recovery	ON
binlog_max_flush_queue_time	0
binlog_order_commits	ON
binlog_row_image	FULL
binlog_row_metadata	MINIMAL
binlog_row_value_options	OFF
binlog_rows_query_log_events	32768
binlog_stmt_cache_size	25000
binlog_transaction_dependency_history_size	25000
binlog_transaction_dependency_tracking	COMMIT_ORDER
block_encryption_mode	aes-128-ecb
bulk_insert_buffer_size	8388608
 ...	
max_allowed_packet	67108864
max_binlog_cache_size	18446744073709547520
max_binlog_size	1073741824
max_binlog_stmt_cache_size	18446744073709547520
max_connect_errors	100
max_connections	151
max_delayed_threads	20
max_digest_length	1024
max_error_count	1024
max_execution_time	0
max_heap_table_size	16777216
max_insert_delayed_threads	20
max_join_size	18446744073709551615
 ...	
thread_handling	one-thread-per-connection
thread_stack	286720
time_zone	SYSTEM
timestamp	1530906638.765316
tls_version	TLSv1.2, TLSv1.3
tmp_table_size	16777216
tmpdir	/tmp
transaction_alloc_block_size	8192
transaction_allow_batching	OFF
transaction_isolation	REPEATABLE-READ
transaction_prealloc_size	4096
transaction_read_only	OFF
transaction_write_set_extraction	XXHASH64
unique_checks	ON
updatable_views_with_limit	YES
version	8.0.12
version_comment	MySQL Community Server - GPL
version_compile_machine	x86_64
version_compile_os	Linux
version_compile_zlib	1.2.11
wait_timeout	28800
warning_count	0
windowing_use_high_precision	ON

With a `LIKE` clause, the statement displays only rows for those variables with names that match the pattern. To obtain the row for a specific variable, use a `LIKE` clause as shown:

```
SHOW VARIABLES LIKE 'max_join_size';
SHOW SESSION VARIABLES LIKE 'max_join_size';
```

To get a list of variables whose name match a pattern, use the `%` wildcard character in a [LIKE clause](#):

```
SHOW VARIABLES LIKE '%size%';
SHOW GLOBAL VARIABLES LIKE '%size%';
```

Wildcard characters can be used in any position within the pattern to be matched. Strictly speaking, because `_` is a wildcard that matches any single character, you should escape it as `_` to match it literally. In practice, this is rarely necessary.

13.7.7.42 SHOW WARNINGS Statement

```
SHOW WARNINGS [LIMIT [offset,] row_count]
SHOW COUNT(*) WARNINGS
```

`SHOW WARNINGS` is a diagnostic statement that displays information about the conditions (errors, warnings, and notes) resulting from executing a statement in the current session. Warnings are generated for DML statements such as `INSERT`, `UPDATE`, and `LOAD DATA` as well as DDL statements such as `CREATE TABLE` and `ALTER TABLE`.

The `LIMIT` clause has the same syntax as for the `SELECT` statement. See [Section 13.2.13, “SELECT Statement”](#).

`SHOW WARNINGS` is also used following `EXPLAIN`, to display the extended information generated by `EXPLAIN`. See [Section 8.8.3, “Extended EXPLAIN Output Format”](#).

`SHOW WARNINGS` displays information about the conditions resulting from execution of the most recent nondiagnostic statement in the current session. If the most recent statement resulted in an error during parsing, `SHOW WARNINGS` shows the resulting conditions, regardless of statement type (diagnostic or nondiagnostic).

The `SHOW COUNT(*) WARNINGS` diagnostic statement displays the total number of errors, warnings, and notes. You can also retrieve this number from the `warning_count` system variable:

```
SHOW COUNT(*) WARNINGS;
SELECT @@warning_count;
```

A difference in these statements is that the first is a diagnostic statement that does not clear the message list. The second, because it is a `SELECT` statement is considered nondiagnostic and does clear the message list.

A related diagnostic statement, `SHOW ERRORS`, shows only error conditions (it excludes warnings and notes), and `SHOW COUNT(*) ERRORS` statement displays the total number of errors. See [Section 13.7.7.17, “SHOW ERRORS Statement”](#). `GET DIAGNOSTICS` can be used to examine information for individual conditions. See [Section 13.6.7.3, “GET DIAGNOSTICS Statement”](#).

Here is a simple example that shows data-conversion warnings for `INSERT`. The example assumes that strict SQL mode is disabled. With strict mode enabled, the warnings would become errors and terminate the `INSERT`.

```
mysql> CREATE TABLE t1 (a TINYINT NOT NULL, b CHAR(4));
Query OK, 0 rows affected (0.05 sec)

mysql> INSERT INTO t1 VALUES(10,'mysql'), (NULL,'test'), (300,'xyz');
Query OK, 3 rows affected, 3 warnings (0.00 sec)
Records: 3  Duplicates: 0  Warnings: 3

mysql> SHOW WARNINGS\G
***** 1. row *****
  Level: Warning
  Code: 1265
Message: Data truncated for column 'b' at row 1
***** 2. row *****
  Level: Warning
  Code: 1048
```

```

Message: Column 'a' cannot be null
***** 3. row *****
Level: Warning
Code: 1264
Message: Out of range value for column 'a' at row 3
3 rows in set (0.00 sec)

```

The `max_error_count` system variable controls the maximum number of error, warning, and note messages for which the server stores information, and thus the number of messages that `SHOW WARNINGS` displays. To change the number of messages the server can store, change the value of `max_error_count`.

`max_error_count` controls only how many messages are stored, not how many are counted. The value of `warning_count` is not limited by `max_error_count`, even if the number of messages generated exceeds `max_error_count`. The following example demonstrates this. The `ALTER TABLE` statement produces three warning messages (strict SQL mode is disabled for the example to prevent an error from occurring after a single conversion issue). Only one message is stored and displayed because `max_error_count` has been set to 1, but all three are counted (as shown by the value of `warning_count`):

```

mysql> SHOW VARIABLES LIKE 'max_error_count';
+-----+-----+
| Variable_name | Value |
+-----+-----+
| max_error_count | 1024 |
+-----+-----+
1 row in set (0.00 sec)

mysql> SET max_error_count=1, sql_mode = '';
Query OK, 0 rows affected (0.00 sec)

mysql> ALTER TABLE t1 MODIFY b CHAR;
Query OK, 3 rows affected, 3 warnings (0.00 sec)
Records: 3  Duplicates: 0  Warnings: 3

mysql> SHOW WARNINGS;
+-----+-----+-----+
| Level | Code | Message |
+-----+-----+-----+
| Warning | 1263 | Data truncated for column 'b' at row 1 |
+-----+-----+-----+
1 row in set (0.00 sec)

mysql> SELECT @@warning_count;
+-----+
| @@warning_count |
+-----+
|            3 |
+-----+
1 row in set (0.01 sec)

```

To disable message storage, set `max_error_count` to 0. In this case, `warning_count` still indicates how many warnings occurred, but messages are not stored and cannot be displayed.

The `sql_notes` system variable controls whether note messages increment `warning_count` and whether the server stores them. By default, `sql_notes` is 1, but if set to 0, notes do not increment `warning_count` and the server does not store them:

```

mysql> SET sql_notes = 1;
mysql> DROP TABLE IF EXISTS test.no_such_table;
Query OK, 0 rows affected, 1 warning (0.00 sec)
mysql> SHOW WARNINGS;
+-----+-----+-----+
| Level | Code | Message |
+-----+-----+-----+
| Note  | 1051 | Unknown table 'test.no_such_table' |
+-----+-----+-----+
1 row in set (0.00 sec)

```

```
mysql> SET sql_notes = 0;
mysql> DROP TABLE IF EXISTS test.no_such_table;
Query OK, 0 rows affected (0.00 sec)
mysql> SHOW WARNINGS;
Empty set (0.00 sec)
```

The MySQL server sends to each client a count indicating the total number of errors, warnings, and notes resulting from the most recent statement executed by that client. From the C API, this value can be obtained by calling `mysql_warning_count()`. See [mysql_warning_count\(\)](#).

In the `mysql` client, you can enable and disable automatic warnings display using the `warnings` and `nowarning` commands, respectively, or their shortcuts, `\w` and `\W` (see [Section 4.5.1.2, “mysql Client Commands”](#)). For example:

```
mysql> \w
Show warnings enabled.
mysql> SELECT 1/0;
+-----+
| 1/0  |
+-----+
| NULL |
+-----+
1 row in set, 1 warning (0.03 sec)

Warning (Code 1365): Division by 0
mysql> \W
Show warnings disabled.
```

13.7.8 Other Administrative Statements

13.7.8.1 BINLOG Statement

```
BINLOG 'str'
```

`BINLOG` is an internal-use statement. It is generated by the `mysqlbinlog` program as the printable representation of certain events in binary log files. (See [Section 4.6.9, “mysqlbinlog — Utility for Processing Binary Log Files”](#).) The '`str`' value is a base 64-encoded string that the server decodes to determine the data change indicated by the corresponding event.

To execute `BINLOG` statements when applying `mysqlbinlog` output, a user account requires the `BINLOG_ADMIN` privilege (or the deprecated `SUPER` privilege), or the `REPLICATION_APPLIER` privilege plus the appropriate privileges to execute each log event.

This statement can execute only format description events and row events.

13.7.8.2 CACHE INDEX Statement

```
CACHE INDEX {
    tbl_index_list [, tbl_index_list] ...
    | tbl_name PARTITION (partition_list)
}
IN key_cache_name

tbl_index_list:
tbl_name [{INDEX|KEY} (index_name[, index_name] ...)]

partition_list:
    partition_name[, partition_name] ...
    | ALL
}
```

The `CACHE INDEX` statement assigns table indexes to a specific key cache. It applies only to `MyISAM` tables, including partitioned `MyISAM` tables. After the indexes have been assigned, they can be preloaded into the cache if desired with `LOAD INDEX INTO CACHE`.

The following statement assigns indexes from the tables `t1`, `t2`, and `t3` to the key cache named `hot_cache`:

```
mysql> CACHE INDEX t1, t2, t3 IN hot_cache;
+-----+-----+-----+
| Table | Op      | Msg_type | Msg_text |
+-----+-----+-----+
| test.t1 | assign_to_keycache | status   | OK       |
| test.t2 | assign_to_keycache | status   | OK       |
| test.t3 | assign_to_keycache | status   | OK       |
+-----+-----+-----+
```

The syntax of `CACHE INDEX` enables you to specify that only particular indexes from a table should be assigned to the cache. However, the implementation assigns all the table's indexes to the cache, so there is no reason to specify anything other than the table name.

The key cache referred to in a `CACHE INDEX` statement can be created by setting its size with a parameter setting statement or in the server parameter settings. For example:

```
SET GLOBAL keycache1.key_buffer_size=128*1024;
```

Key cache parameters are accessed as members of a structured system variable. See [Section 5.1.9.5, “Structured System Variables”](#).

A key cache must exist before you assign indexes to it, or an error occurs:

```
mysql> CACHE INDEX t1 IN non_existent_cache;
ERROR 1284 (HY000): Unknown key cache 'non_existent_cache'
```

By default, table indexes are assigned to the main (default) key cache created at the server startup. When a key cache is destroyed, all indexes assigned to it are reassigned to the default key cache.

Index assignment affects the server globally: If one client assigns an index to a given cache, this cache is used for all queries involving the index, no matter which client issues the queries.

`CACHE INDEX` is supported for partitioned `MyISAM` tables. You can assign one or more indexes for one, several, or all partitions to a given key cache. For example, you can do the following:

```
CREATE TABLE pt (c1 INT, c2 VARCHAR(50), INDEX i(c1))
  ENGINE=MyISAM
  PARTITION BY HASH(c1)
  PARTITIONS 4;

SET GLOBAL kc_fast.key_buffer_size = 128 * 1024;
SET GLOBAL kc_slow.key_buffer_size = 128 * 1024;

CACHE INDEX pt PARTITION (p0) IN kc_fast;
CACHE INDEX pt PARTITION (p1, p3) IN kc_slow;
```

The previous set of statements performs the following actions:

- Creates a partitioned table with 4 partitions; these partitions are automatically named `p0`, ..., `p3`; this table has an index named `i` on column `c1`.
- Creates 2 key caches named `kc_fast` and `kc_slow`
- Assigns the index for partition `p0` to the `kc_fast` key cache and the index for partitions `p1` and `p3` to the `kc_slow` key cache; the index for the remaining partition (`p2`) uses the server's default key cache.

If you wish instead to assign the indexes for all partitions in table `pt` to a single key cache named `kc_all`, you can use either of the following two statements:

```
CACHE INDEX pt PARTITION (ALL) IN kc_all;
```

```
CACHE INDEX pt IN kc_all;
```

The two statements just shown are equivalent, and issuing either one has exactly the same effect. In other words, if you wish to assign indexes for all partitions of a partitioned table to the same key cache, the `PARTITION (ALL)` clause is optional.

When assigning indexes for multiple partitions to a key cache, the partitions need not be contiguous, and you need not list their names in any particular order. Indexes for any partitions not explicitly assigned to a key cache automatically use the server default key cache.

Index preloading is also supported for partitioned MyISAM tables. For more information, see [Section 13.7.8.5, “LOAD INDEX INTO CACHE Statement”](#).

13.7.8.3 FLUSH Statement

```
FLUSH [NO_WRITE_TO_BINLOG | LOCAL] {
    flush_option [, flush_option] ...
    | tables_option
}

flush_option: {
    BINARY LOGS
    | ENGINE LOGS
    | ERROR LOGS
    | GENERAL LOGS
    | HOSTS
    | LOGS
    | PRIVILEGES
    | OPTIMIZER_COSTS
    | RELAY LOGS [FOR CHANNEL channel]
    | SLOW LOGS
    | STATUS
    | USER_RESOURCES
}

tables_option: {
    TABLES
    | TABLES tbl_name [, tbl_name] ...
    | TABLES WITH READ LOCK
    | TABLES tbl_name [, tbl_name] ... WITH READ LOCK
    | TABLES tbl_name [, tbl_name] ... FOR EXPORT
}
```

The `FLUSH` statement has several variant forms that clear or reload various internal caches, flush tables, or acquire locks. Each `FLUSH` operation requires the privileges indicated in its description.



Note

It is not possible to issue `FLUSH` statements within stored functions or triggers. However, you may use `FLUSH` in stored procedures, so long as these are not called from stored functions or triggers. See [Section 25.8, “Restrictions on Stored Programs”](#).

By default, the server writes `FLUSH` statements to the binary log so that they replicate to replicas. To suppress logging, specify the optional `NO_WRITE_TO_BINLOG` keyword or its alias `LOCAL`.



Note

`FLUSH LOGS`, `FLUSH BINARY LOGS`, `FLUSH TABLES WITH READ LOCK` (with or without a table list), and `FLUSH TABLES tbl_name ... FOR EXPORT` are not written to the binary log in any case because they would cause problems if replicated to a replica.

The `FLUSH` statement causes an implicit commit. See [Section 13.3.3, “Statements That Cause an Implicit Commit”](#).

The `mysqladmin` utility provides a command-line interface to some flush operations, using commands such as `flush-hosts`, `flush-logs`, `flush-privileges`, `flush-status`, and `flush-tables`. See [Section 4.5.2, “mysqladmin — A MySQL Server Administration Program”](#).

Sending a `SIGHUP` or `SIGUSR1` signal to the server causes several flush operations to occur that are similar to various forms of the `FLUSH` statement. Signals can be sent by the `root` system account or the system account that owns the server process. This enables the flush operations to be performed without having to connect to the server, which requires a MySQL account that has privileges sufficient for those operations. See [Section 4.10, “Unix Signal Handling in MySQL”](#).

The `RESET` statement is similar to `FLUSH`. See [Section 13.7.8.6, “RESET Statement”](#), for information about using `RESET` with replication.

The following list describes the permitted `FLUSH` statement `flush_option` values. For descriptions of the permitted `tables_option` values, see [FLUSH TABLES Syntax](#).

- [FLUSH BINARY LOGS](#)

Closes and reopens any binary log file to which the server is writing. If binary logging is enabled, the sequence number of the binary log file is incremented by one relative to the previous file.

This operation requires the `RELOAD` privilege.

- [FLUSH ENGINE LOGS](#)

Closes and reopens any flushable logs for installed storage engines. This causes `InnoDB` to flush its logs to disk.

This operation requires the `RELOAD` privilege.

- [FLUSH ERROR LOGS](#)

Closes and reopens any error log file to which the server is writing.

This operation requires the `RELOAD` privilege.

- [FLUSH GENERAL LOGS](#)

Closes and reopens any general query log file to which the server is writing.

This operation requires the `RELOAD` privilege.

This operation has no effect on tables used for the general query log (see [Section 5.4.1, “Selecting General Query Log and Slow Query Log Output Destinations”](#)).

- [FLUSH HOSTS](#)

Empties the host cache and the Performance Schema `host_cache` table that exposes the cache contents, and unblocks any blocked hosts.

This operation requires the `RELOAD` privilege.

For information about why host cache flushing might be advisable or desirable, see [Section 5.1.12.3, “DNS Lookups and the Host Cache”](#).



Note

`FLUSH HOSTS` is deprecated as of MySQL 8.0.23; expect it to be removed in a future MySQL release. Instead, truncate the Performance Schema `host_cache` table:

```
TRUNCATE TABLE performance_schema.host_cache;
```

The `TRUNCATE TABLE` operation requires the `DROP` privilege for the table rather than the `RELOAD` privilege. You should be aware that the `TRUNCATE TABLE` statement is not written to the binary log. To obtain the same behavior from `FLUSH HOSTS`, specify `NO_WRITE_TO_BINLOG` or `LOCAL` as part of the statement.

- **`FLUSH LOGS`**

Closes and reopens any log file to which the server is writing.

This operation requires the `RELOAD` privilege.

The effect of this operation is equivalent to the combined effects of these operations:

```
FLUSH BINARY LOGS
FLUSH ENGINE LOGS
FLUSH ERROR LOGS
FLUSH GENERAL LOGS
FLUSH RELAY LOGS
FLUSH SLOW LOGS
```

- **`FLUSH OPTIMIZER_COSTS`**

Re-reads the cost model tables so that the optimizer starts using the current cost estimates stored in them.

This operation requires the `FLUSH_OPTIMIZER_COSTS` or `RELOAD` privilege.

The server writes a warning to the error log for any unrecognized cost model table entries. For information about these tables, see [Section 8.9.5, “The Optimizer Cost Model”](#). This operation affects only sessions that begin subsequent to the flush. Existing sessions continue to use the cost estimates that were current when they began.

- **`FLUSH PRIVILEGES`**

Re-reads the privileges from the grant tables in the `mysql` system schema. As part of this operation, the server reads the `global_grants` table containing dynamic privilege assignments and registers any unregistered privileges found there.

Reloading the grant tables is necessary to enable updates to MySQL privileges and users only if you make such changes directly to the grant tables; it is not needed for account management statements such as `GRANT` or `REVOKE`, which take effect immediately. See [Section 6.2.13, “When Privilege Changes Take Effect”](#), for more information.

This operation requires the `RELOAD` privilege.

If the `--skip-grant-tables` option was specified at server startup to disable the MySQL privilege system, `FLUSH PRIVILEGES` provides a way to enable the privilege system at runtime.

Resets failed-login tracking (or enables it if the server was started with `--skip-grant-tables`) and unlocks any temporarily locked accounts. See [Section 6.2.15, “Password Management”](#).

Frees memory cached by the server as a result of `GRANT`, `CREATE USER`, `CREATE SERVER`, and `INSTALL PLUGIN` statements. This memory is not released by the corresponding `REVOKE`, `DROP USER`, `DROP SERVER`, and `UNINSTALL PLUGIN` statements, so for a server that executes many instances of the statements that cause caching, there is an increase in cached memory use unless it is freed with `FLUSH PRIVILEGES`.

Clears the in-memory cache used by the `caching_sha2_password` authentication plugin. See [Cache Operation for SHA-2 Pluggable Authentication](#).

- **`FLUSH RELAY LOGS [FOR CHANNEL channel]`**

Closes and reopens any relay log file to which the server is writing. If relay logging is enabled, the sequence number of the relay log file is incremented by one relative to the previous file.

This operation requires the `RELOAD` privilege.

The `FOR CHANNEL channel` clause enables you to name which replication channel the operation applies to. Execute `FLUSH RELAY LOGS FOR CHANNEL channel` to flush the relay log for a specific replication channel. If no channel is named and no extra replication channels exist, the operation applies to the default channel. If no channel is named and multiple replication channels exist, the operation applies to all replication channels. For more information, see [Section 17.2.2, “Replication Channels”](#).

- **`FLUSH SLOW LOGS`**

Closes and reopens any slow query log file to which the server is writing.

This operation requires the `RELOAD` privilege.

This operation has no effect on tables used for the slow query log (see [Section 5.4.1, “Selecting General Query Log and Slow Query Log Output Destinations”](#)).

- **`FLUSH STATUS`**

Flushes status indicators.

This operation adds the current thread's session status variable values to the global values and resets the session values to zero. Some global variables may be reset to zero as well. It also resets the counters for key caches (default and named) to zero and sets `Max_used_connections` to the current number of open connections. This information may be of use when debugging a query. See [Section 1.5, “How to Report Bugs or Problems”](#).

`FLUSH STATUS` is unaffected by `read_only` or `super_read_only`, and is always written to the binary log.

This operation requires the `FLUSH_STATUS` or `RELOAD` privilege.

- **`FLUSH_USER_RESOURCES`**

Resets all per-hour user resource indicators to zero.

This operation requires the `FLUSH_USER_RESOURCES` or `RELOAD` privilege.

Resetting resource indicators enables clients that have reached their hourly connection, query, or update limits to resume activity immediately. `FLUSH_USER_RESOURCES` does not apply to the limit on maximum simultaneous connections that is controlled by the `max_user_connections` system variable. See [Section 6.2.21, “Setting Account Resource Limits”](#).

FLUSH TABLES Syntax

`FLUSH TABLES` flushes tables, and, depending on the variant used, acquires locks. Any `TABLES` variant used in a `FLUSH` statement must be the only option used. `FLUSH TABLE` is a synonym for `FLUSH TABLES`.



Note

The descriptions here that indicate tables are flushed by closing them apply differently for `InnoDB`, which flushes table contents to disk but leaves them open. This still permits table files to be copied while the tables are open, as long as other activity does not modify them.

- **`FLUSH TABLES`**

Closes all open tables, forces all tables in use to be closed, and flushes the prepared statement cache.

This operation requires the `FLUSH_TABLES` or `RELOAD` privilege.

For information about prepared statement caching, see [Section 8.10.3, “Caching of Prepared Statements and Stored Programs”](#).

`FLUSH TABLES` is not permitted when there is an active `LOCK TABLES ... READ`. To flush and lock tables, use `FLUSH TABLES tbl_name ... WITH READ LOCK` instead.

- `FLUSH TABLES tbl_name [, tbl_name] ...`

With a list of one or more comma-separated table names, this operation is like `FLUSH TABLES` with no names except that the server flushes only the named tables. If a named table does not exist, no error occurs.

This operation requires the `FLUSH_TABLES` or `RELOAD` privilege.

- `FLUSH TABLES WITH READ LOCK`

Closes all open tables and locks all tables for all databases with a global read lock.

This operation requires the `FLUSH_TABLES` or `RELOAD` privilege.

This operation is a very convenient way to get backups if you have a file system such as Veritas or ZFS that can take snapshots in time. Use `UNLOCK TABLES` to release the lock.

`FLUSH TABLES WITH READ LOCK` acquires a global read lock rather than table locks, so it is not subject to the same behavior as `LOCK TABLES` and `UNLOCK TABLES` with respect to table locking and implicit commits:

- `UNLOCK TABLES` implicitly commits any active transaction only if any tables currently have been locked with `LOCK TABLES`. The commit does not occur for `UNLOCK TABLES` following `FLUSH TABLES WITH READ LOCK` because the latter statement does not acquire table locks.
- Beginning a transaction causes table locks acquired with `LOCK TABLES` to be released, as though you had executed `UNLOCK TABLES`. Beginning a transaction does not release a global read lock acquired with `FLUSH TABLES WITH READ LOCK`.

`FLUSH TABLES WITH READ LOCK` does not prevent the server from inserting rows into the log tables (see [Section 5.4.1, “Selecting General Query Log and Slow Query Log Output Destinations”](#)).

- `FLUSH TABLES tbl_name [, tbl_name] ... WITH READ LOCK`

Flushes and acquires read locks for the named tables.

This operation requires the `FLUSH_TABLES` or `RELOAD` privilege. Because it acquires table locks, it also requires the `LOCK_TABLES` privilege for each table.

The operation first acquires exclusive metadata locks for the tables, so it waits for transactions that have those tables open to complete. Then the operation flushes the tables from the table cache, reopens the tables, acquires table locks (like `LOCK TABLES ... READ`), and downgrades the

metadata locks from exclusive to shared. After the operation acquires locks and downgrades the metadata locks, other sessions can read but not modify the tables.

This operation applies only to existing base (non-[TEMPORARY](#)) tables. If a name refers to a base table, that table is used. If it refers to a [TEMPORARY](#) table, it is ignored. If a name applies to a view, an [ER_WRONG_OBJECT](#) error occurs. Otherwise, an [ER_NO_SUCH_TABLE](#) error occurs.

Use [UNLOCK TABLES](#) to release the locks, [LOCK TABLES](#) to release the locks and acquire other locks, or [START TRANSACTION](#) to release the locks and begin a new transaction.

This [FLUSH TABLES](#) variant enables tables to be flushed and locked in a single operation. It provides a workaround for the restriction that [FLUSH TABLES](#) is not permitted when there is an active [LOCK TABLES ... READ](#).

This operation does not perform an implicit [UNLOCK TABLES](#), so an error results if you perform the operation while there is any active [LOCK TABLES](#) or use it a second time without first releasing the locks acquired.

If a flushed table was opened with [HANDLER](#), the handler is implicitly flushed and loses its position.

- [FLUSH TABLES *tbl_name \[, ...\]* ... FOR EXPORT](#)

This [FLUSH TABLES](#) variant applies to [InnoDB](#) tables. It ensures that changes to the named tables have been flushed to disk so that binary table copies can be made while the server is running.

This operation requires the [FLUSH_TABLES](#) or [RELOAD](#) privilege. Because it acquires locks on tables in preparation for exporting them, it also requires the [LOCK TABLES](#) and [SELECT](#) privileges for each table.

The operation works like this:

1. It acquires shared metadata locks for the named tables. The operation blocks as long as other sessions have active transactions that have modified those tables or hold table locks for them. When the locks have been acquired, the operation blocks transactions that attempt to update the tables, while permitting read-only operations to continue.
2. It checks whether all storage engines for the tables support [FOR EXPORT](#). If any do not, an [ER_ILLEGAL_HA](#) error occurs and the operation fails.
3. The operation notifies the storage engine for each table to make the table ready for export. The storage engine must ensure that any pending changes are written to disk.
4. The operation puts the session in lock-tables mode so that the metadata locks acquired earlier are not released when the [FOR EXPORT](#) operation completes.

This operation applies only to existing base (non-[TEMPORARY](#)) tables. If a name refers to a base table, that table is used. If it refers to a [TEMPORARY](#) table, it is ignored. If a name applies to a view, an [ER_WRONG_OBJECT](#) error occurs. Otherwise, an [ER_NO_SUCH_TABLE](#) error occurs.

[InnoDB](#) supports [FOR EXPORT](#) for tables that have their own [.ibd](#) file file (that is, tables created with the [innodb_file_per_table](#) setting enabled). [InnoDB](#) ensures when notified by the [FOR EXPORT](#) operation that any changes have been flushed to disk. This permits a binary copy of table contents to be made while the [FOR EXPORT](#) operation is in effect because the [.ibd](#) file is transaction consistent and can be copied while the server is running. [FOR EXPORT](#) does not apply to [InnoDB](#) system tablespace files, or to [InnoDB](#) tables that have [FULLTEXT](#) indexes.

[FLUSH TABLES ...FOR EXPORT](#) is supported for partitioned [InnoDB](#) tables.

When notified by [FOR EXPORT](#), [InnoDB](#) writes to disk certain kinds of data that is normally held in memory or in separate disk buffers outside the tablespace files. For each table, [InnoDB](#) also

produces a file named `table_name.cfg` in the same database directory as the table. The `.cfg` file contains metadata needed to reimport the tablespace files later, into the same or different server.

When the `FOR EXPORT` operation completes, InnoDB has flushed all `dirty pages` to the table data files. Any `change buffer` entries are merged prior to flushing. At this point, the tables are locked and quiescent: The tables are in a transactionally consistent state on disk and you can copy the `.ibd` tablespace files along with the corresponding `.cfg` files to get a consistent snapshot of those tables.

For the procedure to reimport the copied table data into a MySQL instance, see [Section 15.6.1.3, “Importing InnoDB Tables”](#).

After you are done with the tables, use `UNLOCK TABLES` to release the locks, `LOCK TABLES` to release the locks and acquire other locks, or `START TRANSACTION` to release the locks and begin a new transaction.

While any of these statements is in effect within the session, attempts to use `FLUSH TABLES ... FOR EXPORT` produce an error:

```
FLUSH TABLES ... WITH READ LOCK
FLUSH TABLES ... FOR EXPORT
LOCK TABLES ... READ
LOCK TABLES ... WRITE
```

While `FLUSH TABLES ... FOR EXPORT` is in effect within the session, attempts to use any of these statements produce an error:

```
FLUSH TABLES WITH READ LOCK
FLUSH TABLES ... WITH READ LOCK
FLUSH TABLES ... FOR EXPORT
```

13.7.8.4 KILL Statement

```
KILL [CONNECTION | QUERY] processlist_id
```

Each connection to `mysqld` runs in a separate thread. You can kill a thread with the `KILL processlist_id` statement.

Thread processlist identifiers can be determined from the `ID` column of the `INFORMATION_SCHEMA PROCESSLIST` table, the `Id` column of `SHOW PROCESSLIST` output, and the `PROCESSLIST_ID` column of the Performance Schema `threads` table. The value for the current thread is returned by the `CONNECTION_ID()` function.

`KILL` permits an optional `CONNECTION` or `QUERY` modifier:

- `KILL CONNECTION` is the same as `KILL` with no modifier: It terminates the connection associated with the given `processlist_id`, after terminating any statement the connection is executing.
- `KILL QUERY` terminates the statement the connection is currently executing, but leaves the connection itself intact.

The ability to see which threads are available to be killed depends on the `PROCESS` privilege:

- Without `PROCESS`, you can see only your own threads.
- With `PROCESS`, you can see all threads.

The ability to kill threads and statements depends on the `CONNECTION_ADMIN` privilege and the deprecated `SUPER` privilege:

- Without `CONNECTION_ADMIN` or `SUPER`, you can kill only your own threads and statements.
- With `CONNECTION_ADMIN` or `SUPER`, you can kill all threads and statements, except that to affect a thread or statement that is executing with the `SYSTEM_USER` privilege, your own session must additionally have the `SYSTEM_USER` privilege.

You can also use the `mysqladmin processlist` and `mysqladmin kill` commands to examine and kill threads.

When you use `KILL`, a thread-specific kill flag is set for the thread. In most cases, it might take some time for the thread to die because the kill flag is checked only at specific intervals:

- During `SELECT` operations, for `ORDER BY` and `GROUP BY` loops, the flag is checked after reading a block of rows. If the kill flag is set, the statement is aborted.
- `ALTER TABLE` operations that make a table copy check the kill flag periodically for each few copied rows read from the original table. If the kill flag was set, the statement is aborted and the temporary table is deleted.

The `KILL` statement returns without waiting for confirmation, but the kill flag check aborts the operation within a reasonably small amount of time. Aborting the operation to perform any necessary cleanup also takes some time.

- During `UPDATE` or `DELETE` operations, the kill flag is checked after each block read and after each updated or deleted row. If the kill flag is set, the statement is aborted. If you are not using transactions, the changes are not rolled back.
- `GET_LOCK()` aborts and returns `NULL`.
- If the thread is in the table lock handler (state: `Locked`), the table lock is quickly aborted.
- If the thread is waiting for free disk space in a write call, the write is aborted with a “disk full” error message.
- `EXPLAIN ANALYZE` aborts and prints the first row of output. This works in MySQL 8.0.20 and later.



Warning

Killing a `REPAIR TABLE` or `OPTIMIZE TABLE` operation on a `MyISAM` table results in a table that is corrupted and unusable. Any reads or writes to such a table fail until you optimize or repair it again (without interruption).

13.7.8.5 LOAD INDEX INTO CACHE Statement

```
LOAD INDEX INTO CACHE
  tbl_index_list [, tbl_index_list] ...

tbl_index_list:
  tbl_name
    [PARTITION (partition_list)]
    [{INDEX|KEY} (index_name[, index_name] ...)]
    [IGNORE LEAVES]

partition_list:
  partition_name[, partition_name] ...
  | ALL
}
```

The `LOAD INDEX INTO CACHE` statement preloads a table index into the key cache to which it has been assigned by an explicit `CACHE INDEX` statement, or into the default key cache otherwise.

`LOAD INDEX INTO CACHE` applies only to `MyISAM` tables, including partitioned `MyISAM` tables. In addition, indexes on partitioned tables can be preloaded for one, several, or all partitions.

The `IGNORE LEAVES` modifier causes only blocks for the nonleaf nodes of the index to be preloaded.

`IGNORE LEAVES` is also supported for partitioned `MyISAM` tables.

The following statement preloads nodes (index blocks) of indexes for the tables `t1` and `t2`:

```
mysql> LOAD INDEX INTO CACHE t1, t2 IGNORE LEAVES;
```

Table	Op	Msg_type	Msg_text
test.t1	preload_keys	status	OK
test.t2	preload_keys	status	OK

This statement preloads all index blocks from `t1`. It preloads only blocks for the nonleaf nodes from `t2`.

The syntax of `LOAD INDEX INTO CACHE` enables you to specify that only particular indexes from a table should be preloaded. However, the implementation preloads all the table's indexes into the cache, so there is no reason to specify anything other than the table name.

It is possible to preload indexes on specific partitions of partitioned MyISAM tables. For example, of the following 2 statements, the first preloads indexes for partition `p0` of a partitioned table `pt`, while the second preloads the indexes for partitions `p1` and `p3` of the same table:

```
LOAD INDEX INTO CACHE pt PARTITION (p0);
LOAD INDEX INTO CACHE pt PARTITION (p1, p3);
```

To preload the indexes for all partitions in table `pt`, you can use either of the following two statements:

```
LOAD INDEX INTO CACHE pt PARTITION (ALL);
LOAD INDEX INTO CACHE pt;
```

The two statements just shown are equivalent, and issuing either one has exactly the same effect. In other words, if you wish to preload indexes for all partitions of a partitioned table, the `PARTITION (ALL)` clause is optional.

When preloading indexes for multiple partitions, the partitions need not be contiguous, and you need not list their names in any particular order.

`LOAD INDEX INTO CACHE ... IGNORE LEAVES` fails unless all indexes in a table have the same block size. To determine index block sizes for a table, use `myisamchk -dv` and check the `Blocksize` column.

13.7.8.6 RESET Statement

```
RESET reset_option [, reset_option] ...
reset_option: {
    MASTER
    | REPLICA
    | SLAVE
}
```

The `RESET` statement is used to clear the state of various server operations. You must have the `RELOAD` privilege to execute `RESET`.

For information about the `RESET PERSIST` statement that removes persisted global system variables, see [Section 13.7.8.7, “RESET PERSIST Statement”](#).

`RESET` acts as a stronger version of the `FLUSH` statement. See [Section 13.7.8.3, “FLUSH Statement”](#).

The `RESET` statement causes an implicit commit. See [Section 13.3.3, “Statements That Cause an Implicit Commit”](#).

The following list describes the permitted `RESET` statement `reset_option` values:

- `RESET MASTER`

Deletes all binary logs listed in the index file, resets the binary log index file to be empty, and creates a new binary log file.

- `RESET REPLICA`

Makes the replica forget its replication position in the source binary logs. Also resets the relay log by deleting any existing relay log files and beginning a new one. Use `RESET REPLICA` in place of `RESET SLAVE` from MySQL 8.0.22.

13.7.8.7 RESET PERSIST Statement

```
RESET PERSIST [[IF EXISTS] system_var_name]
```

`RESET PERSIST` removes persisted global system variable settings from the `mysqld-auto.cnf` option file in the data directory. Removing a persisted system variable causes the variable no longer to be initialized from `mysqld-auto.cnf` at server startup. For more information about persisting system variables and the `mysqld-auto.cnf` file, see [Section 5.1.9.3, “Persisted System Variables”](#).

Prior to MySQL 8.0.32, this statement did not work with variables whose name contained a dot character (.), such as `MyISAM` multiple key cache variables and variables registered by components. (Bug #33417357)

The privileges required for `RESET PERSIST` depend on the type of system variable to be removed:

- For dynamic system variables, this statement requires the `SYSTEM_VARIABLES_ADMIN` privilege (or the deprecated `SUPER` privilege).
- For read-only system variables, this statement requires the `SYSTEM_VARIABLES_ADMIN` and `PERSIST_RO_VARIABLES_ADMIN` privileges.

See [Section 5.1.9.1, “System Variable Privileges”](#).

Depending on whether the variable name and `IF EXISTS` clauses are present, the `RESET PERSIST` statement has these forms:

- To remove all persisted variables from `mysqld-auto.cnf`, use `RESET PERSIST` without naming any system variable:

```
RESET PERSIST;
```

You must have privileges for removing both dynamic and read-only system variables if `mysqld-auto.cnf` contains both kinds of variables.

- To remove a specific persisted variable from `mysqld-auto.cnf`, name it in the statement:

```
RESET PERSIST system_var_name;
```

This includes plugin system variables, even if the plugin is not currently installed. If the variable is not present in the file, an error occurs.

- To remove a specific persisted variable from `mysqld-auto.cnf`, but produce a warning rather than an error if the variable is not present in the file, add an `IF EXISTS` clause to the previous syntax:

```
RESET PERSIST IF EXISTS system_var_name;
```

`RESET PERSIST` is not affected by the value of the `persisted_globals_load` system variable.

`RESET PERSIST` affects the contents of the Performance Schema `persisted_variables` table because the table contents correspond to the contents of the `mysqld-auto.cnf` file. On the other hand, because `RESET PERSIST` does not change variable values, it has no effect on the contents of the Performance Schema `variables_info` table until the server is restarted.

For information about `RESET` statement variants that clear the state of other server operations, see [Section 13.7.8.6, “RESET Statement”](#).

13.7.8.8 RESTART Statement

```
RESTART
```

This statement stops and restarts the MySQL server. It requires the `SHUTDOWN` privilege.

One use for `RESTART` is when it is not possible or convenient to gain command-line access to the MySQL server on the server host to restart it. For example, `SET PERSIST_ONLY` can be used at runtime to make configuration changes to system variables that can be set only at server startup, but the server must still be restarted for those changes to take effect. The `RESTART` statement provides a way to do so from within client sessions, without requiring command-line access on the server host.



Note

After executing a `RESTART` statement, the client can expect the current connection to be lost. If auto-reconnect is enabled, the connection is reestablished after the server restarts. Otherwise, the connection must be reestablished manually.

A successful `RESTART` operation requires `mysqld` to be running in an environment that has a monitoring process available to detect a server shutdown performed for restart purposes:

- In the presence of a monitoring process, `RESTART` causes `mysqld` to terminate such that the monitoring process can determine that it should start a new `mysqld` instance.
- If no monitoring process is present, `RESTART` fails with an error.

These platforms provide the necessary monitoring support for the `RESTART` statement:

- Windows, when `mysqld` is started as a Windows service or standalone. (`mysqld` forks, and one process acts as a monitor to the other, which acts as the server.)
- Unix and Unix-like systems that use `systemd` or `mysqld_safe` to manage `mysqld`.

To configure a monitoring environment such that `mysqld` enables the `RESTART` statement:

1. Set the `MYSQLD_PARENT_PID` environment variable to the value of the process ID of the process that starts `mysqld`, before starting `mysqld`.
2. When `mysqld` performs a shutdown due to use of the `RESTART` statement, it returns exit code 16.
3. When the monitoring process detects an exit code of 16, it starts `mysqld` again. Otherwise, it exits.

Here is a minimal example as implemented in the `bash` shell:

```
#!/bin/bash

export MYSQLD_PARENT_PID=$$

export MYSQLD_RESTART_EXIT=16

while true ; do
    bin/mysqld mysql options here
    if [ $? -ne $MYSQLD_RESTART_EXIT ] ; then
        break
    fi
done
```

On Windows, the forking used to implement `RESTART` makes determining the server process to attach to for debugging more difficult. To alleviate this, starting the server with `--gdb` suppresses forking, in addition to its other actions done to set up a debugging environment. In non-debug settings, `--no-monitor` may be used for the sole purpose of suppressing forking the monitor process. For a server started with either `--gdb` or `--no-monitor`, executing `RESTART` causes the server to simply exit without restarting.

The `Com_restart` status variable tracks the number of `RESTART` statements. Because status variables are initialized for each server startup and do not persist across restarts, `Com_restart` normally has a value of zero, but can be nonzero if `RESTART` statements were executed but failed.

13.7.8.9 SHUTDOWN Statement

```
SHUTDOWN
```

This statement stops the MySQL server. It requires the `SHUTDOWN` privilege.

`SHUTDOWN` provides an SQL-level interface to the same functionality available using the `mysqladmin shutdown` command or the `mysql_shutdown()` C API function. A successful `SHUTDOWN` sequence consists of checking the privileges, validating the arguments, and sending an OK packet to the client. Then the server is shut down.

The `Com_shutdown` status variable tracks the number of `SHUTDOWN` statements. Because status variables are initialized for each server startup and do not persist across restarts, `Com_shutdown` normally has a value of zero, but can be nonzero if `SHUTDOWN` statements were executed but failed.

Another way to stop the server is to send it a `SIGTERM` signal, which can be done by `root` or the account that owns the server process. `SIGTERM` enables server shutdown to be performed without having to connect to the server. See [Section 4.10, “Unix Signal Handling in MySQL”](#).

13.8 Utility Statements

13.8.1 DESCRIBE Statement

The `DESCRIBE` and `EXPLAIN` statements are synonyms, used either to obtain information about table structure or query execution plans. For more information, see [Section 13.7.7.5, “SHOW COLUMNS Statement”](#), and [Section 13.8.2, “EXPLAIN Statement”](#).

13.8.2 EXPLAIN Statement

```
{EXPLAIN | DESCRIBE | DESC}
  tbl_name [col_name | wild]

{EXPLAIN | DESCRIBE | DESC}
  [explain_type]
  {explainable_stmt | FOR CONNECTION connection_id}

{EXPLAIN | DESCRIBE | DESC} ANALYZE [FORMAT = TREE] select_statement

explain_type: {
  FORMAT = format_name
}

format_name: {
  TRADITIONAL
  | JSON
  | TREE
}

explainable_stmt: {
  SELECT statement
  | TABLE statement
  | DELETE statement
  | INSERT statement
  | REPLACE statement
  | UPDATE statement
}
```

The `DESCRIBE` and `EXPLAIN` statements are synonyms. In practice, the `DESCRIBE` keyword is more often used to obtain information about table structure, whereas `EXPLAIN` is used to obtain a query execution plan (that is, an explanation of how MySQL would execute a query).

The following discussion uses the `DESCRIBE` and `EXPLAIN` keywords in accordance with those uses, but the MySQL parser treats them as completely synonymous.

- [Obtaining Table Structure Information](#)

- [Obtaining Execution Plan Information](#)
- [Obtaining Information with EXPLAIN ANALYZE](#)

Obtaining Table Structure Information

`DESCRIBE` provides information about the columns in a table:

Field	Type	Null	Key	Default	Extra
<code>Id</code>	<code>int(11)</code>	<code>NO</code>	<code>PRI</code>	<code>NULL</code>	<code>auto_increment</code>
<code>Name</code>	<code>char(35)</code>	<code>NO</code>			
<code>Country</code>	<code>char(3)</code>	<code>NO</code>	<code>UNI</code>		
<code>District</code>	<code>char(20)</code>	<code>YES</code>	<code>MUL</code>		
<code>Population</code>	<code>int(11)</code>	<code>NO</code>		<code>0</code>	

`DESCRIBE` is a shortcut for `SHOW COLUMNS`. These statements also display information for views. The description for `SHOW COLUMNS` provides more information about the output columns. See [Section 13.7.7.5, “SHOW COLUMNS Statement”](#).

By default, `DESCRIBE` displays information about all columns in the table. `col_name`, if given, is the name of a column in the table. In this case, the statement displays information only for the named column. `wild`, if given, is a pattern string. It can contain the SQL `%` and `_` wildcard characters. In this case, the statement displays output only for the columns with names matching the string. There is no need to enclose the string within quotation marks unless it contains spaces or other special characters.

The `DESCRIBE` statement is provided for compatibility with Oracle.

The `SHOW CREATE TABLE`, `SHOW TABLE STATUS`, and `SHOW INDEX` statements also provide information about tables. See [Section 13.7.7, “SHOW Statements”](#).

The `explain_format` system variable, added in MySQL 8.0.32, has no effect on the output of `EXPLAIN` when used to obtain information about table columns.

Obtaining Execution Plan Information

The `EXPLAIN` statement provides information about how MySQL executes statements:

- `EXPLAIN` works with `SELECT`, `DELETE`, `INSERT`, `REPLACE`, and `UPDATE` statements. In MySQL 8.0.19 and later, it also works with `TABLE` statements.
- When `EXPLAIN` is used with an explainable statement, MySQL displays information from the optimizer about the statement execution plan. That is, MySQL explains how it would process the statement, including information about how tables are joined and in which order. For information about using `EXPLAIN` to obtain execution plan information, see [Section 8.8.2, “EXPLAIN Output Format”](#).
- When `EXPLAIN` is used with `FOR CONNECTION connection_id` rather than an explainable statement, it displays the execution plan for the statement executing in the named connection. See [Section 8.8.4, “Obtaining Execution Plan Information for a Named Connection”](#).
- For explainable statements, `EXPLAIN` produces additional execution plan information that can be displayed using `SHOW WARNINGS`. See [Section 8.8.3, “Extended EXPLAIN Output Format”](#).
- `EXPLAIN` is useful for examining queries involving partitioned tables. See [Section 24.3.5, “Obtaining Information About Partitions”](#).
- The `FORMAT` option can be used to select the output format. `TRADITIONAL` presents the output in tabular format. This is the default if no `FORMAT` option is present. `JSON` format displays the

information in JSON format. In MySQL 8.0.16 and later, `TREE` provides tree-like output with more precise descriptions of query handling than the `TRADITIONAL` format; it is the only format which shows hash join usage (see [Section 8.2.1.4, “Hash Join Optimization”](#)) and is always used for `EXPLAIN ANALYZE`.

As of MySQL 8.0.32, the default output format used by `EXPLAIN` (that is, when it has no `FORMAT` option) is determined by the value of the `explain_format` system variable. The precise effects of this variable are described later in this section.

`EXPLAIN` requires the same privileges required to execute the explained statement. Additionally, `EXPLAIN` also requires the `SHOW VIEW` privilege for any explained view. `EXPLAIN ... FOR CONNECTION` also requires the `PROCESS` privilege if the specified connection belongs to a different user.

The `explain_format` system variable introduced in MySQL 8.0.32 determines the format of the output from `EXPLAIN` when used to display a query execution plan. This variable can take any of the values used with the `FORMAT` option, with the addition of `DEFAULT` as a synonym for `TRADITIONAL`. The following example uses the `country` table from the `world` database which can be obtained from [MySQL: Other Downloads](#):

```
mysql> USE world; # Make world the current database
Database changed
```

Checking the value of `explain_format`, we see that it has the default value, and that `EXPLAIN` (with no `FORMAT` option) therefore uses the traditional tabular output:

```
mysql> SELECT @@explain_format;
+-----+
| @@explain_format |
+-----+
| TRADITIONAL |
+-----+
1 row in set (0.00 sec)

mysql> EXPLAIN SELECT Name FROM country WHERE Code Like 'A%';
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| id | select_type | table | partitions | type | possible_keys | key | key_len | ref | rows | f
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1 | SIMPLE      | country | NULL       | range | PRIMARY        | PRIMARY | 12      | NULL | 17   | f
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
1 row in set, 1 warning (0.00 sec)
```

If we set the value of `explain_format` to `TREE`, then rerun the same `EXPLAIN` statement, the output uses the tree-like format:

```
mysql> SET @@explain_format=TREE;
Query OK, 0 rows affected (0.00 sec)

mysql> SELECT @@explain_format;
+-----+
| @@explain_format |
+-----+
| TREE           |
+-----+
1 row in set (0.00 sec)

mysql> EXPLAIN SELECT Name FROM country WHERE Code LIKE 'A%';
+-----+
| EXPLAIN
+-----+
| -> Filter: (country.`Code` like 'A%') (cost=3.67 rows=17)
    | -> Index range scan on country using PRIMARY over ('A' <= Code <= 'A????????') (cost=3.67 rows=17)
+-----+
1 row in set, 1 warning (0.00 sec)
```

As stated previously, the `FORMAT` option overrides this setting. Executing the same `EXPLAIN` statement using `FORMAT=JSON` instead of `FORMAT=TREE` shows that this is the case:

EXPLAIN Statement

```
mysql> EXPLAIN FORMAT=JSON SELECT Name FROM country WHERE Code LIKE 'A%';
+-----+
| EXPLAIN |
+-----+
| {
  "query_block": {
    "select_id": 1,
    "cost_info": {
      "query_cost": "3.67"
    },
    "table": {
      "table_name": "country",
      "access_type": "range",
      "possible_keys": [
        "PRIMARY"
      ],
      "key": "PRIMARY",
      "used_key_parts": [
        "Code"
      ],
      "key_length": "12",
      "rows_examined_per_scan": 17,
      "rows_produced_per_join": 17,
      "filtered": "100.00",
      "cost_info": {
        "read_cost": "1.97",
        "eval_cost": "1.70",
        "prefix_cost": "3.67",
        "data_read_per_join": "16K"
      },
      "used_columns": [
        "Code",
        "Name"
      ],
      "attached_condition": "(`world`.`country`.`Code` like 'A%')"
    }
  }
}
+-----+
1 row in set, 1 warning (0.00 sec)
```

To return the default output of `EXPLAIN` to the tabular format, set `explain_format` to `TRADITIONAL`. Alternatively, you can set it to `DEFAULT`, which has the same effect, as shown here:

```
mysql> SET @@explain_format=DEFAULT;
Query OK, 0 rows affected (0.00 sec)

mysql> SELECT @@explain_format;
+-----+
| @@explain_format |
+-----+
| TRADITIONAL     |
+-----+
1 row in set (0.00 sec)
```

With the help of `EXPLAIN`, you can see where you should add indexes to tables so that the statement executes faster by using indexes to find rows. You can also use `EXPLAIN` to check whether the optimizer joins the tables in an optimal order. To give a hint to the optimizer to use a join order corresponding to the order in which the tables are named in a `SELECT` statement, begin the statement with `SELECT STRAIGHT_JOIN` rather than just `SELECT`. (See [Section 13.2.13, “SELECT Statement”](#).)

The optimizer trace may sometimes provide information complementary to that of `EXPLAIN`. However, the optimizer trace format and content are subject to change between versions. For details, see [MySQL Internals: Tracing the Optimizer](#).

If you have a problem with indexes not being used when you believe that they should be, run `ANALYZE TABLE` to update table statistics, such as cardinality of keys, that can affect the choices the optimizer makes. See [Section 13.7.3.1, “ANALYZE TABLE Statement”](#).

**Note**

MySQL Workbench has a Visual Explain capability that provides a visual representation of `EXPLAIN` output. See [Tutorial: Using Explain to Improve Query Performance](#).

Obtaining Information with EXPLAIN ANALYZE

MySQL 8.0.18 introduces `EXPLAIN ANALYZE`, which runs a statement and produces `EXPLAIN` output along with timing and additional, iterator-based, information about how the optimizer's expectations matched the actual execution. For each iterator, the following information is provided:

- Estimated execution cost
(Some iterators are not accounted for by the cost model, and so are not included in the estimate.)
- Estimated number of returned rows
- Time to return first row
- Time spent executing this iterator (including child iterators, but not parent iterators), in milliseconds.
(When there are multiple loops, this figure shows the average time per loop.)
- Number of rows returned by the iterator
- Number of loops

The query execution information is displayed using the `TREE` output format, in which nodes represent iterators. `EXPLAIN ANALYZE` always uses the `TREE` output format. In MySQL 8.0.21 and later, this can optionally be specified explicitly using `FORMAT=TREE`; formats other than `TREE` remain unsupported.

`EXPLAIN ANALYZE` can be used with `SELECT` statements, as well as with multi-table `UPDATE` and `DELETE` statements. Beginning with MySQL 8.0.19, it can also be used with `TABLE` statements.

Beginning with MySQL 8.0.20, you can terminate this statement using `KILL QUERY` or **CTRL-C**.

`EXPLAIN ANALYZE` cannot be used with `FOR CONNECTION`.

Example output:

```
mysql> EXPLAIN ANALYZE SELECT * FROM t1 JOIN t2 ON (t1.c1 = t2.c2)\G
***** 1. row *****
EXPLAIN: -> Inner hash join (t2.c2 = t1.c1) (cost=4.70 rows=6)
(actual time=0.032..0.035 rows=6 loops=1)
    -> Table scan on t2 (cost=0.06 rows=6)
(actual time=0.003..0.005 rows=6 loops=1)
    -> Hash
        -> Table scan on t1 (cost=0.85 rows=6)
(actual time=0.018..0.022 rows=6 loops=1)

mysql> EXPLAIN ANALYZE SELECT * FROM t3 WHERE i > 8\G
***** 1. row *****
EXPLAIN: -> Filter: (t3.i > 8) (cost=1.75 rows=5)
(actual time=0.019..0.021 rows=6 loops=1)
    -> Table scan on t3 (cost=1.75 rows=15)
(actual time=0.017..0.019 rows=15 loops=1)

mysql> EXPLAIN ANALYZE SELECT * FROM t3 WHERE pk > 17\G
***** 1. row *****
EXPLAIN: -> Filter: (t3.pk > 17) (cost=1.26 rows=5)
(actual time=0.013..0.016 rows=5 loops=1)
    -> Index range scan on t3 using PRIMARY (cost=1.26 rows=5)
(actual time=0.012..0.014 rows=5 loops=1)
```

The tables used in the example output were created by the statements shown here:

```

CREATE TABLE t1 (
    c1 INTEGER DEFAULT NULL,
    c2 INTEGER DEFAULT NULL
);

CREATE TABLE t2 (
    c1 INTEGER DEFAULT NULL,
    c2 INTEGER DEFAULT NULL
);

CREATE TABLE t3 (
    pk INTEGER NOT NULL PRIMARY KEY,
    i INTEGER DEFAULT NULL
);

```

Values shown for `actual time` in the output of this statement are expressed in milliseconds.

As of MySQL 8.0.32, the `explain_format` system variable has the following effects on `EXPLAIN ANALYZE`:

- If this value of this variable is `TRADITIONAL` or `TREE`, `EXPLAIN ANALYZE` uses the `TREE` format. This ensures that this statement continues to use the `TREE` format by default, as it did prior to the introduction of `explain_format`.
- If the value of `explain_format` is `JSON`, `EXPLAIN ANALYZE` returns an error unless `FORMAT=TREE` is specified as part of the statement. This is due to the fact that `EXPLAIN ANALYZE` supports only the `TREE` output format.

We illustrate the behavior described in the second point here, re-using the last `EXPLAIN ANALYZE` statement from the previous example:

```

mysql> SET @@explain_format=JSON;
Query OK, 0 rows affected (0.00 sec)

mysql> SELECT @@explain_format;
+-----+
| @@explain_format |
+-----+
| JSON           |
+-----+
1 row in set (0.00 sec)

mysql> EXPLAIN ANALYZE SELECT * FROM t3 WHERE pk > 17\G
ERROR 1235 (42000): This version of MySQL doesn't yet support 'EXPLAIN ANALYZE with JSON format'

mysql> EXPLAIN ANALYZE FORMAT=TRADITIONAL SELECT * FROM t3 WHERE pk > 17\G
ERROR 1235 (42000): This version of MySQL doesn't yet support 'EXPLAIN ANALYZE with TRADITIONAL format'

mysql> EXPLAIN ANALYZE FORMAT=TREE SELECT * FROM t3 WHERE pk > 17\G
***** 1. row *****
EXPLAIN: -> Filter: (t3.pk > 17) (cost=1.26 rows=5)
(actual time=0.013..0.016 rows=5 loops=1)
    -> Index range scan on t3 using PRIMARY (cost=1.26 rows=5)
(actual time=0.012..0.014 rows=5 loops=1)

```

Using `FORMAT=TRADITIONAL` or `FORMAT=JSON` with `EXPLAIN ANALYZE` always raises an error, regardless of the value of `explain_format`.

Beginning with MySQL 8.0.33, numbers in the output of `EXPLAIN ANALYZE` and `EXPLAIN FORMAT=TREE` are formatted according to the following rules:

- Numbers in the range 0.001-999999.5 are printed as decimal numbers.
Decimal numbers less than 1000 have three significant digits; the remainder have four, five, or six.
- Numbers outside the range 0.001-999999.5 are printed in engineering format. Examples of such values are `1.23e+9` and `934e-6`.

- No trailing zeros are printed. For example, we print `2.3` rather than `2.30`, and `1.2e+6` rather than `1.20e+6`.
- Numbers less than `1e-12` are printed as `0`.

13.8.3 HELP Statement

```
HELP 'search_string'
```

The `HELP` statement returns online information from the MySQL Reference Manual. Its proper operation requires that the help tables in the `mysql` database be initialized with help topic information (see [Section 5.1.17, “Server-Side Help Support”](#)).

The `HELP` statement searches the help tables for the given search string and displays the result of the search. The search string is not case-sensitive.

The search string can contain the wildcard characters `%` and `_`. These have the same meaning as for pattern-matching operations performed with the `LIKE` operator. For example, `HELP 'rep%`' returns a list of topics that begin with `rep`.

The `HELP` statement understands several types of search strings:

- At the most general level, use `contents` to retrieve a list of the top-level help categories:

```
HELP 'contents'
```

- For a list of topics in a given help category, such as `Data Types`, use the category name:

```
HELP 'data types'
```

- For help on a specific help topic, such as the `ASCII()` function or the `CREATE TABLE` statement, use the associated keyword or keywords:

```
HELP 'ascii'  
HELP 'create table'
```

In other words, the search string matches a category, many topics, or a single topic. You cannot necessarily tell in advance whether a given search string returns a list of items or the help information for a single help topic. However, you can tell what kind of response `HELP` returned by examining the number of rows and columns in the result set.

The following descriptions indicate the forms that the result set can take. Output for the example statements is shown using the familiar “tabular” or “vertical” format that you see when using the `mysql` client, but note that `mysql` itself reformats `HELP` result sets in a different way.

- Empty result set

No match could be found for the search string.

- Result set containing a single row with three columns

This means that the search string yielded a hit for the help topic. The result has three columns:

- `name`: The topic name.
- `description`: Descriptive help text for the topic.
- `example`: Usage example or examples. This column might be blank.

Example: `HELP 'replace'`

Yields:

```

name: REPLACE
description: Syntax:
REPLACE(str,from_str,to_str)

Returns the string str with all occurrences of the string from_str
replaced by the string to_str. REPLACE() performs a case-sensitive
match when searching for from_str.
example: mysql> SELECT REPLACE('www.mysql.com', 'w', 'Ww');
-> 'WwWwWw.mysql.com'

```

- Result set containing multiple rows with two columns

This means that the search string matched many help topics. The result set indicates the help topic names:

- `name`: The help topic name.
- `is_it_category`: Y if the name represents a help category, N if it does not. If it does not, the `name` value when specified as the argument to the `HELP` statement should yield a single-row result set containing a description for the named item.

Example: `HELP 'status'`

Yields:

<code>name</code>	<code>is_it_category</code>
SHOW	N
SHOW ENGINE	N
SHOW MASTER STATUS	N
SHOW PROCEDURE STATUS	N
SHOW SLAVE STATUS	N
SHOW STATUS	N
SHOW TABLE STATUS	N

- Result set containing multiple rows with three columns

This means the search string matches a category. The result set contains category entries:

- `source_category_name`: The help category name.
- `name`: The category or topic name
- `is_it_category`: Y if the name represents a help category, N if it does not. If it does not, the `name` value when specified as the argument to the `HELP` statement should yield a single-row result set containing a description for the named item.

Example: `HELP 'functions'`

Yields:

<code>source_category_name</code>	<code>name</code>	<code>is_it_category</code>
Functions	CREATE FUNCTION	N
Functions	DROP FUNCTION	N
Functions	Bit Functions	Y
Functions	Comparison operators	Y
Functions	Control flow functions	Y
Functions	Date and Time Functions	Y
Functions	Encryption Functions	Y
Functions	Information Functions	Y
Functions	Logical operators	Y
Functions	Miscellaneous Functions	Y
Functions	Numeric Functions	Y

Functions	String Functions	Y

13.8.4 USE Statement

```
USE db_name
```

The `USE` statement tells MySQL to use the named database as the default (current) database for subsequent statements. This statement requires some privilege for the database or some object within it.

The named database remains the default until the end of the session or another `USE` statement is issued:

```
USE db1;
SELECT COUNT(*) FROM mytable;    # selects from db1.mytable
USE db2;
SELECT COUNT(*) FROM mytable;    # selects from db2.mytable
```

The database name must be specified on a single line. Newlines in database names are not supported.

Making a particular database the default by means of the `USE` statement does not preclude accessing tables in other databases. The following example accesses the `author` table from the `db1` database and the `editor` table from the `db2` database:

```
USE db1;
SELECT author_name,editor_name FROM author,db2.editor
 WHERE author.editor_id = db2.editor.editor_id;
```

Chapter 14 MySQL Data Dictionary

Table of Contents

14.1 Data Dictionary Schema	3023
14.2 Removal of File-based Metadata Storage	3024
14.3 Transactional Storage of Dictionary Data	3025
14.4 Dictionary Object Cache	3025
14.5 INFORMATION_SCHEMA and Data Dictionary Integration	3026
14.6 Serialized Dictionary Information (SDI)	3028
14.7 Data Dictionary Usage Differences	3028
14.8 Data Dictionary Limitations	3030

MySQL Server incorporates a transactional data dictionary that stores information about database objects. In previous MySQL releases, dictionary data was stored in metadata files, nontransactional tables, and storage engine-specific data dictionaries.

This chapter describes the main features, benefits, usage differences, and limitations of the data dictionary. For other implications of the data dictionary feature, refer to the “Data Dictionary Notes” section in the [MySQL 8.0 Release Notes](#).

Benefits of the MySQL data dictionary include:

- Simplicity of a centralized data dictionary schema that uniformly stores dictionary data. See [Section 14.1, “Data Dictionary Schema”](#).
- Removal of file-based metadata storage. See [Section 14.2, “Removal of File-based Metadata Storage”](#).
- Transactional, crash-safe storage of dictionary data. See [Section 14.3, “Transactional Storage of Dictionary Data”](#).
- Uniform and centralized caching for dictionary objects. See [Section 14.4, “Dictionary Object Cache”](#).
- A simpler and improved implementation for some `INFORMATION_SCHEMA` tables. See [Section 14.5, “INFORMATION_SCHEMA and Data Dictionary Integration”](#).
- Atomic DDL. See [Section 13.1.1, “Atomic Data Definition Statement Support”](#).



Important

A data dictionary-enabled server entails some general operational differences compared to a server that does not have a data dictionary; see [Section 14.7, “Data Dictionary Usage Differences”](#). Also, for upgrades to MySQL 8.0, the upgrade procedure differs somewhat from previous MySQL releases and requires that you verify the upgrade readiness of your installation by checking specific prerequisites. For more information, see [Section 2.10, “Upgrading MySQL”](#), particularly [Section 2.10.5, “Preparing Your Installation for Upgrade”](#).

14.1 Data Dictionary Schema

Data dictionary tables are protected and may only be accessed in debug builds of MySQL. However, MySQL supports access to data stored in data dictionary tables through `INFORMATION_SCHEMA` tables and `SHOW` statements. For an overview of the tables that comprise the data dictionary, see [Data Dictionary Tables](#).

MySQL system tables still exist in MySQL 8.0 and can be viewed by issuing a `SHOW TABLES` statement on the `mysql` system database. Generally, the difference between MySQL data dictionary

tables and system tables is that data dictionary tables contain metadata required to execute SQL queries, whereas system tables contain auxiliary data such as time zone and help information. MySQL system tables and data dictionary tables also differ in how they are upgraded. The MySQL server manages data dictionary upgrades. See [How the Data Dictionary is Upgraded](#). Upgrading MySQL system tables requires running the full MySQL upgrade procedure. See [Section 2.10.3, “What the MySQL Upgrade Process Upgrades”](#).

How the Data Dictionary is Upgraded

New versions of MySQL may include changes to data dictionary table definitions. Such changes are present in newly installed versions of MySQL, but when performing an in-place upgrade of MySQL binaries, changes are applied when the MySQL server is restarted using the new binaries. At startup, the data dictionary version of the server is compared to the version information stored in the data dictionary to determine if data dictionary tables should be upgraded. If an upgrade is necessary and supported, the server creates data dictionary tables with updated definitions, copies persisted metadata to the new tables, atomically replaces the old tables with the new ones, and reinitializes the data dictionary. If an upgrade is not necessary, startup continues without updating the data dictionary tables.

Upgrade of data dictionary tables is an atomic operation, which means that all of the data dictionary tables are upgraded as necessary or the operation fails. If the upgrade operation fails, server startup fails with an error. In this case, the old server binaries can be used with the old data directory to start the server. When the new server binaries are used again to start the server, the data dictionary upgrade is reattempted.

Generally, after data dictionary tables are successfully upgraded, it is not possible to restart the server using the old server binaries. As a result, downgrading MySQL server binaries to a previous MySQL version is not supported after data dictionary tables are upgraded.

The `mysqld --no-dd-upgrade` option can be used to prevent automatic upgrade of data dictionary tables at startup. When `--no-dd-upgrade` is specified, and the server finds that the data dictionary version of the server is different from the version stored in the data dictionary, startup fails with an error stating that the data dictionary upgrade is prohibited.

Viewing Data Dictionary Tables Using a Debug Build of MySQL

Data dictionary tables are protected by default but can be accessed by compiling MySQL with debugging support (using the `-DWITH_DEBUG=1 CMake` option) and specifying the `+d,skip_dd_table_access_check` debug option and modifier. For information about compiling debug builds, see [Section 5.9.1.1, “Compiling MySQL for Debugging”](#).



Warning

Modifying or writing to data dictionary tables directly is not recommended and may render your MySQL instance inoperable.

After compiling MySQL with debugging support, use this `SET` statement to make data dictionary tables visible to the `mysql` client session:

```
mysql> SET SESSION debug='+d,skip_dd_table_access_check';
```

Use this query to retrieve a list of data dictionary tables:

```
mysql> SELECT name, schema_id, hidden, type FROM mysql.tables WHERE schema_id=1 AND hidden='System';
```

Use `SHOW CREATE TABLE` to view data dictionary table definitions. For example:

```
mysql> SHOW CREATE TABLE mysql.catalogs\G
```

14.2 Removal of File-based Metadata Storage

In previous MySQL releases, dictionary data was partially stored in metadata files. Issues with file-based metadata storage included expensive file scans, susceptibility to file system-related bugs,

complex code for handling of replication and crash recovery failure states, and a lack of extensibility that made it difficult to add metadata for new features and relational objects.

The metadata files listed below are removed from MySQL. Unless otherwise noted, data previously stored in metadata files is now stored in data dictionary tables.

- `.frm` files: Table metadata files. With the removal of `.frm` files:
 - The 64KB table definition size limit imposed by the `.frm` file structure is removed.
 - The Information Schema `TABLES` table's `VERSION` column reports a hardcoded value of `10`, which is the last `.frm` file version used in MySQL 5.7.
- `.par` files: Partition definition files. `InnoDB` stopped using partition definition files in MySQL 5.7 with the introduction of native partitioning support for `InnoDB` tables.
- `.TRN` files: Trigger namespace files.
- `.TRG` files: Trigger parameter files.
- `.isl` files: `InnoDB` Symbolic Link files containing the location of `file-per-table` tablespace files created outside of the data directory.
- `db.opt` files: Database configuration files. These files, one per database directory, contained database default character set attributes.
- `ddl_log.log` file: The file contained records of metadata operations generated by data definition statements such as `DROP TABLE` and `ALTER TABLE`.

14.3 Transactional Storage of Dictionary Data

The data dictionary schema stores dictionary data in transactional (`InnoDB`) tables. Data dictionary tables are located in the `mysql` database together with non-data dictionary system tables.

Data dictionary tables are created in a single `InnoDB` tablespace named `mysql.ibd`, which resides in the MySQL data directory. The `mysql.ibd` tablespace file must reside in the MySQL data directory and its name cannot be modified or used by another tablespace.

Dictionary data is protected by the same commit, rollback, and crash-recovery capabilities that protect user data that is stored in `InnoDB` tables.

14.4 Dictionary Object Cache

The dictionary object cache is a shared global cache that stores previously accessed data dictionary objects in memory to enable object reuse and minimize disk I/O. Similar to other cache mechanisms used by MySQL, the dictionary object cache uses an `LRU`-based eviction strategy to evict least recently used objects from memory.

The dictionary object cache comprises cache partitions that store different object types. Some cache partition size limits are configurable, whereas others are hardcoded.

- **tablespace definition cache partition:** Stores tablespace definition objects. The `tablespace_definition_cache` option sets a limit for the number of tablespace definition objects that can be stored in the dictionary object cache. The default value is 256.
- **schema definition cache partition:** Stores schema definition objects. The `schema_definition_cache` option sets a limit for the number of schema definition objects that can be stored in the dictionary object cache. The default value is 256.
- **table definition cache partition:** Stores table definition objects. The object limit is set to the value of `max_connections`, which has a default value of 151.

The table definition cache partition exists in parallel with the table definition cache that is configured using the `table_definition_cache` configuration option. Both caches store table definitions but serve different parts of the MySQL server. Objects in one cache have no dependence on the existence of objects in the other.

- **stored program definition cache partition:** Stores stored program definition objects. The `stored_program_definition_cache` option sets a limit for the number of stored program definition objects that can be stored in the dictionary object cache. The default value is 256.

The stored program definition cache partition exists in parallel with the stored procedure and stored function caches that are configured using the `stored_program_cache` option.

The `stored_program_cache` option sets a soft upper limit for the number of cached stored procedures or functions per connection, and the limit is checked each time a connection executes a stored procedure or function. The stored program definition cache partition, on the other hand, is a shared cache that stores stored program definition objects for other purposes. The existence of objects in the stored program definition cache partition has no dependence on the existence of objects in the stored procedure cache or stored function cache, and vice versa.

- **character set definition cache partition:** Stores character set definition objects and has a hardcoded object limit of 256.
- **collation definition cache partition:** Stores collation definition objects and has a hardcoded object limit of 256.

For information about valid values for dictionary object cache configuration options, refer to [Section 5.1.8, “Server System Variables”](#).

14.5 INFORMATION_SCHEMA and Data Dictionary Integration

With the introduction of the data dictionary, the following `INFORMATION_SCHEMA` tables are implemented as views on data dictionary tables:

- `CHARACTER_SETS`
- `CHECK_CONSTRAINTS`
- `COLLATIONS`
- `COLLATION_CHARACTER_SET_APPLICABILITY`
- `COLUMNS`
- `COLUMN_STATISTICS`
- `EVENTS`
- `FILES`
- `INNODB_COLUMNS`
- `INNODB_DATAFILES`
- `INNODB_FIELDS`
- `INNODB_FOREIGN`
- `INNODB_FOREIGN_COLS`
- `INNODB_INDEXES`
- `INNODB_TABLES`
- `INNODB_TABLESPACES`

- [INNODB_TABLESPACES_BRIEF](#)
- [INNODB_TABLESTATS](#)
- [KEY_COLUMN_USAGE](#)
- [KEYWORDS](#)
- [PARAMETERS](#)
- [PARTITIONS](#)
- [REFERENTIAL_CONSTRAINTS](#)
- [RESOURCE_GROUPS](#)
- [ROUTINES](#)
- [SCHEMATA](#)
- [STATISTICS](#)
- [ST_GEOMETRY_COLUMNS](#)
- [ST_SPATIAL_REFERENCE_SYSTEMS](#)
- [TABLES](#)
- [TABLE_CONSTRAINTS](#)
- [TRIGGERS](#)
- [VIEWS](#)
- [VIEW_ROUTINE_USAGE](#)
- [VIEW_TABLE_USAGE](#)

Queries on those tables are now more efficient because they obtain information from data dictionary tables rather than by other, slower means. In particular, for each [INFORMATION_SCHEMA](#) table that is a view on data dictionary tables:

- The server no longer must create a temporary table for each query of the [INFORMATION_SCHEMA](#) table.
- When the underlying data dictionary tables store values previously obtained by directory scans (for example, to enumerate database names or table names within databases) or file-opening operations (for example, to read information from `.frm` files), [INFORMATION_SCHEMA](#) queries for those values now use table lookups instead. (Additionally, even for a non-view [INFORMATION_SCHEMA](#) table, values such as database and table names are retrieved by lookups from the data dictionary and do not require directory or file scans.)
- Indexes on the underlying data dictionary tables permit the optimizer to construct efficient query execution plans, something not true for the previous implementation that processed the [INFORMATION_SCHEMA](#) table using a temporary table per query.

The preceding improvements also apply to [SHOW](#) statements that display information corresponding to the [INFORMATION_SCHEMA](#) tables that are views on data dictionary tables. For example, [SHOW DATABASES](#) displays the same information as the [SCHEMATA](#) table.

In addition to the introduction of views on data dictionary tables, table statistics contained in the [STATISTICS](#) and [TABLES](#) tables is now cached to improve [INFORMATION_SCHEMA](#) query performance. The `information_schema_stats_expiry` system variable defines the period of time before cached table statistics expire. The default is 86400 seconds (24 hours). If there are no cached statistics or statistics have expired, statistics are retrieved from storage engine when querying table statistics columns. To update cached values at any time for a given table, use [ANALYZE TABLE](#)

`information_schema_stats_expiry` can be set to 0 to have `INFORMATION_SCHEMA` queries retrieve the latest statistics directly from the storage engine, which is not as fast as retrieving cached statistics.

For more information, see [Section 8.2.3, “Optimizing INFORMATION_SCHEMA Queries”](#).

`INFORMATION_SCHEMA` tables in MySQL 8.0 are closely tied to the data dictionary, resulting in several usage differences. See [Section 14.7, “Data Dictionary Usage Differences”](#).

14.6 Serialized Dictionary Information (SDI)

In addition to storing metadata about database objects in the data dictionary, MySQL stores it in serialized form. This data is referred to as serialized dictionary information (SDI). `InnoDB` stores SDI data within its tablespace files. `NDBCLUSTER` stores SDI data in the NDB dictionary. Other storage engines store SDI data in `.sdi` files that are created for a given table in the table's database directory. SDI data is generated in a compact `JSON` format.

Serialized dictionary information (SDI) is present in all `InnoDB` tablespace files except for temporary tablespace and undo tablespace files. SDI records in an `InnoDB` tablespace file only describe table and tablespace objects contained within the tablespace.

SDI data is updated by DDL operations on a table or `CHECK TABLE FOR UPGRADE`. SDI data is not updated when the MySQL server is upgraded to a new release or version.

The presence of SDI data provides metadata redundancy. For example, if the data dictionary becomes unavailable, object metadata can be extracted directly from `InnoDB` tablespace files using the `ibd2sdi` tool.

For `InnoDB`, an SDI record requires a single index page, which is 16KB in size by default. However, SDI data is compressed to reduce the storage footprint.

For partitioned `InnoDB` tables comprised of multiple tablespaces, SDI data is stored in the tablespace file of the first partition.

The MySQL server uses an internal API that is accessed during `DDL` operations to create and maintain SDI records.

The `IMPORT TABLE` statement imports `MyISAM` tables based on information contained in `.sdi` files. For more information, see [Section 13.2.6, “IMPORT TABLE Statement”](#).

14.7 Data Dictionary Usage Differences

Use of a data dictionary-enabled MySQL server entails some operational differences compared to a server that does not have a data dictionary:

- Previously, enabling the `innodb_read_only` system variable prevented creating and dropping tables only for the `InnoDB` storage engine. As of MySQL 8.0, enabling `innodb_read_only` prevents these operations for all storage engines. Table creation and drop operations for any storage engine modify data dictionary tables in the `mysql` system database, but those tables use the `InnoDB` storage engine and cannot be modified when `innodb_read_only` is enabled. The same principle applies to other table operations that require modifying data dictionary tables. Examples:
 - `ANALYZE TABLE` fails because it updates table statistics, which are stored in the data dictionary.
 - `ALTER TABLE tbl_name ENGINE=engine_name` fails because it updates the storage engine designation, which is stored in the data dictionary.



Note

Enabling `innodb_read_only` also has important implications for non-data dictionary tables in the `mysql` system database. For details, see the

description of `innodb_read_only` in [Section 15.14, “InnoDB Startup Options and System Variables”](#)

- Previously, tables in the `mysql` system database were visible to DML and DDL statements. As of MySQL 8.0, data dictionary tables are invisible and cannot be modified or queried directly. However, in most cases there are corresponding `INFORMATION_SCHEMA` tables that can be queried instead. This enables the underlying data dictionary tables to be changed as server development proceeds, while maintaining a stable `INFORMATION_SCHEMA` interface for application use.
- `INFORMATION_SCHEMA` tables in MySQL 8.0 are closely tied to the data dictionary, resulting in several usage differences:
 - Previously, `INFORMATION_SCHEMA` queries for table statistics in the `STATISTICS` and `TABLES` tables retrieved statistics directly from storage engines. As of MySQL 8.0, cached table statistics are used by default. The `information_schema_stats_expiry` system variable defines the period of time before cached table statistics expire. The default is 86400 seconds (24 hours). (To update the cached values at any time for a given table, use `ANALYZE TABLE`.) If there are no cached statistics or statistics have expired, statistics are retrieved from storage engines when querying table statistics columns. To always retrieve the latest statistics directly from storage engines, set `information_schema_stats_expiry` to 0. For more information, see [Section 8.2.3, “Optimizing INFORMATION_SCHEMA Queries”](#).
 - Several `INFORMATION_SCHEMA` tables are views on data dictionary tables, which enables the optimizer to use indexes on those underlying tables. Consequently, depending on optimizer choices, the row order of results for `INFORMATION_SCHEMA` queries might differ from previous results. If a query result must have specific row ordering characteristics, include an `ORDER BY` clause.
 - Queries on `INFORMATION_SCHEMA` tables may return column names in a different lettercase than in earlier MySQL series. Applications should test result set column names in case-insensitive fashion. If that is not feasible, a workaround is to use column aliases in the select list that return column names in the required lettercase. For example:

```
SELECT TABLE_SCHEMA AS table_schema, TABLE_NAME AS table_name
FROM INFORMATION_SCHEMA.TABLES WHERE TABLE_NAME = 'users';
```

- `mysqldump` and `mysqlpump` no longer dump the `INFORMATION_SCHEMA` database, even if explicitly named on the command line.
- `CREATE TABLE dst_tbl LIKE src_tbl` requires that `src_tbl` be a base table and fails if it is an `INFORMATION_SCHEMA` table that is a view on data dictionary tables.
- Previously, result set headers of columns selected from `INFORMATION_SCHEMA` tables used the capitalization specified in the query. This query produces a result set with a header of `table_name`:

```
SELECT table_name FROM INFORMATION_SCHEMA.TABLES;
```

As of MySQL 8.0, these headers are capitalized; the preceding query produces a result set with a header of `TABLE_NAME`. If necessary, a column alias can be used to achieve a different lettercase. For example:

```
SELECT table_name AS 'table_name' FROM INFORMATION_SCHEMA.TABLES;
```

- The data directory affects how `mysqldump` and `mysqlpump` dump information from the `mysql` system database:
 - Previously, it was possible to dump all tables in the `mysql` system database. As of MySQL 8.0, `mysqldump` and `mysqlpump` dump only non-data dictionary tables in that database.
 - Previously, the `--routines` and `--events` options were not required to include stored routines and events when using the `--all-databases` option: The dump included the `mysql` system database, and therefore also the `proc` and `event` tables containing stored routine and event definitions. As of MySQL 8.0, the `event` and `proc` tables are not used. Definitions for the corresponding objects are stored in data dictionary tables, but those tables are not dumped. To include stored routines and events in a dump made using `--all-databases`, use the `--routines` and `--events` options explicitly.
 - Previously, the `--routines` option required the `SELECT` privilege for the `proc` table. As of MySQL 8.0, that table is not used; `--routines` requires the global `SELECT` privilege instead.
 - Previously, it was possible to dump stored routine and event definitions together with their creation and modification timestamps, by dumping the `proc` and `event` tables. As of MySQL 8.0, those tables are not used, so it is not possible to dump timestamps.
 - Previously, creating a stored routine that contains illegal characters produced a warning. As of MySQL 8.0, this is an error.

14.8 Data Dictionary Limitations

This section describes temporary limitations introduced with the MySQL data dictionary.

- Manual creation of database directories under the data directory (for example, with `mkdir`) is unsupported. Manually created database directories are not recognized by the MySQL Server.
- DDL operations take longer due to writing to storage, undo logs, and redo logs instead of `.frm` files.

Chapter 15 The InnoDB Storage Engine

Table of Contents

15.1	Introduction to InnoDB	3032
15.1.1	Benefits of Using InnoDB Tables	3034
15.1.2	Best Practices for InnoDB Tables	3035
15.1.3	Verifying that InnoDB is the Default Storage Engine	3035
15.1.4	Testing and Benchmarking with InnoDB	3036
15.2	InnoDB and the ACID Model	3036
15.3	InnoDB Multi-Versioning	3037
15.4	InnoDB Architecture	3039
15.5	InnoDB In-Memory Structures	3039
15.5.1	Buffer Pool	3040
15.5.2	Change Buffer	3045
15.5.3	Adaptive Hash Index	3049
15.5.4	Log Buffer	3050
15.6	InnoDB On-Disk Structures	3050
15.6.1	Tables	3050
15.6.2	Indexes	3074
15.6.3	Tablespaces	3081
15.6.4	Doublewrite Buffer	3104
15.6.5	Redo Log	3105
15.6.6	Undo Logs	3112
15.7	InnoDB Locking and Transaction Model	3113
15.7.1	InnoDB Locking	3114
15.7.2	InnoDB Transaction Model	3118
15.7.3	Locks Set by Different SQL Statements in InnoDB	3127
15.7.4	Phantom Rows	3130
15.7.5	Deadlocks in InnoDB	3131
15.7.6	Transaction Scheduling	3136
15.8	InnoDB Configuration	3137
15.8.1	InnoDB Startup Configuration	3137
15.8.2	Configuring InnoDB for Read-Only Operation	3143
15.8.3	InnoDB Buffer Pool Configuration	3145
15.8.4	Configuring Thread Concurrency for InnoDB	3159
15.8.5	Configuring the Number of Background InnoDB I/O Threads	3160
15.8.6	Using Asynchronous I/O on Linux	3161
15.8.7	Configuring InnoDB I/O Capacity	3161
15.8.8	Configuring Spin Lock Polling	3163
15.8.9	Purge Configuration	3164
15.8.10	Configuring Optimizer Statistics for InnoDB	3165
15.8.11	Configuring the Merge Threshold for Index Pages	3176
15.8.12	Enabling Automatic Configuration for a Dedicated MySQL Server	3178
15.9	InnoDB Table and Page Compression	3181
15.9.1	InnoDB Table Compression	3181
15.9.2	InnoDB Page Compression	3195
15.10	InnoDB Row Formats	3198
15.11	InnoDB Disk I/O and File Space Management	3204
15.11.1	InnoDB Disk I/O	3205
15.11.2	File Space Management	3205
15.11.3	InnoDB Checkpoints	3207
15.11.4	Defragmenting a Table	3207
15.11.5	Reclaiming Disk Space with TRUNCATE TABLE	3208
15.12	InnoDB and Online DDL	3208
15.12.1	Online DDL Operations	3209

15.12.2 Online DDL Performance and Concurrency	3224
15.12.3 Online DDL Space Requirements	3227
15.12.4 Online DDL Memory Management	3228
15.12.5 Configuring Parallel Threads for Online DDL Operations	3228
15.12.6 Simplifying DDL Statements with Online DDL	3229
15.12.7 Online DDL Failure Conditions	3229
15.12.8 Online DDL Limitations	3230
15.13 InnoDB Data-at-Rest Encryption	3230
15.14 InnoDB Startup Options and System Variables	3239
15.15 InnoDB INFORMATION_SCHEMA Tables	3329
15.15.1 InnoDB INFORMATION_SCHEMA Tables about Compression	3330
15.15.2 InnoDB INFORMATION_SCHEMA Transaction and Locking Information	3331
15.15.3 InnoDB INFORMATION_SCHEMA Schema Object Tables	3338
15.15.4 InnoDB INFORMATION_SCHEMA FULLTEXT Index Tables	3343
15.15.5 InnoDB INFORMATION_SCHEMA Buffer Pool Tables	3346
15.15.6 InnoDB INFORMATION_SCHEMA Metrics Table	3350
15.15.7 InnoDB INFORMATION_SCHEMA Temporary Table Info Table	3360
15.15.8 Retrieving InnoDB Tablespace Metadata from INFORMATION_SCHEMA.FILES	3360
15.16 InnoDB Integration with MySQL Performance Schema	3362
15.16.1 Monitoring ALTER TABLE Progress for InnoDB Tables Using Performance Schema	3363
15.16.2 Monitoring InnoDB Mutex Waits Using Performance Schema	3365
15.17 InnoDB Monitors	3369
15.17.1 InnoDB Monitor Types	3369
15.17.2 Enabling InnoDB Monitors	3369
15.17.3 InnoDB Standard Monitor and Lock Monitor Output	3371
15.18 InnoDB Backup and Recovery	3375
15.18.1 InnoDB Backup	3375
15.18.2 InnoDB Recovery	3376
15.19 InnoDB and MySQL Replication	3378
15.20 InnoDB memcached Plugin	3380
15.20.1 Benefits of the InnoDB memcached Plugin	3381
15.20.2 InnoDB memcached Architecture	3382
15.20.3 Setting Up the InnoDB memcached Plugin	3383
15.20.4 InnoDB memcached Multiple get and Range Query Support	3388
15.20.5 Security Considerations for the InnoDB memcached Plugin	3391
15.20.6 Writing Applications for the InnoDB memcached Plugin	3392
15.20.7 The InnoDB memcached Plugin and Replication	3404
15.20.8 InnoDB memcached Plugin Internals	3408
15.20.9 Troubleshooting the InnoDB memcached Plugin	3412
15.21 InnoDB Troubleshooting	3414
15.21.1 Troubleshooting InnoDB I/O Problems	3415
15.21.2 Troubleshooting Recovery Failures	3415
15.21.3 Forcing InnoDB Recovery	3416
15.21.4 Troubleshooting InnoDB Data Dictionary Operations	3417
15.21.5 InnoDB Error Handling	3418
15.22 InnoDB Limits	3419
15.23 InnoDB Restrictions and Limitations	3420

15.1 Introduction to InnoDB

InnoDB is a general-purpose storage engine that balances high reliability and high performance. In MySQL 8.0, InnoDB is the default MySQL storage engine. Unless you have configured a different default storage engine, issuing a `CREATE TABLE` statement without an `ENGINE` clause creates an InnoDB table.

Key Advantages of InnoDB

- Its DML operations follow the ACID model, with transactions featuring commit, rollback, and crash-recovery capabilities to protect user data. See [Section 15.2, “InnoDB and the ACID Model”](#).
- Row-level locking and Oracle-style consistent reads increase multi-user concurrency and performance. See [Section 15.7, “InnoDB Locking and Transaction Model”](#).
- [InnoDB](#) tables arrange your data on disk to optimize queries based on primary keys. Each [InnoDB](#) table has a primary key index called the clustered index that organizes the data to minimize I/O for primary key lookups. See [Section 15.6.2.1, “Clustered and Secondary Indexes”](#).
- To maintain data integrity, [InnoDB](#) supports `FOREIGN KEY` constraints. With foreign keys, inserts, updates, and deletes are checked to ensure they do not result in inconsistencies across related tables. See [Section 13.1.20.5, “FOREIGN KEY Constraints”](#).

Table 15.1 InnoDB Storage Engine Features

Feature	Support
B-tree indexes	Yes
Backup/point-in-time recovery (Implemented in the server, rather than in the storage engine.)	Yes
Cluster database support	No
Clustered indexes	Yes
Compressed data	Yes
Data caches	Yes
Encrypted data	Yes (Implemented in the server via encryption functions; In MySQL 5.7 and later, data-at-rest encryption is supported.)
Foreign key support	Yes
Full-text search indexes	Yes (Support for FULLTEXT indexes is available in MySQL 5.6 and later.)
Geospatial data type support	Yes
Geospatial indexing support	Yes (Support for geospatial indexing is available in MySQL 5.7 and later.)
Hash indexes	No (InnoDB utilizes hash indexes internally for its Adaptive Hash Index feature.)
Index caches	Yes
Locking granularity	Row
MVCC	Yes
Replication support (Implemented in the server, rather than in the storage engine.)	Yes
Storage limits	64TB
T-tree indexes	No
Transactions	Yes
Update statistics for data dictionary	Yes

To compare the features of [InnoDB](#) with other storage engines provided with MySQL, see the *Storage Engine Features* table in [Chapter 16, Alternative Storage Engines](#).

InnoDB Enhancements and New Features

For information about InnoDB enhancements and new features, refer to:

- The InnoDB enhancements list in [Section 1.3, “What Is New in MySQL 8.0”](#).
- The [Release Notes](#).

Additional InnoDB Information and Resources

- For InnoDB-related terms and definitions, see the [MySQL Glossary](#).
- For a forum dedicated to the InnoDB storage engine, see [MySQL Forums::InnoDB](#).
- InnoDB is published under the same GNU GPL License Version 2 (of June 1991) as MySQL. For more information on MySQL licensing, see <http://www.mysql.com/company/legal/licensing/>.

15.1.1 Benefits of Using InnoDB Tables

InnoDB tables have the following benefits:

- If the server unexpectedly exits because of a hardware or software issue, regardless of what was happening in the database at the time, you don't need to do anything special after restarting the database. InnoDB crash recovery automatically finalizes changes that were committed before the time of the crash, and undoes changes that were in process but not committed, permitting you to restart and continue from where you left off. See [Section 15.18.2, “InnoDB Recovery”](#).
- The InnoDB storage engine maintains its own buffer pool that caches table and index data in main memory as data is accessed. Frequently used data is processed directly from memory. This cache applies to many types of information and speeds up processing. On dedicated database servers, up to 80% of physical memory is often assigned to the buffer pool. See [Section 15.5.1, “Buffer Pool”](#).
- If you split up related data into different tables, you can set up foreign keys that enforce referential integrity. See [Section 13.1.20.5, “FOREIGN KEY Constraints”](#).
- If data becomes corrupted on disk or in memory, a checksum mechanism alerts you to the bogus data before you use it. The `innodb_checksum_algorithm` variable defines the checksum algorithm used by InnoDB.
- When you design a database with appropriate primary key columns for each table, operations involving those columns are automatically optimized. It is very fast to reference the primary key columns in `WHERE` clauses, `ORDER BY` clauses, `GROUP BY` clauses, and join operations. See [Section 15.6.2.1, “Clustered and Secondary Indexes”](#).
- Inserts, updates, and deletes are optimized by an automatic mechanism called change buffering. InnoDB not only allows concurrent read and write access to the same table, it caches changed data to streamline disk I/O. See [Section 15.5.2, “Change Buffer”](#).
- Performance benefits are not limited to large tables with long-running queries. When the same rows are accessed over and over from a table, the Adaptive Hash Index takes over to make these lookups even faster, as if they came out of a hash table. See [Section 15.5.3, “Adaptive Hash Index”](#).
- You can compress tables and associated indexes. See [Section 15.9, “InnoDB Table and Page Compression”](#).
- You can encrypt your data. See [Section 15.13, “InnoDB Data-at-Rest Encryption”](#).
- You can create and drop indexes and perform other DDL operations with much less impact on performance and availability. See [Section 15.12.1, “Online DDL Operations”](#).
- Truncating a file-per-table tablespace is very fast and can free up disk space for the operating system to reuse rather than only InnoDB. See [Section 15.6.3.2, “File-Per-Table Tablespace”](#).

- The storage layout for table data is more efficient for `BLOB` and long text fields, with the `DYNAMIC` row format. See [Section 15.10, “InnoDB Row Formats”](#).
- You can monitor the internal workings of the storage engine by querying `INFORMATION_SCHEMA` tables. See [Section 15.15, “InnoDB INFORMATION_SCHEMA Tables”](#).
- You can monitor the performance details of the storage engine by querying Performance Schema tables. See [Section 15.16, “InnoDB Integration with MySQL Performance Schema”](#).
- You can mix `InnoDB` tables with tables from other MySQL storage engines, even within the same statement. For example, you can use a join operation to combine data from `InnoDB` and `MEMORY` tables in a single query.
- `InnoDB` has been designed for CPU efficiency and maximum performance when processing large data volumes.
- `InnoDB` tables can handle large quantities of data, even on operating systems where file size is limited to 2GB.

For `InnoDB`-specific tuning techniques you can apply to your MySQL server and application code, see [Section 8.5, “Optimizing for InnoDB Tables”](#).

15.1.2 Best Practices for InnoDB Tables

This section describes best practices when using `InnoDB` tables.

- Specify a primary key for every table using the most frequently queried column or columns, or an auto-increment value if there is no obvious primary key.
- Use joins wherever data is pulled from multiple tables based on identical ID values from those tables. For fast join performance, define foreign keys on the join columns, and declare those columns with the same data type in each table. Adding foreign keys ensures that referenced columns are indexed, which can improve performance. Foreign keys also propagate deletes and updates to all affected tables, and prevent insertion of data in a child table if the corresponding IDs are not present in the parent table.
- Turn off autocommit. Committing hundreds of times a second puts a cap on performance (limited by the write speed of your storage device).
- Group sets of related DML operations into transactions by bracketing them with `START TRANSACTION` and `COMMIT` statements. While you don't want to commit too often, you also don't want to issue huge batches of `INSERT`, `UPDATE`, or `DELETE` statements that run for hours without committing.
- Do not use `LOCK TABLES` statements. `InnoDB` can handle multiple sessions all reading and writing to the same table at once without sacrificing reliability or high performance. To get exclusive write access to a set of rows, use the `SELECT ... FOR UPDATE` syntax to lock just the rows you intend to update.
- Enable the `innodb_file_per_table` variable or use general tablespaces to put the data and indexes for tables into separate files instead of the system tablespace. The `innodb_file_per_table` variable is enabled by default.
- Evaluate whether your data and access patterns benefit from the `InnoDB` table or page compression features. You can compress `InnoDB` tables without sacrificing read/write capability.
- Run the server with the `--sql_mode=NO_ENGINE_SUBSTITUTION` option to prevent tables from being created with storage engines that you do not want to use.

15.1.3 Verifying that InnoDB is the Default Storage Engine

Issue the `SHOW ENGINES` statement to view the available MySQL storage engines. Look for `DEFAULT` in the `SUPPORT` column.

```
mysql> SHOW ENGINES;
```

Alternatively, query the Information Schema `ENGINES` table.

```
mysql> SELECT * FROM INFORMATION_SCHEMA.ENGINES;
```

15.1.4 Testing and Benchmarking with InnoDB

If `InnoDB` is not the default storage engine, you can determine if your database server and applications work correctly with `InnoDB` by restarting the server with `--default-storage-engine=InnoDB` defined on the command line or with `default-storage-engine=innodb` defined in the `[mysqld]` section of the MySQL server option file.

Since changing the default storage engine only affects newly created tables, run your application installation and setup steps to confirm that everything installs properly, then exercise the application features to make sure the data loading, editing, and querying features work. If a table relies on a feature that is specific to another storage engine, you receive an error. In this case, add the `ENGINE=other_engine_name` clause to the `CREATE TABLE` statement to avoid the error.

If you did not make a deliberate decision about the storage engine, and you want to preview how certain tables work when created using `InnoDB`, issue the command `ALTER TABLE table_name ENGINE=InnoDB;` for each table. Alternatively, to run test queries and other statements without disturbing the original table, make a copy:

```
CREATE TABLE ... ENGINE=InnoDB AS SELECT * FROM other_engine_table;
```

To assess performance with a full application under a realistic workload, install the latest MySQL server and run benchmarks.

Test the full application lifecycle, from installation, through heavy usage, and server restart. Kill the server process while the database is busy to simulate a power failure, and verify that the data is recovered successfully when you restart the server.

Test any replication configurations, especially if you use different MySQL versions and options on the source server and replicas.

15.2 InnoDB and the ACID Model

The `ACID` model is a set of database design principles that emphasize aspects of reliability that are important for business data and mission-critical applications. MySQL includes components such as the `InnoDB` storage engine that adhere closely to the ACID model so that data is not corrupted and results are not distorted by exceptional conditions such as software crashes and hardware malfunctions. When you rely on ACID-compliant features, you do not need to reinvent the wheel of consistency checking and crash recovery mechanisms. In cases where you have additional software safeguards, ultra-reliable hardware, or an application that can tolerate a small amount of data loss or inconsistency, you can adjust MySQL settings to trade some of the ACID reliability for greater performance or throughput.

The following sections discuss how MySQL features, in particular the `InnoDB` storage engine, interact with the categories of the ACID model:

- **A:** atomicity.
- **C:** consistency.
- **I:** isolation.
- **D:** durability.

Atomicity

The **atomicity** aspect of the ACID model mainly involves `InnoDB transactions`. Related MySQL features include:

- The `autocommit` setting.
- The `COMMIT` statement.
- The `ROLLBACK` statement.

Consistency

The **consistency** aspect of the ACID model mainly involves internal `InnoDB` processing to protect data from crashes. Related MySQL features include:

- The `InnoDB` doublewrite buffer. See [Section 15.6.4, “Doublewrite Buffer”](#).
- `InnoDB` crash recovery. See [InnoDB Crash Recovery](#).

Isolation

The **isolation** aspect of the ACID model mainly involves `InnoDB transactions`, in particular the **isolation level** that applies to each transaction. Related MySQL features include:

- The `autocommit` setting.
- Transaction isolation levels and the `SET TRANSACTION` statement. See [Section 15.7.2.1, “Transaction Isolation Levels”](#).
- The low-level details of `InnoDB locking`. Details can be viewed in the `INFORMATION_SCHEMA` tables (see [Section 15.15.2, “InnoDB INFORMATION_SCHEMA Transaction and Locking Information”](#)) and Performance Schema `data_locks` and `data_lock_waits` tables.

Durability

The **durability** aspect of the ACID model involves MySQL software features interacting with your particular hardware configuration. Because of the many possibilities depending on the capabilities of your CPU, network, and storage devices, this aspect is the most complicated to provide concrete guidelines for. (And those guidelines might take the form of “buy new hardware”.) Related MySQL features include:

- The `InnoDB` doublewrite buffer. See [Section 15.6.4, “Doublewrite Buffer”](#).
- The `innodb_flush_log_at_trx_commit` variable.
- The `sync_binlog` variable.
- The `innodb_file_per_table` variable.
- The write buffer in a storage device, such as a disk drive, SSD, or RAID array.
- A battery-backed cache in a storage device.
- The operating system used to run MySQL, in particular its support for the `fsync()` system call.
- An uninterruptible power supply (UPS) protecting the electrical power to all computer servers and storage devices that run MySQL servers and store MySQL data.
- Your backup strategy, such as frequency and types of backups, and backup retention periods.
- For distributed or hosted data applications, the particular characteristics of the data centers where the hardware for the MySQL servers is located, and network connections between the data centers.

15.3 InnoDB Multi-Versioning

[InnoDB](#) is a multi-version storage engine. It keeps information about old versions of changed rows to support transactional features such as concurrency and rollback. This information is stored in undo tablespaces in a data structure called a rollback segment. See [Section 15.6.3.4, “Undo Tablespaces”](#). [InnoDB](#) uses the information in the rollback segment to perform the undo operations needed in a transaction rollback. It also uses the information to build earlier versions of a row for a consistent read. See [Section 15.7.2.3, “Consistent Nonlocking Reads”](#).

Internally, [InnoDB](#) adds three fields to each row stored in the database:

- A 6-byte `DB_TRX_ID` field indicates the transaction identifier for the last transaction that inserted or updated the row. Also, a deletion is treated internally as an update where a special bit in the row is set to mark it as deleted.
- A 7-byte `DB_ROLL_PTR` field called the roll pointer. The roll pointer points to an undo log record written to the rollback segment. If the row was updated, the undo log record contains the information necessary to rebuild the content of the row before it was updated.
- A 6-byte `DB_ROW_ID` field contains a row ID that increases monotonically as new rows are inserted. If [InnoDB](#) generates a clustered index automatically, the index contains row ID values. Otherwise, the `DB_ROW_ID` column does not appear in any index.

Undo logs in the rollback segment are divided into insert and update undo logs. Insert undo logs are needed only in transaction rollback and can be discarded as soon as the transaction commits. Update undo logs are used also in consistent reads, but they can be discarded only after there is no transaction present for which [InnoDB](#) has assigned a snapshot that in a consistent read could require the information in the update undo log to build an earlier version of a database row. For additional information about undo logs, see [Section 15.6.6, “Undo Logs”](#).

It is recommended that you commit transactions regularly, including transactions that issue only consistent reads. Otherwise, [InnoDB](#) cannot discard data from the update undo logs, and the rollback segment may grow too big, filling up the undo tablespace in which it resides. For information about managing undo tablespaces, see [Section 15.6.3.4, “Undo Tablespaces”](#).

The physical size of an undo log record in the rollback segment is typically smaller than the corresponding inserted or updated row. You can use this information to calculate the space needed for your rollback segment.

In the [InnoDB](#) multi-versioning scheme, a row is not physically removed from the database immediately when you delete it with an SQL statement. [InnoDB](#) only physically removes the corresponding row and its index records when it discards the update undo log record written for the deletion. This removal operation is called a purge, and it is quite fast, usually taking the same order of time as the SQL statement that did the deletion.

If you insert and delete rows in smallish batches at about the same rate in the table, the purge thread can start to lag behind and the table can grow bigger and bigger because of all the “dead” rows, making everything disk-bound and very slow. In such cases, throttle new row operations, and allocate more resources to the purge thread by tuning the `innodb_max_purge_lag` system variable. For more information, see [Section 15.8.9, “Purge Configuration”](#).

Multi-Versioning and Secondary Indexes

[InnoDB](#) multiversion concurrency control (MVCC) treats secondary indexes differently than clustered indexes. Records in a clustered index are updated in-place, and their hidden system columns point undo log entries from which earlier versions of records can be reconstructed. Unlike clustered index records, secondary index records do not contain hidden system columns nor are they updated in-place.

When a secondary index column is updated, old secondary index records are delete-marked, new records are inserted, and delete-marked records are eventually purged. When a secondary index

record is delete-marked or the secondary index page is updated by a newer transaction, InnoDB looks up the database record in the clustered index. In the clustered index, the record's DB_TRX_ID is checked, and the correct version of the record is retrieved from the undo log if the record was modified after the reading transaction was initiated.

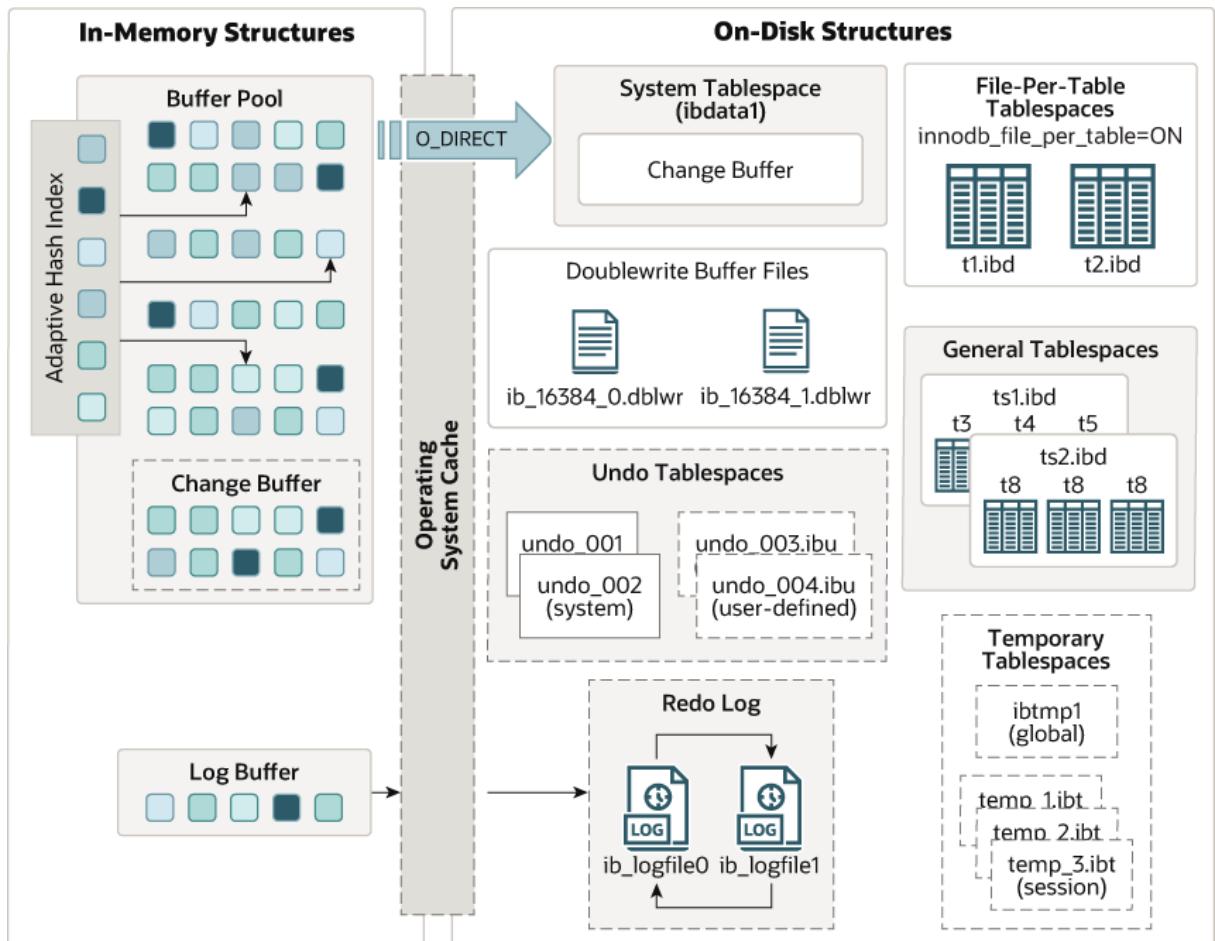
If a secondary index record is marked for deletion or the secondary index page is updated by a newer transaction, the [covering index](#) technique is not used. Instead of returning values from the index structure, InnoDB looks up the record in the clustered index.

However, if the [index condition pushdown \(ICP\)](#) optimization is enabled, and parts of the WHERE condition can be evaluated using only fields from the index, the MySQL server still pushes this part of the WHERE condition down to the storage engine where it is evaluated using the index. If no matching records are found, the clustered index lookup is avoided. If matching records are found, even among delete-marked records, InnoDB looks up the record in the clustered index.

15.4 InnoDB Architecture

The following diagram shows in-memory and on-disk structures that comprise the InnoDB storage engine architecture. For information about each structure, see [Section 15.5, “InnoDB In-Memory Structures”](#), and [Section 15.6, “InnoDB On-Disk Structures”](#).

Figure 15.1 InnoDB Architecture



15.5 InnoDB In-Memory Structures

This section describes InnoDB in-memory structures and related topics.

15.5.1 Buffer Pool

The buffer pool is an area in main memory where [InnoDB](#) caches table and index data as it is accessed. The buffer pool permits frequently used data to be accessed directly from memory, which speeds up processing. On dedicated servers, up to 80% of physical memory is often assigned to the buffer pool.

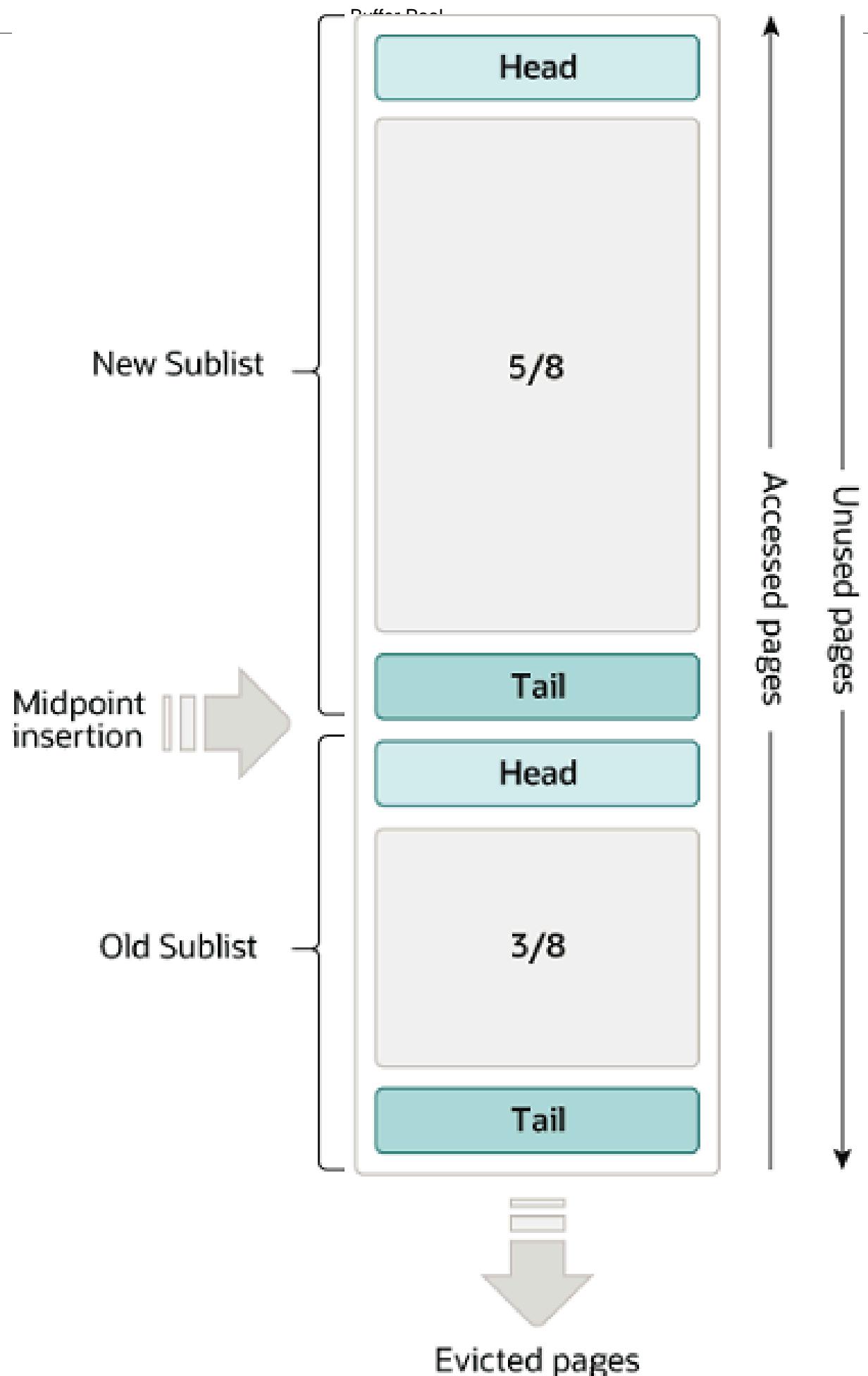
For efficiency of high-volume read operations, the buffer pool is divided into pages that can potentially hold multiple rows. For efficiency of cache management, the buffer pool is implemented as a linked list of pages; data that is rarely used is aged out of the cache using a variation of the least recently used (LRU) algorithm.

Knowing how to take advantage of the buffer pool to keep frequently accessed data in memory is an important aspect of MySQL tuning.

Buffer Pool LRU Algorithm

The buffer pool is managed as a list using a variation of the LRU algorithm. When room is needed to add a new page to the buffer pool, the least recently used page is evicted and a new page is added to the middle of the list. This midpoint insertion strategy treats the list as two sublists:

- At the head, a sublist of new (“young”) pages that were accessed recently
- At the tail, a sublist of old pages that were accessed less recently



The algorithm keeps frequently used pages in the new sublist. The old sublist contains less frequently used pages; these pages are candidates for [eviction](#).

By default, the algorithm operates as follows:

- 3/8 of the buffer pool is devoted to the old sublist.
- The midpoint of the list is the boundary where the tail of the new sublist meets the head of the old sublist.
- When [InnoDB](#) reads a page into the buffer pool, it initially inserts it at the midpoint (the head of the old sublist). A page can be read because it is required for a user-initiated operation such as an SQL query, or as part of a [read-ahead](#) operation performed automatically by [InnoDB](#).
- Accessing a page in the old sublist makes it “young”, moving it to the head of the new sublist. If the page was read because it was required by a user-initiated operation, the first access occurs immediately and the page is made young. If the page was read due to a read-ahead operation, the first access does not occur immediately and might not occur at all before the page is evicted.
- As the database operates, pages in the buffer pool that are not accessed “age” by moving toward the tail of the list. Pages in both the new and old sublists age as other pages are made new. Pages in the old sublist also age as pages are inserted at the midpoint. Eventually, a page that remains unused reaches the tail of the old sublist and is evicted.

By default, pages read by queries are immediately moved into the new sublist, meaning they stay in the buffer pool longer. A table scan, performed for a [mysqldump](#) operation or a [SELECT](#) statement with no [WHERE](#) clause, for example, can bring a large amount of data into the buffer pool and evict an equivalent amount of older data, even if the new data is never used again. Similarly, pages that are loaded by the read-ahead background thread and accessed only once are moved to the head of the new list. These situations can push frequently used pages to the old sublist where they become subject to eviction. For information about optimizing this behavior, see [Section 15.8.3.3, “Making the Buffer Pool Scan Resistant”](#), and [Section 15.8.3.4, “Configuring InnoDB Buffer Pool Prefetching \(Read-Ahead\)”](#).

[InnoDB](#) Standard Monitor output contains several fields in the [BUFFER POOL AND MEMORY](#) section regarding operation of the buffer pool LRU algorithm. For details, see [Monitoring the Buffer Pool Using the InnoDB Standard Monitor](#).

Buffer Pool Configuration

You can configure the various aspects of the buffer pool to improve performance.

- Ideally, you set the size of the buffer pool to as large a value as practical, leaving enough memory for other processes on the server to run without excessive paging. The larger the buffer pool, the more [InnoDB](#) acts like an in-memory database, reading data from disk once and then accessing the data from memory during subsequent reads. See [Section 15.8.3.1, “Configuring InnoDB Buffer Pool Size”](#).
- On 64-bit systems with sufficient memory, you can split the buffer pool into multiple parts to minimize contention for memory structures among concurrent operations. For details, see [Section 15.8.3.2, “Configuring Multiple Buffer Pool Instances”](#).
- You can keep frequently accessed data in memory regardless of sudden spikes of activity from operations that would bring large amounts of infrequently accessed data into the buffer pool. For details, see [Section 15.8.3.3, “Making the Buffer Pool Scan Resistant”](#).
- You can control how and when to perform read-ahead requests to prefetch pages into the buffer pool asynchronously in anticipation of impending need for them. For details, see [Section 15.8.3.4, “Configuring InnoDB Buffer Pool Prefetching \(Read-Ahead\)”](#).
- You can control when background flushing occurs and whether or not the rate of flushing is dynamically adjusted based on workload. For details, see [Section 15.8.3.5, “Configuring Buffer Pool Flushing”](#).

- You can configure how InnoDB preserves the current buffer pool state to avoid a lengthy warmup period after a server restart. For details, see [Section 15.8.3.6, “Saving and Restoring the Buffer Pool State”](#).

Monitoring the Buffer Pool Using the InnoDB Standard Monitor

InnoDB Standard Monitor output, which can be accessed using `SHOW ENGINE INNODB STATUS`, provides metrics regarding operation of the buffer pool. Buffer pool metrics are located in the `BUFFER POOL AND MEMORY` section of InnoDB Standard Monitor output:

```
-----
BUFFER POOL AND MEMORY
-----
Total large memory allocated 2198863872
Dictionary memory allocated 776332
Buffer pool size    131072
Free buffers        124908
Database pages      5720
Old database pages  2071
Modified db pages   910
Pending reads       0
Pending writes: LRU 0, flush list 0, single page 0
Pages made young 4, not young 0
0.10 youngs/s, 0.00 non-youngs/s
Pages read 197, created 5523, written 5060
0.00 reads/s, 190.89 creates/s, 244.94 writes/s
Buffer pool hit rate 1000 / 1000, young-making rate 0 / 1000 not
0 / 1000
Pages read ahead 0.00/s, evicted without access 0.00/s, Random read
ahead 0.00/s
LRU len: 5720, unzip_LRU len: 0
I/O sum[0]:cur[0], unzip sum[0]:cur[0]
```

The following table describes buffer pool metrics reported by the InnoDB Standard Monitor.

Per second averages provided in InnoDB Standard Monitor output are based on the elapsed time since InnoDB Standard Monitor output was last printed.

Table 15.2 InnoDB Buffer Pool Metrics

Name	Description
Total memory allocated	The total memory allocated for the buffer pool in bytes.
Dictionary memory allocated	The total memory allocated for the InnoDB data dictionary in bytes.
Buffer pool size	The total size in pages allocated to the buffer pool.
Free buffers	The total size in pages of the buffer pool free list.
Database pages	The total size in pages of the buffer pool LRU list.
Old database pages	The total size in pages of the buffer pool old LRU sublist.
Modified db pages	The current number of pages modified in the buffer pool.
Pending reads	The number of buffer pool pages waiting to be read into the buffer pool.
Pending writes LRU	The number of old dirty pages within the buffer pool to be written from the bottom of the LRU list.
Pending writes flush list	The number of buffer pool pages to be flushed during checkpointing.
Pending writes single page	The number of pending independent page writes within the buffer pool.

Name	Description
Pages made young	The total number of pages made young in the buffer pool LRU list (moved to the head of sublist of “new” pages).
Pages made not young	The total number of pages not made young in the buffer pool LRU list (pages that have remained in the “old” sublist without being made young).
youngs/s	The per second average of accesses to old pages in the buffer pool LRU list that have resulted in making pages young. See the notes that follow this table for more information.
non-youngs/s	The per second average of accesses to old pages in the buffer pool LRU list that have resulted in not making pages young. See the notes that follow this table for more information.
Pages read	The total number of pages read from the buffer pool.
Pages created	The total number of pages created within the buffer pool.
Pages written	The total number of pages written from the buffer pool.
reads/s	The per second average number of buffer pool page reads per second.
creates/s	The average number of buffer pool pages created per second.
writes/s	The average number of buffer pool page writes per second.
Buffer pool hit rate	The buffer pool page hit rate for pages read from the buffer pool vs from disk storage.
young-making rate	The average hit rate at which page accesses have resulted in making pages young. See the notes that follow this table for more information.
not (young-making rate)	The average hit rate at which page accesses have not resulted in making pages young. See the notes that follow this table for more information.
Pages read ahead	The per second average of read ahead operations.
Pages evicted without access	The per second average of the pages evicted without being accessed from the buffer pool.
Random read ahead	The per second average of random read ahead operations.
LRU len	The total size in pages of the buffer pool LRU list.
unzip_LRU len	The length (in pages) of the buffer pool unzip_LRU list.
I/O sum	The total number of buffer pool LRU list pages accessed.
I/O cur	The total number of buffer pool LRU list pages accessed in the current interval.

Name	Description
I/O unzip sum	The total number of buffer pool unzip_LRU list pages decompressed.
I/O unzip cur	The total number of buffer pool unzip_LRU list pages decompressed in the current interval.

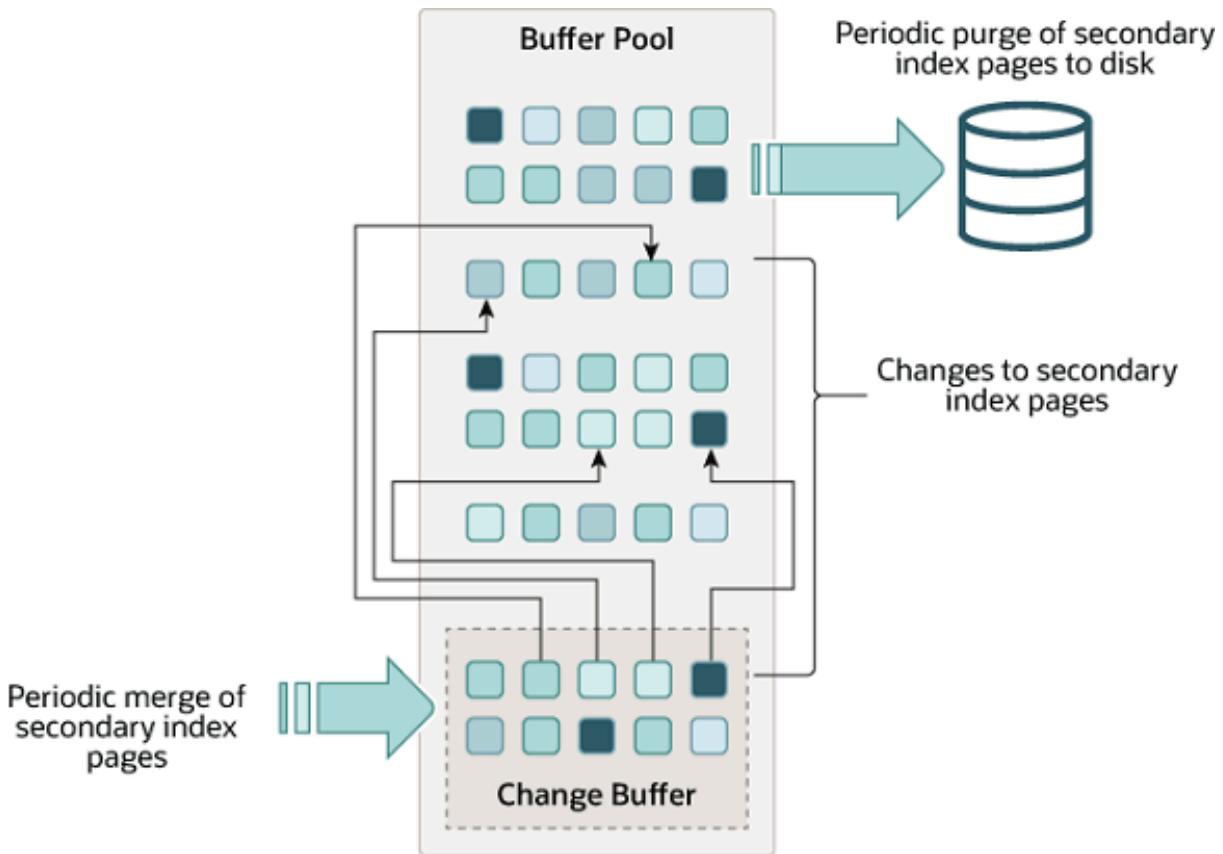
Notes:

- The `youngs/s` metric is applicable only to old pages. It is based on the number of page accesses. There can be multiple accesses for a given page, all of which are counted. If you see very low `youngs/s` values when there are no large scans occurring, consider reducing the delay time or increasing the percentage of the buffer pool used for the old sublist. Increasing the percentage makes the old sublist larger so that it takes longer for pages in that sublist to move to the tail, which increases the likelihood that those pages are accessed again and made young. See [Section 15.8.3.3, “Making the Buffer Pool Scan Resistant”](#).
- The `non-youngs/s` metric is applicable only to old pages. It is based on the number of page accesses. There can be multiple accesses for a given page, all of which are counted. If you do not see a higher `non-youngs/s` value when performing large table scans (and a higher `youngs/s` value), increase the delay value. See [Section 15.8.3.3, “Making the Buffer Pool Scan Resistant”](#).
- The `young-making` rate accounts for all buffer pool page accesses, not just accesses for pages in the old sublist. The `young-making` rate and `not` rate do not normally add up to the overall buffer pool hit rate. Page hits in the old sublist cause pages to move to the new sublist, but page hits in the new sublist cause pages to move to the head of the list only if they are a certain distance from the head.
- `not (young-making rate)` is the average hit rate at which page accesses have not resulted in making pages young due to the delay defined by `innodb_old_blocks_time` not being met, or due to page hits in the new sublist that did not result in pages being moved to the head. This rate accounts for all buffer pool page accesses, not just accesses for pages in the old sublist.

Buffer pool [server status variables](#) and the `INNODB_BUFFER_POOL_STATS` table provide many of the same buffer pool metrics found in [InnoDB Standard Monitor output](#). For more information, see [Example 15.10, “Querying the INNODB_BUFFER_POOL_STATS Table”](#).

15.5.2 Change Buffer

The change buffer is a special data structure that caches changes to [secondary index](#) pages when those pages are not in the [buffer pool](#). The buffered changes, which may result from [INSERT](#), [UPDATE](#), or [DELETE](#) operations (DML), are merged later when the pages are loaded into the buffer pool by other read operations.

Figure 15.3 Change Buffer

Unlike [clustered indexes](#), secondary indexes are usually nonunique, and inserts into secondary indexes happen in a relatively random order. Similarly, deletes and updates may affect secondary index pages that are not adjacently located in an index tree. Merging cached changes at a later time, when affected pages are read into the buffer pool by other operations, avoids substantial random access I/O that would be required to read secondary index pages into the buffer pool from disk.

Periodically, the purge operation that runs when the system is mostly idle, or during a slow shutdown, writes the updated index pages to disk. The purge operation can write disk blocks for a series of index values more efficiently than if each value were written to disk immediately.

Change buffer merging may take several hours when there are many affected rows and numerous secondary indexes to update. During this time, disk I/O is increased, which can cause a significant slowdown for disk-bound queries. Change buffer merging may also continue to occur after a transaction is committed, and even after a server shutdown and restart (see [Section 15.21.3, “Forcing InnoDB Recovery”](#) for more information).

In memory, the change buffer occupies part of the buffer pool. On disk, the change buffer is part of the system tablespace, where index changes are buffered when the database server is shut down.

The type of data cached in the change buffer is governed by the `innodb_change_buffering` variable. For more information, see [Configuring Change Buffering](#). You can also configure the maximum change buffer size. For more information, see [Configuring the Change Buffer Maximum Size](#).

Change buffering is not supported for a secondary index if the index contains a descending index column or if the primary key includes a descending index column.

For answers to frequently asked questions about the change buffer, see [Section A.16, “MySQL 8.0 FAQ: InnoDB Change Buffer”](#).

Configuring Change Buffering

When `INSERT`, `UPDATE`, and `DELETE` operations are performed on a table, the values of indexed columns (particularly the values of secondary keys) are often in an unsorted order, requiring substantial I/O to bring secondary indexes up to date. The `change buffer` caches changes to secondary index entries when the relevant `page` is not in the `buffer pool`, thus avoiding expensive I/O operations by not immediately reading in the page from disk. The buffered changes are merged when the page is loaded into the buffer pool, and the updated page is later flushed to disk. 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, change buffering is most valuable for workloads that are I/O-bound; for example, applications with a high volume of DML operations such as bulk inserts benefit from change buffering.

However, the change 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 change buffering. If the working data set fits entirely within the buffer pool, change buffering does not impose extra overhead, because it only applies to pages that are not in the buffer pool.

The `innodb_change_buffering` variable controls the extent to which `InnoDB` performs change buffering. You can enable or disable 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 a combination of an insert and a delete. The default `innodb_change_buffering` value is `all`.

Permitted `innodb_change_buffering` values include:

- `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 operations.

- `purges`

Buffer the physical deletion operations that happen in the background.

You can set the `innodb_change_buffering` variable in the MySQL option file (`my.cnf` or `my.ini`) or change it dynamically with the `SET GLOBAL` statement, which requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#). Changing the setting affects the buffering of new operations; the merging of existing buffered entries is not affected.

Configuring the Change Buffer Maximum Size

The `innodb_change_buffer_max_size` variable permits configuring the maximum size of the change buffer as a percentage of the total size of the buffer pool. By default, `innodb_change_buffer_max_size` is set to 25. The maximum setting is 50.

Consider increasing `innodb_change_buffer_max_size` on a MySQL server with heavy insert, update, and delete activity, where change buffer merging does not keep pace with new change buffer entries, causing the change buffer to reach its maximum size limit.

Consider decreasing `innodb_change_buffer_max_size` on a MySQL server with static data used for reporting, or if the change buffer consumes too much of the memory space shared with the buffer pool, causing pages to age out of the buffer pool sooner than desired.

Test different settings with a representative workload to determine an optimal configuration. The `innodb_change_buffer_max_size` variable is dynamic, which permits modifying the setting without restarting the server.

Monitoring the Change Buffer

The following options are available for change buffer monitoring:

- InnoDB Standard Monitor output includes change buffer status information. To view monitor data, issue the `SHOW ENGINE INNODB STATUS` statement.

```
mysql> SHOW ENGINE INNODB STATUS\G
```

Change buffer status information is located under the `INSERT BUFFER AND ADAPTIVE HASH INDEX` heading and appears similar to the following:

```
-----
INSERT BUFFER AND ADAPTIVE HASH INDEX
-----
Ibuf: size 1, free list len 0, seg size 2, 0 merges
merged operations:
  insert 0, delete mark 0, delete 0
discarded operations:
  insert 0, delete mark 0, delete 0
Hash table size 4425293, used cells 32, node heap has 1 buffer(s)
13577.57 hash searches/s, 202.47 non-hash searches/s
```

For more information, see [Section 15.17.3, “InnoDB Standard Monitor and Lock Monitor Output”](#).

- The Information Schema `INNODB_METRICS` table provides most of the data points found in InnoDB Standard Monitor output plus other data points. To view change buffer metrics and a description of each, issue the following query:

```
mysql> SELECT NAME, COMMENT FROM INFORMATION_SCHEMA.INNODB_METRICS WHERE NAME LIKE '%ibuf%'\G
```

See [Section 15.15.6, “InnoDB INFORMATION_SCHEMA Metrics Table”](#).

- The Information Schema `INNODB_BUFFER_PAGE` table provides metadata about each page in the buffer pool, including change buffer index and change buffer bitmap pages. Change buffer pages are identified by `PAGE_TYPE`. `IBUF_INDEX` is the page type for change buffer index pages, and `IBUF_BITMAP` is the page type for change buffer bitmap pages.



Warning

Querying the `INNODB_BUFFER_PAGE` table can introduce significant performance overhead. To avoid impacting performance, reproduce the issue you want to investigate on a test instance and run your queries on the test instance.

For example, you can query the `INNODB_BUFFER_PAGE` table to determine the approximate number of `IBUF_INDEX` and `IBUF_BITMAP` pages as a percentage of total buffer pool pages.

```
mysql> SELECT (SELECT COUNT(*) FROM INFORMATION_SCHEMA.INNODB_BUFFER_PAGE
    WHERE PAGE_TYPE LIKE 'IBUF%') AS change_buffer_pages,
    (SELECT COUNT(*) FROM INFORMATION_SCHEMA.INNODB_BUFFER_PAGE) AS total_pages,
    ((change_buffer_pages/total_pages)*100)
```

AS change_buffer_page_percentage;		
change_buffer_pages	total_pages	change_buffer_page_percentage
25	8192	0.3052

For information about other data provided by the `INNODB_BUFFER_PAGE` table, see [Section 26.4.2, “The INFORMATION_SCHEMA INNODB_BUFFER_PAGE Table”](#). For related usage information, see [Section 15.15.5, “InnoDB INFORMATION_SCHEMA Buffer Pool Tables”](#).

- **Performance Schema** provides change buffer mutex wait instrumentation for advanced performance monitoring. To view change buffer instrumentation, issue the following query:

mysql> SELECT * FROM performance_schema.setup_instruments WHERE NAME LIKE '%wait/synch/mutex/innodb/ibuf%';		
NAME	ENABLED	TIMED
wait/synch/mutex/innodb/ibuf_bitmap_mutex	YES	YES
wait/synch/mutex/innodb/ibuf_mutex	YES	YES
wait/synch/mutex/innodb/ibuf_pessimistic_insert_mutex	YES	YES

For information about monitoring InnoDB mutex waits, see [Section 15.16.2, “Monitoring InnoDB Mutex Waits Using Performance Schema”](#).

15.5.3 Adaptive Hash Index

The adaptive hash index enables InnoDB to perform more like an in-memory database on systems with appropriate combinations of workload and sufficient memory for the buffer pool without sacrificing transactional features or reliability. The adaptive hash index is enabled by the `innodb_adaptive_hash_index` variable, or turned off at server startup by `--skip-innodb-adaptive-hash-index`.

Based on the observed pattern of searches, a hash index is built using a prefix of the index key. The prefix can be any length, and it may be that only some values in the B-tree appear in the hash index. Hash indexes are built on demand for the pages of the index that are accessed often.

If a table fits almost entirely in main memory, a hash index speeds up queries by enabling direct lookup of any element, turning the index value into a sort of pointer. InnoDB has a mechanism that monitors index searches. If InnoDB notices that queries could benefit from building a hash index, it does so automatically.

With some workloads, the speedup from hash index lookups greatly outweighs the extra work to monitor index lookups and maintain the hash index structure. Access to the adaptive hash index can sometimes become a source of contention under heavy workloads, such as multiple concurrent joins. Queries with `LIKE` operators and `%` wildcards also tend not to benefit. For workloads that do not benefit from the adaptive hash index, turning it off reduces unnecessary performance overhead. Because it is difficult to predict in advance whether the adaptive hash index is appropriate for a particular system and workload, consider running benchmarks with it enabled and disabled.

The adaptive hash index feature is partitioned. Each index is bound to a specific partition, and each partition is protected by a separate latch. Partitioning is controlled by the `innodb_adaptive_hash_index_parts` variable. The `innodb_adaptive_hash_index_parts` variable is set to 8 by default. The maximum setting is 512.

You can monitor adaptive hash index use and contention in the `SEMAPHORES` section of `SHOW ENGINE INNODB STATUS` output. If there are numerous threads waiting on rw-latches created in `btr0sea.c`, consider increasing the number of adaptive hash index partitions or disabling the adaptive hash index.

For information about the performance characteristics of hash indexes, see [Section 8.3.9, “Comparison of B-Tree and Hash Indexes”](#).

15.5.4 Log Buffer

The log buffer is the memory area that holds data to be written to the log files on disk. Log buffer size is defined by the `innodb_log_buffer_size` variable. The default size is 16MB. The contents of the log buffer are periodically flushed to disk. A large log buffer enables large transactions to run without the need to write redo log data to disk before the transactions commit. Thus, if you have transactions that update, insert, or delete many rows, increasing the size of the log buffer saves disk I/O.

The `innodb_flush_log_at_trx_commit` variable controls how the contents of the log buffer are written and flushed to disk. The `innodb_flush_log_at_timeout` variable controls log flushing frequency.

For related information, see [Memory Configuration](#), and [Section 8.5.4, “Optimizing InnoDB Redo Logging”](#).

15.6 InnoDB On-Disk Structures

This section describes [InnoDB](#) on-disk structures and related topics.

15.6.1 Tables

This section covers topics related to [InnoDB](#) tables.

15.6.1.1 Creating InnoDB Tables

[InnoDB](#) tables are created using the `CREATE TABLE` statement; for example:

```
CREATE TABLE t1 (a INT, b CHAR (20), PRIMARY KEY (a)) ENGINE=InnoDB;
```

The `ENGINE=InnoDB` clause is not required when [InnoDB](#) is defined as the default storage engine, which it is by default. However, the `ENGINE` clause is useful if the `CREATE TABLE` statement is to be replayed on a different MySQL Server instance where the default storage engine is not [InnoDB](#) or is unknown. You can determine the default storage engine on a MySQL Server instance by issuing the following statement:

```
mysql> SELECT @@default_storage_engine;
+-----+
| @@default_storage_engine |
+-----+
| InnoDB                   |
+-----+
```

[InnoDB](#) tables are created in file-per-table tablespaces by default. To create an [InnoDB](#) table in the [InnoDB](#) system tablespace, disable the `innodb_file_per_table` variable before creating the table. To create an [InnoDB](#) table in a general tablespace, use `CREATE TABLE ... TABLESPACE` syntax. For more information, see [Section 15.6.3, “Tablespaces”](#).

Row Formats

The row format of an [InnoDB](#) table determines how its rows are physically stored on disk. [InnoDB](#) supports four row formats, each with different storage characteristics. Supported row formats include `REDUNDANT`, `COMPACT`, `DYNAMIC`, and `COMPRESSED`. The `DYNAMIC` row format is the default. For information about row format characteristics, see [Section 15.10, “InnoDB Row Formats”](#).

The `innodb_default_row_format` variable defines the default row format. The row format of a table can also be defined explicitly using the `ROW_FORMAT` table option in a `CREATE TABLE` or `ALTER TABLE` statement. See [Defining the Row Format of a Table](#).

Primary Keys

It is recommended that you define a primary key for each table that you create. When selecting primary key columns, choose columns with the following characteristics:

- Columns that are referenced by the most important queries.
- Columns that are never left blank.
- Columns that never have duplicate values.
- Columns that rarely if ever change value once inserted.

For example, in a table containing information about people, you would not create a primary key on `(firstname, lastname)` because more than one person can have the same name, a name column may be left blank, and sometimes people change their names. With so many constraints, often there is not an obvious set of columns to use as a primary key, so you create a new column with a numeric ID to serve as all or part of the primary key. You can declare an `auto-increment` column so that ascending values are filled in automatically as rows are inserted:

```
# The value of ID can act like a pointer between related items in different tables.
CREATE TABLE t5 (id INT AUTO_INCREMENT, b CHAR (20), PRIMARY KEY (id));

# The primary key can consist of more than one column. Any autoinc column must come first.
CREATE TABLE t6 (id INT AUTO_INCREMENT, a INT, b CHAR (20), PRIMARY KEY (id,a));
```

For more information about auto-increment columns, see [Section 15.6.1.6, “AUTO_INCREMENT Handling in InnoDB”](#).

Although a table works correctly without defining a primary key, the primary key is involved with many aspects of performance and is a crucial design aspect for any large or frequently used table. It is recommended that you always specify a primary key in the `CREATE TABLE` statement. If you create the table, load data, and then run `ALTER TABLE` to add a primary key later, that operation is much slower than defining the primary key when creating the table. For more information about primary keys, see [Section 15.6.2.1, “Clustered and Secondary Indexes”](#).

Viewing InnoDB Table Properties

To view the properties of an `InnoDB` table, issue a `SHOW TABLE STATUS` statement:

```
mysql> SHOW TABLE STATUS FROM test LIKE 't%';
***** 1. row *****
      Name: t1
      Engine: InnoDB
     Version: 10
   Row_format: Dynamic
       Rows: 0
 Avg_row_length: 0
   Data_length: 16384
Max_data_length: 0
   Index_length: 0
    Data_free: 0
Auto_increment: NULL
 Create_time: 2021-02-18 12:18:28
 Update_time: NULL
  Check_time: NULL
   Collation: utf8mb4_0900_ai_ci
    Checksum: NULL
Create_options:
     Comment:
```

For information about `SHOW TABLE STATUS` output, see [Section 13.7.7.38, “SHOW TABLE STATUS Statement”](#).

You can also access `InnoDB` table properties by querying the `InnoDB` Information Schema system tables:

```
mysql> SELECT * FROM INFORMATION_SCHEMA.INNODB_TABLES WHERE NAME='test/t1' \G
```

```
***** 1. row *****
TABLE_ID: 1144
  NAME: test/t1
  FLAG: 33
N_COLS: 5
  SPACE: 30
ROW_FORMAT: Dynamic
ZIP_PAGE_SIZE: 0
  SPACE_TYPE: Single
INSTANT_COLS: 0
```

For more information, see [Section 15.15.3, “InnoDB INFORMATION_SCHEMA Schema Object Tables”](#).

15.6.1.2 Creating Tables Externally

There are different reasons for creating InnoDB tables externally; that is, creating tables outside of the data directory. Those reasons might include space management, I/O optimization, or placing tables on a storage device with particular performance or capacity characteristics, for example.

InnoDB supports the following methods for creating tables externally:

- [Using the DATA DIRECTORY Clause](#)
- [Using CREATE TABLE ... TABLESPACE Syntax](#)
- [Creating a Table in an External General Tablespace](#)

Using the DATA DIRECTORY Clause

You can create an InnoDB table in an external directory by specifying a `DATA DIRECTORY` clause in the `CREATE TABLE` statement.

```
CREATE TABLE t1 (c1 INT PRIMARY KEY) DATA DIRECTORY = '/external/directory';
```

The `DATA DIRECTORY` clause is supported for tables created in file-per-table tablespaces. Tables are implicitly created in file-per-table tablespaces when the `innodb_file_per_table` variable is enabled, which it is by default.

```
mysql> SELECT @@innodb_file_per_table;
+-----+
| @@innodb_file_per_table |
+-----+
| 1 |
+-----+
```

For more information about file-per-table tablespaces, see [Section 15.6.3.2, “File-Per-Table Tablespaces”](#).

When you specify a `DATA DIRECTORY` clause in a `CREATE TABLE` statement, the table's data file (`table_name.ibd`) is created in a schema directory under the specified directory.

As of MySQL 8.0.21, tables and table partitions created outside of the data directory using the `DATA DIRECTORY` clause are restricted to directories known to InnoDB. This requirement permits database administrators to control where tablespace data files are created and ensures that data files can be found during recovery (see [Tablespace Discovery During Crash Recovery](#)). Known directories are those defined by the `datadir`, `innodb_data_home_dir`, and `innodb_directories` variables. You can use the following statement to check those settings:

```
mysql> SELECT @@datadir,@@innodb_data_home_dir,@@innodb_directories;
```

If the directory you want to use is unknown, add it to the `innodb_directories` setting before you create the table. The `innodb_directories` variable is read-only. Configuring it requires restarting

the server. For general information about setting system variables, see [Section 5.1.9, “Using System Variables”](#).

The following example demonstrates creating a table in an external directory using the `DATA DIRECTORY` clause. It is assumed that the `innodb_file_per_table` variable is enabled and that the directory is known to `InnoDB`.

```
mysql> USE test;
Database changed

mysql> CREATE TABLE t1 (c1 INT PRIMARY KEY) DATA DIRECTORY = '/external/directory';

# MySQL creates the table's data file in a schema directory
# under the external directory

$> cd /external/directory/test
$> ls
t1.ibd
```

Usage Notes:

- MySQL initially holds the tablespace data file open, preventing you from dismounting the device, but might eventually close the file if the server is busy. Be careful not to accidentally dismount an external device while MySQL is running, or start MySQL while the device is disconnected. Attempting to access a table when the associated data file is missing causes a serious error that requires a server restart.
- A server restart might fail if the data file is not found at the expected path. In this case, you can restore the tablespace data file from a backup or drop the table to remove the information about it from the [data dictionary](#).
- Before placing a table on an NFS-mounted volume, review potential issues outlined in [Using NFS with MySQL](#).
 - If using an LVM snapshot, file copy, or other file-based mechanism to back up the table's data file, always use the `FLUSH TABLES ... FOR EXPORT` statement first to ensure that all changes buffered in memory are [flushed](#) to disk before the backup occurs.
 - Using the `DATA DIRECTORY` clause to create a table in an external directory is an alternative to using [symbolic links](#), which `InnoDB` does not support.
 - The `DATA DIRECTORY` clause is not supported in a replication environment where the source and replica reside on the same host. The `DATA DIRECTORY` clause requires a full directory path. Replicating the path in this case would cause the source and replica to create the table in same location.
 - As of MySQL 8.0.21, tables created in file-per-table tablespaces can no longer be created in the undo tablespace directory (`innodb_undo_directory`) unless that directly is known to `InnoDB`. Known directories are those defined by the `datadir`, `innodb_data_home_dir`, and `innodb_directories` variables.

Using `CREATE TABLE ... TABLESPACE` Syntax

`CREATE TABLE ... TABLESPACE` syntax can be used in combination with the `DATA DIRECTORY` clause to create a table in an external directory. To do so, specify `innodb_file_per_table` as the tablespace name.

```
mysql> CREATE TABLE t2 (c1 INT PRIMARY KEY) TABLESPACE = innodb_file_per_table
        DATA DIRECTORY = '/external/directory';
```

This method is supported only for tables created in file-per-table tablespaces, but does not require the `innodb_file_per_table` variable to be enabled. In all other respects, this method is equivalent to the `CREATE TABLE ... DATA DIRECTORY` method described above. The same usage notes apply.

Creating a Table in an External General Tablespace

You can create a table in a general tablespace that resides in an external directory.

- For information about creating a general tablespace in an external directory, see [Creating a General Tablespace](#).
- For information about creating a table in a general tablespace, see [Adding Tables to a General Tablespace](#).

15.6.1.3 Importing InnoDB Tables

This section describes how to import tables using the *Transportable Tablespaces* feature, which permits importing tables, partitioned tables, or individual table partitions that reside in file-per-table tablespaces. There are many reasons why you might want to import tables:

- To run reports on a non-production MySQL server instance to avoid placing extra load on a production server.
- To copy data to a new replica server.
- To restore a table from a backed-up tablespace file.
- As a faster way of moving data than importing a dump file, which requires reinserting data and rebuilding indexes.
- To move a data to a server with storage media that is better suited to your storage requirements. For example, you might move busy tables to an SSD device, or move large tables to a high-capacity HDD device.

The *Transportable Tablespaces* feature is described under the following topics in this section:

- [Prerequisites](#)
- [Importing Tables](#)
- [Importing Partitioned Tables](#)
- [Importing Table Partitions](#)
- [Limitations](#)
- [Usage Notes](#)
- [Internals](#)

Prerequisites

- The `innodb_file_per_table` variable must be enabled, which it is by default.
- The page size of the tablespace must match the page size of the destination MySQL server instance. `InnoDB` page size is defined by the `innodb_page_size` variable, which is configured when initializing a MySQL server instance.
- If the table has a foreign key relationship, `foreign_key_checks` must be disabled before executing `DISCARD TABLESPACE`. Also, you should export all foreign key related tables at the same logical point in time, as `ALTER TABLE ... IMPORT TABLESPACE` does not enforce foreign key constraints on imported data. To do so, stop updating the related tables, commit all transactions, acquire shared locks on the tables, and perform the export operations.
- When importing a table from another MySQL server instance, both MySQL server instances must have General Availability (GA) status and must be the same version. Otherwise, the table must be created on the same MySQL server instance into which it is being imported.

- If the table was created in an external directory by specifying the `DATA DIRECTORY` clause in the `CREATE TABLE` statement, the table that you replace on the destination instance must be defined with the same `DATA DIRECTORY` clause. A schema mismatch error is reported if the clauses do not match. To determine if the source table was defined with a `DATA DIRECTORY` clause, use `SHOW CREATE TABLE` to view the table definition. For information about using the `DATA DIRECTORY` clause, see [Section 15.6.1.2, “Creating Tables Externally”](#).
- If a `ROW_FORMAT` option is not defined explicitly in the table definition or `ROW_FORMAT=DEFAULT` is used, the `innodb_default_row_format` setting must be the same on the source and destination instances. Otherwise, a schema mismatch error is reported when you attempt the import operation. Use `SHOW CREATE TABLE` to check the table definition. Use `SHOW VARIABLES` to check the `innodb_default_row_format` setting. For related information, see [Defining the Row Format of a Table](#).

Importing Tables

This example demonstrates how to import a regular non-partitioned table that resides in a file-per-table tablespace.

1. On the destination instance, create a table with the same definition as the table you intend to import. (You can obtain the table definition using `SHOW CREATE TABLE` syntax.) If the table definition does not match, a schema mismatch error is reported when you attempt the import operation.

```
mysql> USE test;
mysql> CREATE TABLE t1 (c1 INT) ENGINE=INNODB;
```

2. On the destination instance, discard the tablespace of the table that you just created. (Before importing, you must discard the tablespace of the receiving table.)

```
mysql> ALTER TABLE t1 DISCARD TABLESPACE;
```

3. On the source instance, run `FLUSH TABLES ... FOR EXPORT` to quiesce the table you intend to import. When a table is quiesced, only read-only transactions are permitted on the table.

```
mysql> USE test;
mysql> FLUSH TABLES t1 FOR EXPORT;
```

`FLUSH TABLES ... FOR EXPORT` ensures that changes to the named table are flushed to disk so that a binary table copy can be made while the server is running. When `FLUSH TABLES ... FOR EXPORT` is run, InnoDB generates a `.cfg` metadata file in the schema directory of the table. The `.cfg` file contains metadata that is used for schema verification during the import operation.



Note

The connection executing `FLUSH TABLES ... FOR EXPORT` must remain open while the operation is running; otherwise, the `.cfg` file is removed as locks are released upon connection closure.

4. Copy the `.ibd` file and `.cfg` metadata file from the source instance to the destination instance. For example:

```
$> scp /path/to/datadir/test/t1.{ibd,cfg} destination-server:/path/to/datadir/test
```

The `.ibd` file and `.cfg` file must be copied before releasing the shared locks, as described in the next step.



Note

If you are importing a table from an encrypted tablespace, InnoDB generates a `.cfp` file in addition to a `.cfg` metadata file. The `.cfp` file must be copied to the destination instance together with the `.cfg` file. The

.cfp file contains a transfer key and an encrypted tablespace key. On import, InnoDB uses the transfer key to decrypt the tablespace key. For related information, see [Section 15.13, “InnoDB Data-at-Rest Encryption”](#).

5. On the source instance, use `UNLOCK TABLES` to release the locks acquired by the `FLUSH TABLES ... FOR EXPORT` statement:

```
mysql> USE test;
mysql> UNLOCK TABLES;
```

The `UNLOCK TABLES` operation also removes the `.cfg` file.

6. On the destination instance, import the tablespace:

```
mysql> USE test;
mysql> ALTER TABLE t1 IMPORT TABLESPACE;
```

Importing Partitioned Tables

This example demonstrates how to import a partitioned table, where each table partition resides in a file-per-table tablespace.

1. On the destination instance, create a partitioned table with the same definition as the partitioned table that you want to import. (You can obtain the table definition using `SHOW CREATE TABLE` syntax.) If the table definition does not match, a schema mismatch error is reported when you attempt the import operation.

```
mysql> USE test;
mysql> CREATE TABLE t1 (i int) ENGINE = InnoDB PARTITION BY KEY (i) PARTITIONS 3;
```

In the `/datadir/test` directory, there is a tablespace `.ibd` file for each of the three partitions.

```
mysql> ! ls /path/to/datadir/test/
t1#p#p0.ibd  t1#p#p1.ibd  t1#p#p2.ibd
```

2. On the destination instance, discard the tablespace for the partitioned table. (Before the import operation, you must discard the tablespace of the receiving table.)

```
mysql> ALTER TABLE t1 DISCARD TABLESPACE;
```

The three tablespace `.ibd` files of the partitioned table are discarded from the `/datadir/test` directory.

3. On the source instance, run `FLUSH TABLES ... FOR EXPORT` to quiesce the partitioned table that you intend to import. When a table is quiesced, only read-only transactions are permitted on the table.

```
mysql> USE test;
mysql> FLUSH TABLES t1 FOR EXPORT;
```

`FLUSH TABLES ... FOR EXPORT` ensures that changes to the named table are flushed to disk so that binary table copy can be made while the server is running. When `FLUSH TABLES ... FOR EXPORT` is run, InnoDB generates `.cfg` metadata files in the schema directory of the table for each of the table's tablespace files.

```
mysql> ! ls /path/to/datadir/test/
t1#p#p0.ibd  t1#p#p1.ibd  t1#p#p2.ibd
t1#p#p0.cfg  t1#p#p1.cfg  t1#p#p2.cfg
```

The `.cfg` files contain metadata that is used for schema verification when importing the tablespace. `FLUSH TABLES ... FOR EXPORT` can only be run on the table, not on individual table partitions.

4. Copy the `.ibd` and `.cfg` files from the source instance schema directory to the destination instance schema directory. For example:

```
$>scp /path/to/datadir/test/t1*.{ibd,cfg} destination-server:/path/to/datadir/test
```

The `.ibd` and `.cfg` files must be copied before releasing the shared locks, as described in the next step.



Note

If you are importing a table from an encrypted tablespace, InnoDB generates a `.cfp` files in addition to a `.cfg` metadata files. The `.cfp` files must be copied to the destination instance together with the `.cfg` files. The `.cfp` files contain a transfer key and an encrypted tablespace key. On import, InnoDB uses the transfer key to decrypt the tablespace key. For related information, see [Section 15.13, “InnoDB Data-at-Rest Encryption”](#).

5. On the source instance, use `UNLOCK TABLES` to release the locks acquired by `FLUSH TABLES ... FOR EXPORT`:

```
mysql> USE test;
mysql> UNLOCK TABLES;
```

6. On the destination instance, import the tablespace of the partitioned table:

```
mysql> USE test;
mysql> ALTER TABLE t1 IMPORT TABLESPACE;
```

Importing Table Partitions

This example demonstrates how to import individual table partitions, where each partition resides in a file-per-table tablespace file.

In the following example, two partitions (`p2` and `p3`) of a four-partition table are imported.

1. On the destination instance, create a partitioned table with the same definition as the partitioned table that you want to import partitions from. (You can obtain the table definition using `SHOW CREATE TABLE` syntax.) If the table definition does not match, a schema mismatch error is reported when you attempt the import operation.

```
mysql> USE test;
mysql> CREATE TABLE t1 (i int) ENGINE = InnoDB PARTITION BY KEY (i) PARTITIONS 4;
```

In the `/datadir/test` directory, there is a tablespace `.ibd` file for each of the four partitions.

```
mysql> ! ls /path/to/datadir/test/
t1#p#p0.ibd  t1#p#p1.ibd  t1#p#p2.ibd  t1#p#p3.ibd
```

2. On the destination instance, discard the partitions that you intend to import from the source instance. (Before importing partitions, you must discard the corresponding partitions from the receiving partitioned table.)

```
mysql> ALTER TABLE t1 DISCARD PARTITION p2, p3 TABLESPACE;
```

The tablespace `.ibd` files for the two discarded partitions are removed from the `/datadir/test` directory on the destination instance, leaving the following files:

```
mysql> ! ls /path/to/datadir/test/
t1#p#p0.ibd  t1#p#p1.ibd
```



Note

When `ALTER TABLE ... DISCARD PARTITION ... TABLESPACE` is run on subpartitioned tables, both partition and subpartition table names are permitted. When a partition name is specified, subpartitions of that partition are included in the operation.

3. On the source instance, run `FLUSH TABLES ... FOR EXPORT` to quiesce the partitioned table. When a table is quiesced, only read-only transactions are permitted on the table.

```
mysql> USE test;
mysql> FLUSH TABLES t1 FOR EXPORT;
```

`FLUSH TABLES ... FOR EXPORT` ensures that changes to the named table are flushed to disk so that binary table copy can be made while the instance is running. When `FLUSH TABLES ... FOR EXPORT` is run, InnoDB generates a `.cfg` metadata file for each of the table's tablespace files in the schema directory of the table.

```
mysql> \! ls /path/to/datadir/test/
t1#p#p0.ibd t1#p#p1.ibd t1#p#p2.ibd t1#p#p3.ibd
t1#p#p0.cfg t1#p#p1.cfg t1#p#p2.cfg t1#p#p3.cfg
```

The `.cfg` files contain metadata that used for schema verification during the import operation. `FLUSH TABLES ... FOR EXPORT` can only be run on the table, not on individual table partitions.

4. Copy the `.ibd` and `.cfg` files for partition `p2` and partition `p3` from the source instance schema directory to the destination instance schema directory.

```
$> scp t1#p#p2.ibd t1#p#p2.cfg t1#p#p3.ibd t1#p#p3.cfg destination-server:/path/to/datadir/test
```

The `.ibd` and `.cfg` files must be copied before releasing the shared locks, as described in the next step.



Note

If you are importing partitions from an encrypted tablespace, InnoDB generates a `.cfp` files in addition to a `.cfg` metadata files. The `.cfp` files must be copied to the destination instance together with the `.cfg` files. The `.cfp` files contain a transfer key and an encrypted tablespace key. On import, InnoDB uses the transfer key to decrypt the tablespace key. For related information, see [Section 15.13, “InnoDB Data-at-Rest Encryption”](#).

5. On the source instance, use `UNLOCK TABLES` to release the locks acquired by `FLUSH TABLES ... FOR EXPORT`:

```
mysql> USE test;
mysql> UNLOCK TABLES;
```

6. On the destination instance, import table partitions `p2` and `p3`:

```
mysql> USE test;
mysql> ALTER TABLE t1 IMPORT PARTITION p2, p3 TABLESPACE;
```



Note

When `ALTER TABLE ... IMPORT PARTITION ... TABLESPACE` is run on subpartitioned tables, both partition and subpartition table names are permitted. When a partition name is specified, subpartitions of that partition are included in the operation.

Limitations

- The *Transportable Tablespaces* feature is only supported for tables that reside in file-per-table tablespaces. It is not supported for the tables that reside in the system tablespace or general tablespaces. Tables in shared tablespaces cannot be quiesced.
- `FLUSH TABLES ... FOR EXPORT` is not supported on tables with a `FULLTEXT` index, as full-text search auxiliary tables cannot be flushed. After importing a table with a `FULLTEXT` index, run `OPTIMIZE TABLE` to rebuild the `FULLTEXT` indexes. Alternatively, drop `FULLTEXT` indexes before the export operation and recreate the indexes after importing the table on the destination instance.

- Due to a `.cfg` metadata file limitation, schema mismatches are not reported for partition type or partition definition differences when importing a partitioned table. Column differences are reported.
- Prior to MySQL 8.0.19, index key part sort order information is not stored to the `.cfg` metadata file used during a tablespace import operation. The index key part sort order is therefore assumed to be ascending, which is the default. As a result, records could be sorted in an unintended order if one table involved in the import operation is defined with a DESC index key part sort order and the other table is not. The workaround is to drop and recreate affected indexes. For information about index key part sort order, see [Section 13.1.15, “CREATE INDEX Statement”](#).

The `.cfg` file format was updated in MySQL 8.0.19 to include index key part sort order information. The issue described above does not affect import operations between MySQL 8.0.19 server instances or higher.

Usage Notes

- With the exception of tables that contain instantly added or dropped columns, `ALTER TABLE ... IMPORT TABLESPACE` does not require a `.cfg` metadata file to import a table. However, metadata checks are not performed when importing without a `.cfg` file, and a warning similar to the following is issued:

```
Message: InnoDB: IO Read error: (2, No such file or directory) Error opening './test\t.cfg', will attempt to import without schema verification
1 row in set (0.00 sec)
```

Importing a table without a `.cfg` metadata file should only be considered if no schema mismatches are expected and the table does not contain any instantly added or dropped columns. The ability to import without a `.cfg` file could be useful in crash recovery scenarios where metadata is not accessible.

Attempting to import a table with columns that were added or dropped using `ALGORITHM=INSTANT` without using a `.cfg` file can result in undefined behavior.

- On Windows, InnoDB stores database, tablespace, and table names internally in lowercase. To avoid import problems on case-sensitive operating systems such as Linux and Unix, create all databases, tablespaces, and tables using lowercase names. A convenient way to ensure that names are created in lowercase is to set `lower_case_table_names` to 1 before initializing the server. (It is prohibited to start the server with a `lower_case_table_names` setting that is different from the setting used when the server was initialized.)

```
[mysqld]
lower_case_table_names=1
```

- When running `ALTER TABLE ... DISCARD PARTITION ... TABLESPACE` and `ALTER TABLE ... IMPORT PARTITION ... TABLESPACE` on subpartitioned tables, both partition and subpartition table names are permitted. When a partition name is specified, subpartitions of that partition are included in the operation.

Internals

The following information describes internals and messages written to the error log during a table import procedure.

When `ALTER TABLE ... DISCARD TABLESPACE` is run on the destination instance:

- The table is locked in X mode.
- The tablespace is detached from the table.

When `FLUSH TABLES ... FOR EXPORT` is run on the source instance:

- The table being flushed for export is locked in shared mode.

- The purge coordinator thread is stopped.
- Dirty pages are synchronized to disk.
- Table metadata is written to the binary `.cfg` file.

Expected error log messages for this operation:

```
[Note] InnoDB: Sync to disk of '"test"."t1"' started.
[NB] InnoDB: Stopping purge
[NB] InnoDB: Writing table metadata to './test/t1.cfg'
[NB] InnoDB: Table '"test"."t1"' flushed to disk
```

When `UNLOCK TABLES` is run on the source instance:

- The binary `.cfg` file is deleted.
- The shared lock on the table or tables being imported is released and the purge coordinator thread is restarted.

Expected error log messages for this operation:

```
[Note] InnoDB: Deleting the meta-data file './test/t1.cfg'
[NB] InnoDB: Resuming purge
```

When `ALTER TABLE ... IMPORT TABLESPACE` is run on the destination instance, the import algorithm performs the following operations for each tablespace being imported:

- Each tablespace page is checked for corruption.
- The space ID and log sequence numbers (LSNs) on each page are updated.
- Flags are validated and LSN updated for the header page.
- Btree pages are updated.
- The page state is set to dirty so that it is written to disk.

Expected error log messages for this operation:

```
[Note] InnoDB: Importing tablespace for table 'test/t1' that was exported
from host 'host_name'
[NB] InnoDB: Phase I - Update all pages
[NB] InnoDB: Sync to disk
[NB] InnoDB: Sync to disk - done!
[NB] InnoDB: Phase III - Flush changes to disk
[NB] InnoDB: Phase IV - Flush complete
```



Note

You may also receive a warning that a tablespace is discarded (if you discarded the tablespace for the destination table) and a message stating that statistics could not be calculated due to a missing `.ibd` file:

```
[Warning] InnoDB: Table "test"."t1" tablespace is set as discarded.
7f34d9a37700 InnoDB: cannot calculate statistics for table
"test"."t1" because the .ibd file is missing. For help, please refer to
http://dev.mysql.com/doc/refman/8.0/en/innodb-troubleshooting.html
```

15.6.1.4 Moving or Copying InnoDB Tables

This section describes techniques for moving or copying some or all InnoDB tables to a different server or instance. For example, you might move an entire MySQL instance to a larger, faster server; you might clone an entire MySQL instance to a new replica server; you might copy individual tables to another instance to develop and test an application, or to a data warehouse server to produce reports.

On Windows, [InnoDB](#) always stores database and table names internally in lowercase. To move databases in a binary format from Unix to Windows or from Windows to Unix, create all databases and tables using lowercase names. A convenient way to accomplish this is to add the following line to the `[mysqld]` section of your `my.cnf` or `my.ini` file before creating any databases or tables:

```
[mysqld]
lower_case_table_names=1
```

**Note**

It is prohibited to start the server with a `lower_case_table_names` setting that is different from the setting used when the server was initialized.

Techniques for moving or copying [InnoDB](#) tables include:

- [Importing Tables](#)
- [MySQL Enterprise Backup](#)
- [Copying Data Files \(Cold Backup Method\)](#)
- [Restoring from a Logical Backup](#)

Importing Tables

A table that resides in a file-per-table tablespace can be imported from another MySQL server instance or from a backup using the *Transportable Tablespace* feature. See [Section 15.6.1.3, “Importing InnoDB Tables”](#).

MySQL Enterprise Backup

The MySQL Enterprise Backup product lets you back up a running MySQL database with minimal disruption to operations while producing a consistent snapshot of the database. When MySQL Enterprise Backup is copying tables, reads and writes can continue. In addition, MySQL Enterprise Backup can create compressed backup files, and back up subsets of tables. In conjunction with the MySQL binary log, you can perform point-in-time recovery. MySQL Enterprise Backup is included as part of the MySQL Enterprise subscription.

For more details about MySQL Enterprise Backup, see [Section 30.2, “MySQL Enterprise Backup Overview”](#).

Copying Data Files (Cold Backup Method)

You can move an [InnoDB](#) database simply by copying all the relevant files listed under "Cold Backups" in [Section 15.18.1, “InnoDB Backup”](#).

[InnoDB](#) data and log files are binary-compatible on all platforms having the same floating-point number format. If the floating-point formats differ but you have not used `FLOAT` or `DOUBLE` data types in your tables, then the procedure is the same: simply copy the relevant files.

When you move or copy file-per-table `.ibd` files, the database directory name must be the same on the source and destination systems. The table definition stored in the [InnoDB](#) shared tablespace includes the database name. The transaction IDs and log sequence numbers stored in the tablespace files also differ between databases.

To move an `.ibd` file and the associated table from one database to another, use a `RENAME TABLE` statement:

```
RENAME TABLE db1.tbl_name TO db2.tbl_name;
```

If you have a “clean” backup of an `.ibd` file, you can restore it to the MySQL installation from which it originated as follows:

1. The table must not have been dropped or truncated since you copied the `.ibd` file, because doing so changes the table ID stored inside the tablespace.
2. Issue this `ALTER TABLE` statement to delete the current `.ibd` file:

```
ALTER TABLE tbl_name DISCARD TABLESPACE;
```

3. Copy the backup `.ibd` file to the proper database directory.
4. Issue this `ALTER TABLE` statement to tell `InnoDB` to use the new `.ibd` file for the table:

```
ALTER TABLE tbl_name IMPORT TABLESPACE;
```



Note

The `ALTER TABLE ... IMPORT TABLESPACE` feature does not enforce foreign key constraints on imported data.

In this context, a “clean” `.ibd` file backup is one for which the following requirements are satisfied:

- There are no uncommitted modifications by transactions in the `.ibd` file.
- There are no unmerged insert buffer entries in the `.ibd` file.
- Purge has removed all delete-marked index records from the `.ibd` file.
- `mysqld` has flushed all modified pages of the `.ibd` file from the buffer pool to the file.

You can make a clean backup `.ibd` file using the following method:

1. Stop all activity from the `mysqld` server and commit all transactions.
2. Wait until `SHOW ENGINE INNODB STATUS` shows that there are no active transactions in the database, and the main thread status of `InnoDB` is `Waiting for server activity`. Then you can make a copy of the `.ibd` file.

Another method for making a clean copy of an `.ibd` file is to use the MySQL Enterprise Backup product:

1. Use MySQL Enterprise Backup to back up the `InnoDB` installation.
2. Start a second `mysqld` server on the backup and let it clean up the `.ibd` files in the backup.

Restoring from a Logical Backup

You can use a utility such as `mysqldump` to perform a logical backup, which produces a set of SQL statements that can be executed to reproduce the original database object definitions and table data for transfer to another SQL server. Using this method, it does not matter whether the formats differ or if your tables contain floating-point data.

To improve the performance of this method, disable `autocommit` when importing data. Perform a commit only after importing an entire table or segment of a table.

15.6.1.5 Converting Tables from MyISAM to InnoDB

If you have `MyISAM` tables that you want to convert to `InnoDB` for better reliability and scalability, review the following guidelines and tips before converting.



Note

Partitioned `MyISAM` tables created in previous versions of MySQL are not compatible with MySQL 8.0. Such tables must be prepared prior to upgrade, either by removing the partitioning, or by converting them to `InnoDB`. See

Section 24.6.2, “Partitioning Limitations Relating to Storage Engines”, for more information.

- [Adjusting Memory Usage for MyISAM and InnoDB](#)
- [Handling Too-Long Or Too-Short Transactions](#)
- [Handling Deadlocks](#)
- [Storage Layout](#)
- [Converting an Existing Table](#)
- [Cloning the Structure of a Table](#)
- [Transferring Data](#)
- [Storage Requirements](#)
- [Defining Primary Keys](#)
- [Application Performance Considerations](#)
- [Understanding Files Associated with InnoDB Tables](#)

Adjusting Memory Usage for MyISAM and InnoDB

As you transition away from MyISAM tables, lower the value of the `key_buffer_size` configuration option to free memory no longer needed for caching results. Increase the value of the `innodb_buffer_pool_size` configuration option, which performs a similar role of allocating cache memory for InnoDB tables. The InnoDB buffer pool caches both table data and index data, speeding up lookups for queries and keeping query results in memory for reuse. For guidance regarding buffer pool size configuration, see [Section 8.12.3.1, “How MySQL Uses Memory”](#).

Handling Too-Long Or Too-Short Transactions

Because MyISAM tables do not support `transactions`, you might not have paid much attention to the `autocommit` configuration option and the `COMMIT` and `ROLLBACK` statements. These keywords are important to allow multiple sessions to read and write InnoDB tables concurrently, providing substantial scalability benefits in write-heavy workloads.

While a transaction is open, the system keeps a snapshot of the data as seen at the beginning of the transaction, which can cause substantial overhead if the system inserts, updates, and deletes millions of rows while a stray transaction keeps running. Thus, take care to avoid transactions that run for too long:

- If you are using a `mysql` session for interactive experiments, always `COMMIT` (to finalize the changes) or `ROLLBACK` (to undo the changes) when finished. Close down interactive sessions rather than leave them open for long periods, to avoid keeping transactions open for long periods by accident.
- Make sure that any error handlers in your application also `ROLLBACK` incomplete changes or `COMMIT` completed changes.
- `ROLLBACK` is a relatively expensive operation, because `INSERT`, `UPDATE`, and `DELETE` operations are written to InnoDB tables prior to the `COMMIT`, with the expectation that most changes are committed successfully and rollbacks are rare. When experimenting with large volumes of data, avoid making changes to large numbers of rows and then rolling back those changes.
- When loading large volumes of data with a sequence of `INSERT` statements, periodically `COMMIT` the results to avoid having transactions that last for hours. In typical load operations for data warehousing, if something goes wrong, you truncate the table (using `TRUNCATE TABLE`) and start over from the beginning rather than doing a `ROLLBACK`.

The preceding tips save memory and disk space that can be wasted during too-long transactions. When transactions are shorter than they should be, the problem is excessive I/O. With each `COMMIT`, MySQL makes sure each change is safely recorded to disk, which involves some I/O.

- For most operations on `InnoDB` tables, you should use the setting `autocommit=0`. From an efficiency perspective, this avoids unnecessary I/O when you issue large numbers of consecutive `INSERT`, `UPDATE`, or `DELETE` statements. From a safety perspective, this allows you to issue a `ROLLBACK` statement to recover lost or garbled data if you make a mistake on the `mysql` command line, or in an exception handler in your application.
- `autocommit=1` is suitable for `InnoDB` tables when running a sequence of queries for generating reports or analyzing statistics. In this situation, there is no I/O penalty related to `COMMIT` or `ROLLBACK`, and `InnoDB` can automatically optimize the read-only workload.
- If you make a series of related changes, finalize all the changes at once with a single `COMMIT` at the end. For example, if you insert related pieces of information into several tables, do a single `COMMIT` after making all the changes. Or if you run many consecutive `INSERT` statements, do a single `COMMIT` after all the data is loaded; if you are doing millions of `INSERT` statements, perhaps split up the huge transaction by issuing a `COMMIT` every ten thousand or hundred thousand records, so the transaction does not grow too large.
- Remember that even a `SELECT` statement opens a transaction, so after running some report or debugging queries in an interactive `mysql` session, either issue a `COMMIT` or close the `mysql` session.

For related information, see [Section 15.7.2.2, “autocommit, Commit, and Rollback”](#).

Handling Deadlocks

You might see warning messages referring to “deadlocks” in the MySQL error log, or the output of `SHOW ENGINE INNODB STATUS`. A `deadlock` is not a serious issue for `InnoDB` tables, and often does not require any corrective action. When two transactions start modifying multiple tables, accessing the tables in a different order, they can reach a state where each transaction is waiting for the other and neither can proceed. When `deadlock detection` is enabled (the default), MySQL immediately detects this condition and cancels (`rolls back`) the “smaller” transaction, allowing the other to proceed. If deadlock detection is disabled using the `innodb_deadlock_detect` configuration option, `InnoDB` relies on the `innodb_lock_wait_timeout` setting to roll back transactions in case of a deadlock.

Either way, your applications need error-handling logic to restart a transaction that is forcibly cancelled due to a deadlock. When you re-issue the same SQL statements as before, the original timing issue no longer applies. Either the other transaction has already finished and yours can proceed, or the other transaction is still in progress and your transaction waits until it finishes.

If deadlock warnings occur constantly, you might review the application code to reorder the SQL operations in a consistent way, or to shorten the transactions. You can test with the `innodb_print_all_deadlocks` option enabled to see all deadlock warnings in the MySQL error log, rather than only the last warning in the `SHOW ENGINE INNODB STATUS` output.

For more information, see [Section 15.7.5, “Deadlocks in InnoDB”](#).

Storage Layout

To get the best performance from `InnoDB` tables, you can adjust a number of parameters related to storage layout.

When you convert `MyISAM` tables that are large, frequently accessed, and hold vital data, investigate and consider the `innodb_file_per_table` and `innodb_page_size` variables, and the `ROW_FORMAT` and `KEY_BLOCK_SIZE` clauses of the `CREATE TABLE` statement.

During your initial experiments, the most important setting is `innodb_file_per_table`. When this setting is enabled, which is the default, new `InnoDB` tables are implicitly created in `file-per-table`

tablespaces. In contrast with the [InnoDB](#) system tablespace, file-per-table tablespaces allow disk space to be reclaimed by the operating system when a table is truncated or dropped. File-per-table tablespaces also support [DYNAMIC](#) and [COMPRESSED](#) row formats and associated features such as table compression, efficient off-page storage for long variable-length columns, and large index prefixes. For more information, see [Section 15.6.3.2, “File-Per-Table Tablespaces”](#).

You can also store [InnoDB](#) tables in a shared general tablespace, which support multiple tables and all row formats. For more information, see [Section 15.6.3.3, “General Tablespaces”](#).

Converting an Existing Table

To convert a non-[InnoDB](#) table to use [InnoDB](#) use `ALTER TABLE`:

```
ALTER TABLE table_name ENGINE=InnoDB;
```

Cloning the Structure of a Table

You might make an [InnoDB](#) table that is a clone of a MyISAM table, rather than using `ALTER TABLE` to perform conversion, to test the old and new table side-by-side before switching.

Create an empty [InnoDB](#) table with identical column and index definitions. Use `SHOW CREATE TABLE table_name\G` to see the full `CREATE TABLE` statement to use. Change the `ENGINE` clause to `ENGINE=INNODB`.

Transferring Data

To transfer a large volume of data into an empty [InnoDB](#) table created as shown in the previous section, insert the rows with `INSERT INTO innodb_table SELECT * FROM myisam_table ORDER BY primary_key_columns`.

You can also create the indexes for the [InnoDB](#) table after inserting the data. Historically, creating new secondary indexes was a slow operation for [InnoDB](#), but now you can create the indexes after the data is loaded with relatively little overhead from the index creation step.

If you have `UNIQUE` constraints on secondary keys, you can speed up a table import by turning off the uniqueness checks temporarily during the import operation:

```
SET unique_checks=0;
... import operation ...
SET unique_checks=1;
```

For big tables, this saves disk I/O because [InnoDB](#) can use its [change buffer](#) to write secondary index records as a batch. Be certain that the data contains no duplicate keys. `unique_checks` permits but does not require storage engines to ignore duplicate keys.

For better control over the insertion process, you can insert big tables in pieces:

```
INSERT INTO newtable SELECT * FROM oldtable
  WHERE yourkey > something AND yourkey <= somethingelse;
```

After all records are inserted, you can rename the tables.

During the conversion of big tables, increase the size of the [InnoDB](#) buffer pool to reduce disk I/O. Typically, the recommended buffer pool size is 50 to 75 percent of system memory. You can also increase the size of [InnoDB](#) log files.

Storage Requirements

If you intend to make several temporary copies of your data in [InnoDB](#) tables during the conversion process, it is recommended that you create the tables in file-per-table tablespaces so that you can reclaim the disk space when you drop the tables. When the `innodb_file_per_table` configuration option is enabled (the default), newly created [InnoDB](#) tables are implicitly created in file-per-table tablespaces.

Whether you convert the [MyISAM](#) table directly or create a cloned [InnoDB](#) table, make sure that you have sufficient disk space to hold both the old and new tables during the process. [InnoDB](#) tables require more disk space than [MyISAM](#) tables. If an [ALTER TABLE](#) operation runs out of space, it starts a rollback, and that can take hours if it is disk-bound. For inserts, [InnoDB](#) uses the insert buffer to merge secondary index records to indexes in batches. That saves a lot of disk I/O. For rollback, no such mechanism is used, and the rollback can take 30 times longer than the insertion.

In the case of a runaway rollback, if you do not have valuable data in your database, it may be advisable to kill the database process rather than wait for millions of disk I/O operations to complete. For the complete procedure, see [Section 15.21.3, “Forcing InnoDB Recovery”](#).

Defining Primary Keys

The [PRIMARY KEY](#) clause is a critical factor affecting the performance of MySQL queries and the space usage for tables and indexes. The primary key uniquely identifies a row in a table. Every row in the table should have a primary key value, and no two rows can have the same primary key value.

These are guidelines for the primary key, followed by more detailed explanations.

- Declare a [PRIMARY KEY](#) for each table. Typically, it is the most important column that you refer to in [WHERE](#) clauses when looking up a single row.
- Declare the [PRIMARY KEY](#) clause in the original [CREATE TABLE](#) statement, rather than adding it later through an [ALTER TABLE](#) statement.
- Choose the column and its data type carefully. Prefer numeric columns over character or string ones.
- Consider using an auto-increment column if there is not another stable, unique, non-null, numeric column to use.
- An auto-increment column is also a good choice if there is any doubt whether the value of the primary key column could ever change. Changing the value of a primary key column is an expensive operation, possibly involving rearranging data within the table and within each secondary index.

Consider adding a [primary key](#) to any table that does not already have one. Use the smallest practical numeric type based on the maximum projected size of the table. This can make each row slightly more compact, which can yield substantial space savings for large tables. The space savings are multiplied if the table has any [secondary indexes](#), because the primary key value is repeated in each secondary index entry. In addition to reducing data size on disk, a small primary key also lets more data fit into the [buffer pool](#), speeding up all kinds of operations and improving concurrency.

If the table already has a primary key on some longer column, such as a [VARCHAR](#), consider adding a new unsigned [AUTO_INCREMENT](#) column and switching the primary key to that, even if that column is not referenced in queries. This design change can produce substantial space savings in the secondary indexes. You can designate the former primary key columns as [UNIQUE NOT NULL](#) to enforce the same constraints as the [PRIMARY KEY](#) clause, that is, to prevent duplicate or null values across all those columns.

If you spread related information across multiple tables, typically each table uses the same column for its primary key. For example, a personnel database might have several tables, each with a primary key of employee number. A sales database might have some tables with a primary key of customer number, and other tables with a primary key of order number. Because lookups using the primary key are very fast, you can construct efficient join queries for such tables.

If you leave the [PRIMARY KEY](#) clause out entirely, MySQL creates an invisible one for you. It is a 6-byte value that might be longer than you need, thus wasting space. Because it is hidden, you cannot refer to it in queries.

Application Performance Considerations

The reliability and scalability features of [InnoDB](#) require more disk storage than equivalent [MyISAM](#) tables. You might change the column and index definitions slightly, for better space utilization, reduced

I/O and memory consumption when processing result sets, and better query optimization plans making efficient use of index lookups.

If you set up a numeric ID column for the primary key, use that value to cross-reference with related values in any other tables, particularly for [join](#) queries. For example, rather than accepting a country name as input and doing queries searching for the same name, do one lookup to determine the country ID, then do other queries (or a single join query) to look up relevant information across several tables. Rather than storing a customer or catalog item number as a string of digits, potentially using up several bytes, convert it to a numeric ID for storing and querying. A 4-byte unsigned [INT](#) column can index over 4 billion items (with the US meaning of billion: 1000 million). For the ranges of the different integer types, see [Section 11.1.2, “Integer Types \(Exact Value\) - INTEGER, INT, SMALLINT, TINYINT, MEDIUMINT, BIGINT”](#).

Understanding Files Associated with InnoDB Tables

[InnoDB](#) files require more care and planning than [MyISAM](#) files do.

- You must not delete the [ibdata files](#) that represent the [InnoDB system tablespace](#).
- Methods of moving or copying [InnoDB](#) tables to a different server are described in [Section 15.6.1.4, “Moving or Copying InnoDB Tables”](#).

15.6.1.6 AUTO_INCREMENT Handling in InnoDB

[InnoDB](#) provides a configurable locking mechanism that can significantly improve scalability and performance of SQL statements that add rows to tables with [AUTO_INCREMENT](#) columns. To use the [AUTO_INCREMENT](#) mechanism with an [InnoDB](#) table, an [AUTO_INCREMENT](#) column must be defined as the first or only column of some index such that it is possible to perform the equivalent of an indexed `SELECT MAX(ai_col)` lookup on the table to obtain the maximum column value. The index is not required to be a [PRIMARY KEY](#) or [UNIQUE](#), but to avoid duplicate values in the [AUTO_INCREMENT](#) column, those index types are recommended.

This section describes the [AUTO_INCREMENT](#) lock modes, usage implications of different [AUTO_INCREMENT](#) lock mode settings, and how [InnoDB](#) initializes the [AUTO_INCREMENT](#) counter.

- [InnoDB AUTO_INCREMENT Lock Modes](#)
- [InnoDB AUTO_INCREMENT Lock Mode Usage Implications](#)
- [InnoDB AUTO_INCREMENT Counter Initialization](#)
- [Notes](#)

InnoDB AUTO_INCREMENT Lock Modes

This section describes the [AUTO_INCREMENT](#) lock modes used to generate auto-increment values, and how each lock mode affects replication. The auto-increment lock mode is configured at startup using the [innodb_autoinc_lock_mode](#) variable.

The following terms are used in describing [innodb_autoinc_lock_mode](#) settings:

- “[INSERT-like](#)” statements

All statements that generate new rows in a table, including `INSERT`, `INSERT ... SELECT`, `REPLACE`, `REPLACE ... SELECT`, and `LOAD DATA`. Includes “simple-inserts”, “bulk-inserts”, and “mixed-mode” inserts.

- “Simple inserts”

Statements for which the number of rows to be inserted can be determined in advance (when the statement is initially processed). This includes single-row and multiple-row `INSERT` and `REPLACE` statements that do not have a nested subquery, but not `INSERT ... ON DUPLICATE KEY UPDATE`.

- “Bulk inserts”

Statements for which the number of rows to be inserted (and the number of required auto-increment values) is not known in advance. This includes `INSERT ... SELECT, REPLACE ... SELECT`, and `LOAD DATA` statements, but not plain `INSERT`. InnoDB assigns new values for the `AUTO_INCREMENT` column one at a time as each row is processed.

- “Mixed-mode inserts”

These are “simple insert” statements that specify the auto-increment value for some (but not all) of the new rows. An example follows, where `c1` is an `AUTO_INCREMENT` column of table `t1`:

```
INSERT INTO t1 (c1,c2) VALUES (1,'a'), (NULL,'b'), (5,'c'), (NULL,'d');
```

Another type of “mixed-mode insert” is `INSERT ... ON DUPLICATE KEY UPDATE`, which in the worst case is in effect an `INSERT` followed by a `UPDATE`, where the allocated value for the `AUTO_INCREMENT` column may or may not be used during the update phase.

There are three possible settings for the `innodb_autoinc_lock_mode` variable. The settings are 0, 1, or 2, for “traditional”, “consecutive”, or “interleaved” lock mode, respectively. As of MySQL 8.0, interleaved lock mode (`innodb_autoinc_lock_mode=2`) is the default setting. Prior to MySQL 8.0, consecutive lock mode is the default (`innodb_autoinc_lock_mode=1`).

The default setting of interleaved lock mode in MySQL 8.0 reflects the change from statement-based replication to row based replication as the default replication type. Statement-based replication requires the consecutive auto-increment lock mode to ensure that auto-increment values are assigned in a predictable and repeatable order for a given sequence of SQL statements, whereas row-based replication is not sensitive to the execution order of SQL statements.

- `innodb_autoinc_lock_mode = 0` (“traditional” lock mode)

The traditional lock mode provides the same behavior that existed before the `innodb_autoinc_lock_mode` variable was introduced. The traditional lock mode option is provided for backward compatibility, performance testing, and working around issues with “mixed-mode inserts”, due to possible differences in semantics.

In this lock mode, all “`INSERT-like`” statements obtain a special table-level `AUTO-INC` lock for inserts into tables with `AUTO_INCREMENT` columns. This lock is normally held to the end of the statement (not to the end of the transaction) to ensure that auto-increment values are assigned in a predictable and repeatable order for a given sequence of `INSERT` statements, and to ensure that auto-increment values assigned by any given statement are consecutive.

In the case of statement-based replication, this means that when an SQL statement is replicated on a replica server, the same values are used for the auto-increment column as on the source server. The result of execution of multiple `INSERT` statements is deterministic, and the replica reproduces the same data as on the source. If auto-increment values generated by multiple `INSERT` statements were interleaved, the result of two concurrent `INSERT` statements would be nondeterministic, and could not reliably be propagated to a replica server using statement-based replication.

To make this clear, consider an example that uses this table:

```
CREATE TABLE t1 (
  c1 INT(11) NOT NULL AUTO_INCREMENT,
  c2 VARCHAR(10) DEFAULT NULL,
  PRIMARY KEY (c1)
) ENGINE=InnoDB;
```

Suppose that there are two transactions running, each inserting rows into a table with an `AUTO_INCREMENT` column. One transaction is using an `INSERT ... SELECT` statement that inserts 1000 rows, and another is using a simple `INSERT` statement that inserts one row:

```
Tx1: INSERT INTO t1 (c2) SELECT 1000 rows from another table ...
```

```
Tx2: INSERT INTO t1 (c2) VALUES ('xxx');
```

InnoDB cannot tell in advance how many rows are retrieved from the SELECT in the INSERT statement in Tx1, and it assigns the auto-increment values one at a time as the statement proceeds. With a table-level lock, held to the end of the statement, only one INSERT statement referring to table t1 can execute at a time, and the generation of auto-increment numbers by different statements is not interleaved. The auto-increment values generated by the Tx1 INSERT ... SELECT statement are consecutive, and the (single) auto-increment value used by the INSERT statement in Tx2 is either smaller or larger than all those used for Tx1, depending on which statement executes first.

As long as the SQL statements execute in the same order when replayed from the binary log (when using statement-based replication, or in recovery scenarios), the results are the same as they were when Tx1 and Tx2 first ran. Thus, table-level locks held until the end of a statement make INSERT statements using auto-increment safe for use with statement-based replication. However, those table-level locks limit concurrency and scalability when multiple transactions are executing insert statements at the same time.

In the preceding example, if there were no table-level lock, the value of the auto-increment column used for the INSERT in Tx2 depends on precisely when the statement executes. If the INSERT of Tx2 executes while the INSERT of Tx1 is running (rather than before it starts or after it completes), the specific auto-increment values assigned by the two INSERT statements are nondeterministic, and may vary from run to run.

Under the consecutive lock mode, InnoDB can avoid using table-level AUTO-INC locks for “simple insert” statements where the number of rows is known in advance, and still preserve deterministic execution and safety for statement-based replication.

If you are not using the binary log to replay SQL statements as part of recovery or replication, the interleaved lock mode can be used to eliminate all use of table-level AUTO-INC locks for even greater concurrency and performance, at the cost of permitting gaps in auto-increment numbers assigned by a statement and potentially having the numbers assigned by concurrently executing statements interleaved.

- `innodb_autoinc_lock_mode = 1` (“consecutive” lock mode)

In this mode, “bulk inserts” use the special AUTO-INC table-level lock and hold it until the end of the statement. This applies to all INSERT ... SELECT, REPLACE ... SELECT, and LOAD DATA statements. Only one statement holding the AUTO-INC lock can execute at a time. If the source table of the bulk insert operation is different from the target table, the AUTO-INC lock on the target table is taken after a shared lock is taken on the first row selected from the source table. If the source and target of the bulk insert operation are the same table, the AUTO-INC lock is taken after shared locks are taken on all selected rows.

“Simple inserts” (for which the number of rows to be inserted is known in advance) avoid table-level AUTO-INC locks by obtaining the required number of auto-increment values under the control of a mutex (a light-weight lock) that is only held for the duration of the allocation process, *not* until the statement completes. No table-level AUTO-INC lock is used unless an AUTO-INC lock is held by another transaction. If another transaction holds an AUTO-INC lock, a “simple insert” waits for the AUTO-INC lock, as if it were a “bulk insert”.

This lock mode ensures that, in the presence of INSERT statements where the number of rows is not known in advance (and where auto-increment numbers are assigned as the statement progresses), all auto-increment values assigned by any “INSERT-like” statement are consecutive, and operations are safe for statement-based replication.

Simply put, this lock mode significantly improves scalability while being safe for use with statement-based replication. Further, as with “traditional” lock mode, auto-increment numbers assigned by any given statement are consecutive. There is *no change* in semantics compared to “traditional” mode for any statement that uses auto-increment, with one important exception.

The exception is for “mixed-mode inserts”, where the user provides explicit values for an `AUTO_INCREMENT` column for some, but not all, rows in a multiple-row “simple insert”. For such inserts, InnoDB allocates more auto-increment values than the number of rows to be inserted. However, all values automatically assigned are consecutively generated (and thus higher than) the auto-increment value generated by the most recently executed previous statement. “Excess” numbers are lost.

- `innodb_autoinc_lock_mode = 2` (“interleaved” lock mode)

In this lock mode, no “`INSERT`-like” statements use the table-level `AUTO-INC` lock, and multiple statements can execute at the same time. This is the fastest and most scalable lock mode, but it is *not safe* when using statement-based replication or recovery scenarios when SQL statements are replayed from the binary log.

In this lock mode, auto-increment values are guaranteed to be unique and monotonically increasing across all concurrently executing “`INSERT`-like” statements. However, because multiple statements can be generating numbers at the same time (that is, allocation of numbers is *interleaved* across statements), the values generated for the rows inserted by any given statement may not be consecutive.

If the only statements executing are “simple inserts” where the number of rows to be inserted is known ahead of time, there are no gaps in the numbers generated for a single statement, except for “mixed-mode inserts”. However, when “bulk inserts” are executed, there may be gaps in the auto-increment values assigned by any given statement.

InnoDB AUTO_INCREMENT Lock Mode Usage Implications

- Using auto-increment with replication

If you are using statement-based replication, set `innodb_autoinc_lock_mode` to 0 or 1 and use the same value on the source and its replicas. Auto-increment values are not ensured to be the same on the replicas as on the source if you use `innodb_autoinc_lock_mode = 2` (“interleaved”) or configurations where the source and replicas do not use the same lock mode.

If you are using row-based or mixed-format replication, all of the auto-increment lock modes are safe, since row-based replication is not sensitive to the order of execution of the SQL statements (and the mixed format uses row-based replication for any statements that are unsafe for statement-based replication).

- “Lost” auto-increment values and sequence gaps

In all lock modes (0, 1, and 2), if a transaction that generated auto-increment values rolls back, those auto-increment values are “lost”. Once a value is generated for an auto-increment column, it cannot be rolled back, whether or not the “`INSERT`-like” statement is completed, and whether or not the containing transaction is rolled back. Such lost values are not reused. Thus, there may be gaps in the values stored in an `AUTO_INCREMENT` column of a table.

- Specifying NULL or 0 for the `AUTO_INCREMENT` column

In all lock modes (0, 1, and 2), if a user specifies NULL or 0 for the `AUTO_INCREMENT` column in an `INSERT`, InnoDB treats the row as if the value was not specified and generates a new value for it.

- Assigning a negative value to the `AUTO_INCREMENT` column

In all lock modes (0, 1, and 2), the behavior of the auto-increment mechanism is undefined if you assign a negative value to the `AUTO_INCREMENT` column.

- If the `AUTO_INCREMENT` value becomes larger than the maximum integer for the specified integer type

In all lock modes (0, 1, and 2), the behavior of the auto-increment mechanism is undefined if the value becomes larger than the maximum integer that can be stored in the specified integer type.

- Gaps in auto-increment values for “bulk inserts”

With `innodb_autoinc_lock_mode` set to 0 (“traditional”) or 1 (“consecutive”), the auto-increment values generated by any given statement are consecutive, without gaps, because the table-level `AUTO-INC` lock is held until the end of the statement, and only one such statement can execute at a time.

With `innodb_autoinc_lock_mode` set to 2 (“interleaved”), there may be gaps in the auto-increment values generated by “bulk inserts,” but only if there are concurrently executing “`INSERT`-like” statements.

For lock modes 1 or 2, gaps may occur between successive statements because for bulk inserts the exact number of auto-increment values required by each statement may not be known and overestimation is possible.

- Auto-increment values assigned by “mixed-mode inserts”

Consider a “mixed-mode insert,” where a “simple insert” specifies the auto-increment value for some (but not all) resulting rows. Such a statement behaves differently in lock modes 0, 1, and 2. For example, assume `c1` is an `AUTO_INCREMENT` column of table `t1`, and that the most recent automatically generated sequence number is 100.

```
mysql> CREATE TABLE t1 (
->   c1 INT UNSIGNED NOT NULL AUTO_INCREMENT PRIMARY KEY,
->   c2 CHAR(1)
-> ) ENGINE = INNODB;
```

Now, consider the following “mixed-mode insert” statement:

```
mysql> INSERT INTO t1 (c1,c2) VALUES (1,'a'), (NULL,'b'), (5,'c'), (NULL,'d');
```

With `innodb_autoinc_lock_mode` set to 0 (“traditional”), the four new rows are:

```
mysql> SELECT c1, c2 FROM t1 ORDER BY c2;
+----+----+
| c1 | c2 |
+----+----+
| 1  | a  |
| 101 | b  |
| 5  | c  |
| 102 | d  |
+----+----+
```

The next available auto-increment value is 103 because the auto-increment values are allocated one at a time, not all at once at the beginning of statement execution. This result is true whether or not there are concurrently executing “`INSERT`-like” statements (of any type).

With `innodb_autoinc_lock_mode` set to 1 (“consecutive”), the four new rows are also:

```
mysql> SELECT c1, c2 FROM t1 ORDER BY c2;
+----+----+
| c1 | c2 |
+----+----+
| 1  | a  |
| 101 | b  |
| 5  | c  |
| 102 | d  |
+----+----+
```

```
+-----+-----+
```

However, in this case, the next available auto-increment value is 105, not 103 because four auto-increment values are allocated at the time the statement is processed, but only two are used. This result is true whether or not there are concurrently executing “[INSERT](#)-like” statements (of any type).

With `innodb_autoinc_lock_mode` set to 2 (“interleaved”), the four new rows are:

```
mysql> SELECT c1, c2 FROM t1 ORDER BY c2;
+---+---+
| c1 | c2 |
+---+---+
| 1  | a  |
| x  | b  |
| 5  | c  |
| y  | d  |
+---+---+
```

The values of `x` and `y` are unique and larger than any previously generated rows. However, the specific values of `x` and `y` depend on the number of auto-increment values generated by concurrently executing statements.

Finally, consider the following statement, issued when the most-recently generated sequence number is 100:

```
mysql> INSERT INTO t1 (c1,c2) VALUES (1,'a'), (NULL,'b'), (101,'c'), (NULL,'d');
```

With any `innodb_autoinc_lock_mode` setting, this statement generates a duplicate-key error 23000 (`Can't write; duplicate key in table`) because 101 is allocated for the row `(NULL, 'b')` and insertion of the row `(101, 'c')` fails.

- Modifying `AUTO_INCREMENT` column values in the middle of a sequence of [INSERT](#) statements

In MySQL 5.7 and earlier, modifying an `AUTO_INCREMENT` column value in the middle of a sequence of [INSERT](#) statements could lead to “Duplicate entry” errors. For example, if you performed an `UPDATE` operation that changed an `AUTO_INCREMENT` column value to a value larger than the current maximum auto-increment value, subsequent [INSERT](#) operations that did not specify an unused auto-increment value could encounter “Duplicate entry” errors. In MySQL 8.0 and later, if you modify an `AUTO_INCREMENT` column value to a value larger than the current maximum auto-increment value, the new value is persisted, and subsequent [INSERT](#) operations allocate auto-increment values starting from the new, larger value. This behavior is demonstrated in the following example.

```
mysql> CREATE TABLE t1 (
    -> c1 INT NOT NULL AUTO_INCREMENT,
    -> PRIMARY KEY (c1)
    -> ) ENGINE = InnoDB;

mysql> INSERT INTO t1 VALUES(0), (0), (3);

mysql> SELECT c1 FROM t1;
+---+
| c1 |
+---+
| 1  |
| 2  |
| 3  |
+---+

mysql> UPDATE t1 SET c1 = 4 WHERE c1 = 1;

mysql> SELECT c1 FROM t1;
+---+
| c1 |
+---+
| 2  |
+---+
```

```

|   3 |
|   4 |
+----+
mysql> INSERT INTO t1 VALUES(0);

mysql> SELECT c1 FROM t1;
+---+
| c1 |
+---+
| 2 |
| 3 |
| 4 |
| 5 |
+---+

```

InnoDB AUTO_INCREMENT Counter Initialization

This section describes how InnoDB initializes AUTO_INCREMENT counters.

If you specify an AUTO_INCREMENT column for an InnoDB table, the in-memory table object contains a special counter called the auto-increment counter that is used when assigning new values for the column.

In MySQL 5.7 and earlier, the auto-increment counter is stored in main memory, not on disk. To initialize an auto-increment counter after a server restart, InnoDB would execute the equivalent of the following statement on the first insert into a table containing an AUTO_INCREMENT column.

```
SELECT MAX(ai_col) FROM table_name FOR UPDATE;
```

In MySQL 8.0, this behavior is changed. The current maximum auto-increment counter value is written to the redo log each time it changes and saved to the data dictionary on each checkpoint. These changes make the current maximum auto-increment counter value persistent across server restarts.

On a server restart following a normal shutdown, InnoDB initializes the in-memory auto-increment counter using the current maximum auto-increment value stored in the data dictionary.

On a server restart during crash recovery, InnoDB initializes the in-memory auto-increment counter using the current maximum auto-increment value stored in the data dictionary and scans the redo log for auto-increment counter values written since the last checkpoint. If a redo-logged value is greater than the in-memory counter value, the redo-logged value is applied. However, in the case of an unexpected server exit, reuse of a previously allocated auto-increment value cannot be guaranteed. Each time the current maximum auto-increment value is changed due to an INSERT or UPDATE operation, the new value is written to the redo log, but if the unexpected exit occurs before the redo log is flushed to disk, the previously allocated value could be reused when the auto-increment counter is initialized after the server is restarted.

The only circumstance in which InnoDB uses the equivalent of a `SELECT MAX(ai_col) FROM table_name FOR UPDATE` statement to initialize an auto-increment counter is when importing a table without a .cfg metadata file. Otherwise, the current maximum auto-increment counter value is read from the .cfg metadata file if present. Aside from counter value initialization, the equivalent of a `SELECT MAX(ai_col) FROM table_name` statement is used to determine the current maximum auto-increment counter value of the table when attempting to set the counter value to one that is smaller than or equal to the persisted counter value using an `ALTER TABLE ... AUTO_INCREMENT = N` statement. For example, you might try to set the counter value to a lesser value after deleting some records. In this case, the table must be searched to ensure that the new counter value is not less than or equal to the actual current maximum counter value.

In MySQL 5.7 and earlier, a server restart cancels the effect of the `AUTO_INCREMENT = N` table option, which may be used in a CREATE TABLE or ALTER TABLE statement to set an initial counter value or alter the existing counter value, respectively. In MySQL 8.0, a server restart does not cancel the effect of the `AUTO_INCREMENT = N` table option. If you initialize the auto-increment counter to

a specific value, or if you alter the auto-increment counter value to a larger value, the new value is persisted across server restarts.



Note

`ALTER TABLE ... AUTO_INCREMENT = N` can only change the auto-increment counter value to a value larger than the current maximum.

In MySQL 5.7 and earlier, a server restart immediately following a `ROLLBACK` operation could result in the reuse of auto-increment values that were previously allocated to the rolled-back transaction, effectively rolling back the current maximum auto-increment value. In MySQL 8.0, the current maximum auto-increment value is persisted, preventing the reuse of previously allocated values.

If a `SHOW TABLE STATUS` statement examines a table before the auto-increment counter is initialized, `InnoDB` opens the table and initializes the counter value using the current maximum auto-increment value that is stored in the data dictionary. The value is then stored in memory for use by later inserts or updates. Initialization of the counter value uses a normal exclusive-locking read on the table which lasts to the end of the transaction. `InnoDB` follows the same procedure when initializing the auto-increment counter for a newly created table that has a user-specified auto-increment value greater than 0.

After the auto-increment counter is initialized, if you do not explicitly specify an auto-increment value when inserting a row, `InnoDB` implicitly increments the counter and assigns the new value to the column. If you insert a row that explicitly specifies an auto-increment column value, and the value is greater than the current maximum counter value, the counter is set to the specified value.

`InnoDB` uses the in-memory auto-increment counter as long as the server runs. When the server is stopped and restarted, `InnoDB` reinitializes the auto-increment counter, as described earlier.

The `auto_increment_offset` variable determines the starting point for the `AUTO_INCREMENT` column value. The default setting is 1.

The `auto_increment_increment` variable controls the interval between successive column values. The default setting is 1.

Notes

When an `AUTO_INCREMENT` integer column runs out of values, a subsequent `INSERT` operation returns a duplicate-key error. This is general MySQL behavior.

15.6.2 Indexes

This section covers topics related to `InnoDB` indexes.

15.6.2.1 Clustered and Secondary Indexes

Each `InnoDB` table has a special index called the clustered index that stores row data. Typically, the clustered index is synonymous with the primary key. To get the best performance from queries, inserts, and other database operations, it is important to understand how `InnoDB` uses the clustered index to optimize the common lookup and DML operations.

- When you define a `PRIMARY KEY` on a table, `InnoDB` uses it as the clustered index. A primary key should be defined for each table. If there is no logical unique and non-null column or set of columns to use as the primary key, add an auto-increment column. Auto-increment column values are unique and are added automatically as new rows are inserted.
- If you do not define a `PRIMARY KEY` for a table, `InnoDB` uses the first `UNIQUE` index with all key columns defined as `NOT NULL` as the clustered index.
- If a table has no `PRIMARY KEY` or suitable `UNIQUE` index, `InnoDB` generates a hidden clustered index named `GEN_CLUST_INDEX` on a synthetic column that contains row ID values. The rows are

ordered by the row ID that [InnoDB](#) assigns. The row ID is a 6-byte field that increases monotonically as new rows are inserted. Thus, the rows ordered by the row ID are physically in order of insertion.

How the Clustered Index Speeds Up Queries

Accessing a row through the clustered index is fast because the index search leads directly to the page that contains the row data. If a table is large, the clustered index architecture often saves a disk I/O operation when compared to storage organizations that store row data using a different page from the index record.

How Secondary Indexes Relate to the Clustered Index

Indexes other than the clustered index are known as secondary indexes. In [InnoDB](#), each record in a secondary index contains the primary key columns for the row, as well as the columns specified for the secondary index. [InnoDB](#) uses this primary key value to search for the row in the clustered index.

If the primary key is long, the secondary indexes use more space, so it is advantageous to have a short primary key.

For guidelines to take advantage of [InnoDB](#) clustered and secondary indexes, see [Section 8.3, “Optimization and Indexes”](#).

15.6.2.2 The Physical Structure of an InnoDB Index

With the exception of spatial indexes, [InnoDB](#) indexes are [B-tree](#) data structures. Spatial indexes use [R-trees](#), which are specialized data structures for indexing multi-dimensional data. Index records are stored in the leaf pages of their B-tree or R-tree data structure. The default size of an index page is 16KB. The page size is determined by the `innodb_page_size` setting when the MySQL instance is initialized. See [Section 15.8.1, “InnoDB Startup Configuration”](#).

When new records are inserted into an [InnoDB](#) [clustered index](#), [InnoDB](#) tries to leave 1/16 of the page free for future insertions and updates of the index records. If index records are inserted in a sequential order (ascending or descending), the resulting index pages are about 15/16 full. If records are inserted in a random order, the pages are from 1/2 to 15/16 full.

[InnoDB](#) performs a bulk load when creating or rebuilding B-tree indexes. This method of index creation is known as a sorted index build. The `innodb_fill_factor` variable defines the percentage of space on each B-tree page that is filled during a sorted index build, with the remaining space reserved for future index growth. Sorted index builds are not supported for spatial indexes. For more information, see [Section 15.6.2.3, “Sorted Index Builds”](#). An `innodb_fill_factor` setting of 100 leaves 1/16 of the space in clustered index pages free for future index growth.

If the fill factor of an [InnoDB](#) index page drops below the `MERGE_THRESHOLD`, which is 50% by default if not specified, [InnoDB](#) tries to contract the index tree to free the page. The `MERGE_THRESHOLD` setting applies to both B-tree and R-tree indexes. For more information, see [Section 15.8.11, “Configuring the Merge Threshold for Index Pages”](#).

15.6.2.3 Sorted Index Builds

[InnoDB](#) performs a bulk load instead of inserting one index record at a time when creating or rebuilding indexes. This method of index creation is also known as a sorted index build. Sorted index builds are not supported for spatial indexes.

There are three phases to an index build. In the first phase, the [clustered index](#) is scanned, and index entries are generated and added to the sort buffer. When the [sort buffer](#) becomes full, entries are sorted and written out to a temporary intermediate file. This process is also known as a “run”. In the second phase, with one or more runs written to the temporary intermediate file, a merge sort is performed on all entries in the file. In the third and final phase, the sorted entries are inserted into the B-tree.

Prior to the introduction of sorted index builds, index entries were inserted into the B-tree one record at a time using insert APIs. This method involved opening a B-tree [cursor](#) to find the insert position and then inserting entries into a B-tree page using an [optimistic](#) insert. If an insert failed due to a page being full, a [pessimistic](#) insert would be performed, which involves opening a B-tree cursor and splitting and merging B-tree nodes as necessary to find space for the entry. The drawbacks of this “top-down” method of building an index are the cost of searching for an insert position and the constant splitting and merging of B-tree nodes.

Sorted index builds use a “bottom-up” approach to building an index. With this approach, a reference to the right-most leaf page is held at all levels of the B-tree. The right-most leaf page at the necessary B-tree depth is allocated and entries are inserted according to their sorted order. Once a leaf page is full, a node pointer is appended to the parent page and a sibling leaf page is allocated for the next insert. This process continues until all entries are inserted, which may result in inserts up to the root level. When a sibling page is allocated, the reference to the previously pinned leaf page is released, and the newly allocated leaf page becomes the right-most leaf page and new default insert location.

Reserving B-tree Page Space for Future Index Growth

To set aside space for future index growth, you can use the `innodb_fill_factor` variable to reserve a percentage of B-tree page space. For example, setting `innodb_fill_factor` to 80 reserves 20 percent of the space in B-tree pages during a sorted index build. This setting applies to both B-tree leaf and non-leaf pages. It does not apply to external pages used for `TEXT` or `BLOB` entries. The amount of space that is reserved may not be exactly as configured, as the `innodb_fill_factor` value is interpreted as a hint rather than a hard limit.

Sorted Index Builds and Full-Text Index Support

Sorted index builds are supported for [fulltext indexes](#). Previously, SQL was used to insert entries into a fulltext index.

Sorted Index Builds and Compressed Tables

For [compressed tables](#), the previous index creation method appended entries to both compressed and uncompressed pages. When the modification log (representing free space on the compressed page) became full, the compressed page would be recompressed. If compression failed due to a lack of space, the page would be split. With sorted index builds, entries are only appended to uncompressed pages. When an uncompressed page becomes full, it is compressed. Adaptive padding is used to ensure that compression succeeds in most cases, but if compression fails, the page is split and compression is attempted again. This process continues until compression is successful. For more information about compression of B-Tree pages, see [Section 15.9.1.5, “How Compression Works for InnoDB Tables”](#).

Sorted Index Builds and Redo Logging

[Redo logging](#) is disabled during a sorted index build. Instead, there is a [checkpoint](#) to ensure that the index build can withstand an unexpected exit or failure. The checkpoint forces a write of all dirty pages to disk. During a sorted index build, the [page cleaner](#) thread is signaled periodically to flush [dirty pages](#) to ensure that the checkpoint operation can be processed quickly. Normally, the page cleaner thread flushes dirty pages when the number of clean pages falls below a set threshold. For sorted index builds, dirty pages are flushed promptly to reduce checkpoint overhead and to parallelize I/O and CPU activity.

Sorted Index Builds and Optimizer Statistics

Sorted index builds may result in [optimizer](#) statistics that differ from those generated by the previous method of index creation. The difference in statistics, which is not expected to affect workload performance, is due to the different algorithm used to populate the index.

15.6.2.4 InnoDB Full-Text Indexes

Full-text indexes are created on text-based columns (`CHAR`, `VARCHAR`, or `TEXT` columns) to speed up queries and DML operations on data contained within those columns.

A full-text index is defined as part of a `CREATE TABLE` statement or added to an existing table using `ALTER TABLE` or `CREATE INDEX`.

Full-text search is performed using `MATCH() ... AGAINST` syntax. For usage information, see [Section 12.10, “Full-Text Search Functions”](#).

`InnoDB` full-text indexes are described under the following topics in this section:

- [InnoDB Full-Text Index Design](#)
- [InnoDB Full-Text Index Tables](#)
- [InnoDB Full-Text Index Cache](#)
- [InnoDB Full-Text Index DOC_ID and FTS_DOC_ID Column](#)
- [InnoDB Full-Text Index Deletion Handling](#)
- [InnoDB Full-Text Index Transaction Handling](#)
- [Monitoring InnoDB Full-Text Indexes](#)

InnoDB Full-Text Index Design

`InnoDB` full-text indexes have an inverted index design. Inverted indexes store a list of words, and for each word, a list of documents that the word appears in. To support proximity search, position information for each word is also stored, as a byte offset.

InnoDB Full-Text Index Tables

When an `InnoDB` full-text index is created, a set of index tables is created, as shown in the following example:

```
mysql> CREATE TABLE opening_lines (
    id INT UNSIGNED AUTO_INCREMENT NOT NULL PRIMARY KEY,
    opening_line TEXT(500),
    author VARCHAR(200),
    title VARCHAR(200),
    FULLTEXT idx (opening_line)
) ENGINE=InnoDB;

mysql> SELECT table_id, name, space from INFORMATION_SCHEMA.INNODB_TABLES
WHERE name LIKE 'test/%';
+-----+-----+-----+
| table_id | name | space |
+-----+-----+-----+
| 333 | test/fts_0000000000000000147_0000000000000001c9_index_1 | 289 |
| 334 | test/fts_0000000000000000147_0000000000000001c9_index_2 | 290 |
| 335 | test/fts_0000000000000000147_0000000000000001c9_index_3 | 291 |
| 336 | test/fts_0000000000000000147_0000000000000001c9_index_4 | 292 |
| 337 | test/fts_0000000000000000147_0000000000000001c9_index_5 | 293 |
| 338 | test/fts_0000000000000000147_0000000000000001c9_index_6 | 294 |
| 330 | test/fts_0000000000000000147_being_deleted | 286 |
| 331 | test/fts_0000000000000000147_being_deleted_cache | 287 |
| 332 | test/fts_0000000000000000147_config | 288 |
| 328 | test/fts_0000000000000000147_deleted | 284 |
| 329 | test/fts_0000000000000000147_deleted_cache | 285 |
| 327 | test/opening_lines | 283 |
+-----+-----+-----+
```

The first six index tables comprise the inverted index and are referred to as auxiliary index tables. When incoming documents are tokenized, the individual words (also referred to as “tokens”) are

inserted into the index tables along with position information and an associated `DOC_ID`. The words are fully sorted and partitioned among the six index tables based on the character set sort weight of the word's first character.

The inverted index is partitioned into six auxiliary index tables to support parallel index creation. By default, two threads tokenize, sort, and insert words and associated data into the index tables. The number of threads that perform this work is configurable using the `innodb_ft_sort_pll_degree` variable. Consider increasing the number of threads when creating full-text indexes on large tables.

Auxiliary index table names are prefixed with `fts_` and postfixed with `index_#`. Each auxiliary index table is associated with the indexed table by a hex value in the auxiliary index table name that matches the `table_id` of the indexed table. For example, the `table_id` of the `test/opening_lines` table is `327`, for which the hex value is `0x147`. As shown in the preceding example, the “`147`” hex value appears in the names of auxiliary index tables that are associated with the `test/opening_lines` table.

A hex value representing the `index_id` of the full-text index also appears in auxiliary index table names. For example, in the auxiliary table name `test/fts_000000000000147_00000000000001c9_index_1`, the hex value `1c9` has a decimal value of `457`. The index defined on the `opening_lines` table (`idx`) can be identified by querying the Information Schema `INNODB_INDEXES` table for this value (`457`).

```
mysql> SELECT index_id, name, table_id, space FROM INFORMATION_SCHEMA.INNODB_INDEXES
      WHERE index_id=457;
+-----+-----+-----+
| index_id | name   | table_id | space |
+-----+-----+-----+
|     457  | idx    |      327 |    283 |
+-----+-----+-----+
```

Index tables are stored in their own tablespace if the primary table is created in a [file-per-table](#) tablespace. Otherwise, index tables are stored in the tablespace where the indexed table resides.

The other index tables shown in the preceding example are referred to as common index tables and are used for deletion handling and storing the internal state of full-text indexes. Unlike the inverted index tables, which are created for each full-text index, this set of tables is common to all full-text indexes created on a particular table.

Common index tables are retained even if full-text indexes are dropped. When a full-text index is dropped, the `FTS_DOC_ID` column that was created for the index is retained, as removing the `FTS_DOC_ID` column would require rebuilding the previously indexed table. Common index tables are required to manage the `FTS_DOC_ID` column.

- `fts_*_deleted` and `fts_*_deleted_cache`

Contain the document IDs (DOC_ID) for documents that are deleted but whose data is not yet removed from the full-text index. The `fts_*_deleted_cache` is the in-memory version of the `fts_*_deleted` table.

- `fts_*_being_deleted` and `fts_*_being_deleted_cache`

Contain the document IDs (DOC_ID) for documents that are deleted and whose data is currently in the process of being removed from the full-text index. The `fts_*_being_deleted_cache` table is the in-memory version of the `fts_*_being_deleted` table.

- `fts_*_config`

Stores information about the internal state of the full-text index. Most importantly, it stores the `FTS_SYNCED_DOC_ID`, which identifies documents that have been parsed and flushed to disk. In case of crash recovery, `FTS_SYNCED_DOC_ID` values are used to identify documents that have not been flushed to disk so that the documents can be re-parsed and added back to the full-text index cache. To view the data in this table, query the Information Schema `INNODB_FT_CONFIG` table.

InnoDB Full-Text Index Cache

When a document is inserted, it is tokenized, and the individual words and associated data are inserted into the full-text index. This process, even for small documents, can result in numerous small insertions into the auxiliary index tables, making concurrent access to these tables a point of contention. To avoid this problem, InnoDB uses a full-text index cache to temporarily cache index table insertions for recently inserted rows. This in-memory cache structure holds insertions until the cache is full and then batch flushes them to disk (to the auxiliary index tables). You can query the Information Schema `INNODB_FT_INDEX_CACHE` table to view tokenized data for recently inserted rows.

The caching and batch flushing behavior avoids frequent updates to auxiliary index tables, which could result in concurrent access issues during busy insert and update times. The batching technique also avoids multiple insertions for the same word, and minimizes duplicate entries. Instead of flushing each word individually, insertions for the same word are merged and flushed to disk as a single entry, improving insertion efficiency while keeping auxiliary index tables as small as possible.

The `innodb_ft_cache_size` variable is used to configure the full-text index cache size (on a per-table basis), which affects how often the full-text index cache is flushed. You can also define a global full-text index cache size limit for all tables in a given instance using the `innodb_ft_total_cache_size` variable.

The full-text index cache stores the same information as auxiliary index tables. However, the full-text index cache only caches tokenized data for recently inserted rows. The data that is already flushed to disk (to the auxiliary index tables) is not brought back into the full-text index cache when queried. The data in auxiliary index tables is queried directly, and results from the auxiliary index tables are merged with results from the full-text index cache before being returned.

InnoDB Full-Text Index DOC_ID and FTS_DOC_ID Column

InnoDB uses a unique document identifier referred to as the `DOC_ID` to map words in the full-text index to document records where the word appears. The mapping requires an `FTS_DOC_ID` column on the indexed table. If an `FTS_DOC_ID` column is not defined, InnoDB automatically adds a hidden `FTS_DOC_ID` column when the full-text index is created. The following example demonstrates this behavior.

The following table definition does not include an `FTS_DOC_ID` column:

```
mysql> CREATE TABLE opening_lines (
    id INT UNSIGNED AUTO_INCREMENT NOT NULL PRIMARY KEY,
    opening_line TEXT(500),
    author VARCHAR(200),
    title VARCHAR(200)
) ENGINE=InnoDB;
```

When you create a full-text index on the table using `CREATE FULLTEXT INDEX` syntax, a warning is returned which reports that InnoDB is rebuilding the table to add the `FTS_DOC_ID` column.

```
mysql> CREATE FULLTEXT INDEX idx ON opening_lines(opening_line);
Query OK, 0 rows affected, 1 warning (0.19 sec)
Records: 0  Duplicates: 0  Warnings: 1

mysql> SHOW WARNINGS;
+-----+-----+
| Level | Code | Message |
+-----+-----+
| Warning | 124 | InnoDB rebuilding table to add column FTS_DOC_ID |
+-----+-----+
```

The same warning is returned when using `ALTER TABLE` to add a full-text index to a table that does not have an `FTS_DOC_ID` column. If you create a full-text index at `CREATE TABLE` time and do not specify an `FTS_DOC_ID` column, InnoDB adds a hidden `FTS_DOC_ID` column, without warning.

Defining an `FTS_DOC_ID` column at `CREATE TABLE` time is less expensive than creating a full-text index on a table that is already loaded with data. If an `FTS_DOC_ID` column is defined on a table prior

to loading data, the table and its indexes do not have to be rebuilt to add the new column. If you are not concerned with `CREATE FULLTEXT INDEX` performance, leave out the `FTS_DOC_ID` column to have InnoDB create it for you. InnoDB creates a hidden `FTS_DOC_ID` column along with a unique index (`FTS_DOC_ID_INDEX`) on the `FTS_DOC_ID` column. If you want to create your own `FTS_DOC_ID` column, the column must be defined as `BIGINT UNSIGNED NOT NULL` and named `FTS_DOC_ID` (all uppercase), as in the following example:



Note

The `FTS_DOC_ID` column does not need to be defined as an `AUTO_INCREMENT` column, but doing so could make loading data easier.

```
mysql> CREATE TABLE opening_lines (
    FTS_DOC_ID BIGINT UNSIGNED AUTO_INCREMENT NOT NULL PRIMARY KEY,
    opening_line TEXT(500),
    author VARCHAR(200),
    title VARCHAR(200)
) ENGINE=InnoDB;
```

If you choose to define the `FTS_DOC_ID` column yourself, you are responsible for managing the column to avoid empty or duplicate values. `FTS_DOC_ID` values cannot be reused, which means `FTS_DOC_ID` values must be ever increasing.

Optionally, you can create the required unique `FTS_DOC_ID_INDEX` (all uppercase) on the `FTS_DOC_ID` column.

```
mysql> CREATE UNIQUE INDEX FTS_DOC_ID_INDEX ON opening_lines(FTS_DOC_ID);
```

If you do not create the `FTS_DOC_ID_INDEX`, InnoDB creates it automatically.



Note

`FTS_DOC_ID_INDEX` cannot be defined as a descending index because the InnoDB SQL parser does not use descending indexes.

The permitted gap between the largest used `FTS_DOC_ID` value and new `FTS_DOC_ID` value is 65535.

To avoid rebuilding the table, the `FTS_DOC_ID` column is retained when dropping a full-text index.

InnoDB Full-Text Index Deletion Handling

Deleting a record that has a full-text index column could result in numerous small deletions in the auxiliary index tables, making concurrent access to these tables a point of contention. To avoid this problem, the `DOC_ID` of a deleted document is logged in a special `FTS_*_DELETED` table whenever a record is deleted from an indexed table, and the indexed record remains in the full-text index. Before returning query results, information in the `FTS_*_DELETED` table is used to filter out deleted `DOC_IDS`. The benefit of this design is that deletions are fast and inexpensive. The drawback is that the size of the index is not immediately reduced after deleting records. To remove full-text index entries for deleted records, run `OPTIMIZE TABLE` on the indexed table with `innodb_optimize_fulltext_only=ON` to rebuild the full-text index. For more information, see [Optimizing InnoDB Full-Text Indexes](#).

InnoDB Full-Text Index Transaction Handling

InnoDB full-text indexes have special transaction handling characteristics due its caching and batch processing behavior. Specifically, updates and insertions on a full-text index are processed at transaction commit time, which means that a full-text search can only see committed data. The following example demonstrates this behavior. The full-text search only returns a result after the inserted lines are committed.

```
mysql> CREATE TABLE opening_lines (
    id INT UNSIGNED AUTO_INCREMENT NOT NULL PRIMARY KEY,
```

```

opening_line TEXT(500),
author VARCHAR(200),
title VARCHAR(200),
FULLTEXT idx (opening_line)
) ENGINE=InnoDB;

mysql> BEGIN;

mysql> INSERT INTO opening_lines(opening_line,author,title) VALUES
('Call me Ishmael.', 'Herman Melville', 'Moby-Dick'),
('A screaming comes across the sky.', 'Thomas Pynchon', 'Gravity\'s Rainbow'),
('I am an invisible man.', 'Ralph Ellison', 'Invisible Man'),
('Where now? Who now? When now?', 'Samuel Beckett', 'The Unnamable'),
('It was love at first sight.', 'Joseph Heller', 'Catch-22'),
('All this happened, more or less.', 'Kurt Vonnegut', 'Slaughterhouse-Five'),
('Mrs. Dalloway said she would buy the flowers herself.', 'Virginia Woolf', 'Mrs. Dalloway'),
('It was a pleasure to burn.', 'Ray Bradbury', 'Fahrenheit 451');

mysql> SELECT COUNT(*) FROM opening_lines WHERE MATCH(opening_line) AGAINST('Ishmael');

+-----+
| COUNT(*) |
+-----+
|      0 |
+-----+

mysql> COMMIT;

mysql> SELECT COUNT(*) FROM opening_lines WHERE MATCH(opening_line) AGAINST('Ishmael');

+-----+
| COUNT(*) |
+-----+
|      1 |
+-----+

```

Monitoring InnoDB Full-Text Indexes

You can monitor and examine the special text-processing aspects of [InnoDB](#) full-text indexes by querying the following [INFORMATION_SCHEMA](#) tables:

- [INNODB_FT_CONFIG](#)
- [INNODB_FT_INDEX_TABLE](#)
- [INNODB_FT_INDEX_CACHE](#)
- [INNODB_FT_DEFAULT_STOPWORD](#)
- [INNODB_FT_DELETED](#)
- [INNODB_FT_BEING_DELETED](#)

You can also view basic information for full-text indexes and tables by querying [INNODB_INDEXES](#) and [INNODB_TABLES](#).

For more information, see [Section 15.15.4, “InnoDB INFORMATION_SCHEMA FULLTEXT Index Tables”](#).

15.6.3 Tablespaces

This section covers topics related to [InnoDB](#) tablespaces.

15.6.3.1 The System Tablespace

The system tablespace is the storage area for the change buffer. It may also contain table and index data if tables are created in the system tablespace rather than file-per-table or general tablespaces. In previous MySQL versions, the system tablespace contained the [InnoDB](#) data dictionary. In MySQL

8.0, InnoDB stores metadata in the MySQL data dictionary. See [Chapter 14, MySQL Data Dictionary](#). In previous MySQL releases, the system tablespace also contained the doublewrite buffer storage area. This storage area resides in separate doublewrite files as of MySQL 8.0.20. See [Section 15.6.4, “Doublewrite Buffer”](#).

The system tablespace can have one or more data files. By default, a single system tablespace data file, named `ibdata1`, is created in the data directory. The size and number of system tablespace data files is defined by the `innodb_data_file_path` startup option. For configuration information, see [System Tablespace Data File Configuration](#).

Additional information about the system tablespace is provided under the following topics in the section:

- [Resizing the System Tablespace](#)
- [Using Raw Disk Partitions for the System Tablespace](#)

Resizing the System Tablespace

This section describes how to increase or decrease the size of the system tablespace.

Increasing the Size of the System Tablespace

The easiest way to increase the size of the system tablespace is to configure it to be auto-extending. To do so, specify the `autoextend` attribute for the last data file in the `innodb_data_file_path` setting, and restart the server. For example:

```
innodb_data_file_path=ibdata1:10M:autoextend
```

When the `autoextend` attribute is specified, the data file automatically increases in size by 8MB increments as space is required. The `innodb_autoextend_increment` variable controls the increment size.

You can also increase system tablespace size by adding another data file. To do so:

1. Stop the MySQL server.
2. If the last data file in the `innodb_data_file_path` setting is defined with the `autoextend` attribute, remove it, and modify the size attribute to reflect the current data file size. To determine the appropriate data file size to specify, check your file system for the file size, and round that value down to the closest MB value, where a MB is equal to 1024×1024 bytes.
3. Append a new data file to the `innodb_data_file_path` setting, optionally specifying the `autoextend` attribute. The `autoextend` attribute can be specified only for the last data file in the `innodb_data_file_path` setting.
4. Start the MySQL server.

For example, this tablespace has one auto-extending data file:

```
innodb_data_home_dir =
innodb_data_file_path = /ibdata/ibdata1:10M:autoextend
```

Suppose that the data file has grown to 988MB over time. This is the `innodb_data_file_path` setting after modifying the size attribute to reflect the current data file size, and after specifying a new 50MB auto-extending data file:

```
innodb_data_home_dir =
innodb_data_file_path = /ibdata/ibdata1:988M;/disk2/ibdata2:50M:autoextend
```

When adding a new data file, do not specify an existing file name. InnoDB creates and initializes the new data file when you start the server.

**Note**

You cannot increase the size of an existing system tablespace data file by changing its size attribute. For example, changing the `innodb_data_file_path` setting from `ibdata1:10M:autoextend` to `ibdata1:12M:autoextend` produces the following error when starting the server:

```
[ERROR] [MY-012263] [InnoDB] The Auto-extending innodb_system data file './ibdata1' is of a different size 640 pages (rounded down to MB) than specified in the .cnf file: initial 768 pages, max 0 (relevant if non-zero) pages!
```

The error indicates that the existing data file size (expressed in `InnoDB` pages) is different from the data file size specified in the configuration file. If you encounter this error, restore the previous `innodb_data_file_path` setting, and refer to the system tablespace resizing instructions.

Decreasing the Size of the InnoDB System Tablespace

Decreasing the size of an existing system tablespace is not supported. The only option to achieve a smaller system tablespace is to restore your data from a backup to a new MySQL instance created with the desired system tablespace size configuration.

For information about creating backups, see [Section 15.18.1, “InnoDB Backup”](#).

For information about configuring data files for a new system tablespace. See [System Tablespace Data File Configuration](#).

To avoid a large system tablespace, consider using file-per-table tablespaces or general tablespaces for your data. File-per-table tablespaces are the default tablespace type and are used implicitly when creating an `InnoDB` table. Unlike the system tablespace, file-per-table tablespaces return disk space to the operating system when they are truncated or dropped. For more information, see [Section 15.6.3.2, “File-Per-Table Tablespaces”](#). General tablespaces are multi-table tablespaces that can also be used as an alternative to the system tablespace. See [Section 15.6.3.3, “General Tablespaces”](#).

Using Raw Disk Partitions for the System Tablespace

Raw disk partitions can be used as system tablespace data files. This technique enables nonbuffered I/O on Windows and some Linux and Unix systems without file system overhead. Perform tests with and without raw partitions to verify whether they improve performance on your system.

When using a raw disk partition, ensure that the user ID that runs the MySQL server has read and write privileges for that partition. For example, if running the server as the `mysql` user, the partition must be readable and writeable by `mysql`. If running the server with the `--memlock` option, the server must be run as `root`, so the partition must be readable and writeable by `root`.

The procedures described below involve option file modification. For additional information, see [Section 4.2.2.2, “Using Option Files”](#).

Allocating a Raw Disk Partition on Linux and Unix Systems

- When creating a new data file, specify the keyword `newraw` immediately after the data file size for the `innodb_data_file_path` option. The partition must be at least as large as the size that you specify. Note that 1MB in `InnoDB` is 1024×1024 bytes, whereas 1MB in disk specifications usually means 1,000,000 bytes.

```
[mysqld]
innodb_data_home_dir=
innodb_data_file_path=/dev/hdd1:3Gnewraw;/dev/hdd2:2Gnewraw
```

- Restart the server. `InnoDB` notices the `newraw` keyword and initializes the new partition. However, do not create or change any `InnoDB` tables yet. Otherwise, when you next restart the server,

`InnoDB` reinitializes the partition and your changes are lost. (As a safety measure `InnoDB` prevents users from modifying data when any partition with `newraw` is specified.)

3. After `InnoDB` has initialized the new partition, stop the server, change `newraw` in the data file specification to `raw`:

```
[mysqld]
innodb_data_home_dir=
innodb_data_file_path=/dev/hdd1:3Graw;/dev/hdd2:2Graw
```

4. Restart the server. `InnoDB` now permits changes to be made.

Allocating a Raw Disk Partition on Windows

On Windows systems, the same steps and accompanying guidelines described for Linux and Unix systems apply except that the `innodb_data_file_path` setting differs slightly on Windows.

1. When creating a new data file, specify the keyword `newraw` immediately after the data file size for the `innodb_data_file_path` option:

```
[mysqld]
innodb_data_home_dir=
innodb_data_file_path=//./D:::10Gnewraw
```

The `//./` corresponds to the Windows syntax of `\.\` for accessing physical drives. In the example above, `D:` is the drive letter of the partition.

2. Restart the server. `InnoDB` notices the `newraw` keyword and initializes the new partition.
3. After `InnoDB` has initialized the new partition, stop the server, change `newraw` in the data file specification to `raw`:

```
[mysqld]
innodb_data_home_dir=
innodb_data_file_path=//./D:::10Graw
```

4. Restart the server. `InnoDB` now permits changes to be made.

15.6.3.2 File-Per-Table Tablespaces

A file-per-table tablespace contains data and indexes for a single `InnoDB` table, and is stored on the file system in a single data file.

File-per-table tablespace characteristics are described under the following topics in this section:

- [File-Per-Table Tablespace Configuration](#)
- [File-Per-Table Tablespace Data Files](#)
- [File-Per-Table Tablespace Advantages](#)
- [File-Per-Table Tablespace Disadvantages](#)

File-Per-Table Tablespace Configuration

`InnoDB` creates tables in file-per-table tablespaces by default. This behavior is controlled by the `innodb_file_per_table` variable. Disabling `innodb_file_per_table` causes `InnoDB` to create tables in the system tablespace.

An `innodb_file_per_table` setting can be specified in an option file or configured at runtime using a `SET GLOBAL` statement. Changing the setting at runtime requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#).

Option file:

```
[mysqld]
innodb_file_per_table=ON
```

Using `SET GLOBAL` at runtime:

```
mysql> SET GLOBAL innodb_file_per_table=ON;
```

File-Per-Table Tablespace Data Files

A file-per-table tablespace is created in an `.ibd` data file in a schema directory under the MySQL data directory. The `.ibd` file is named for the table (`table_name.ibd`). For example, the data file for table `test.t1` is created in the `test` directory under the MySQL data directory:

```
mysql> USE test;
mysql> CREATE TABLE t1 (
    id INT PRIMARY KEY AUTO_INCREMENT,
    name VARCHAR(100)
) ENGINE = InnoDB;
$> cd /path/to/mysql/data/test
$> ls
t1.ibd
```

You can use the `DATA DIRECTORY` clause of the `CREATE TABLE` statement to implicitly create a file-per-table tablespace data file outside of the data directory. For more information, see [Section 15.6.1.2, “Creating Tables Externally”](#).

File-Per-Table Tablespace Advantages

File-per-table tablespaces have the following advantages over shared tablespaces such as the system tablespace or general tablespaces.

- Disk space is returned to the operating system after truncating or dropping a table created in a file-per-table tablespace. Truncating or dropping a table stored in a shared tablespace creates free space within the shared tablespace data file, which can only be used for `InnoDB` data. In other words, a shared tablespace data file does not shrink in size after a table is truncated or dropped.
- A table-copying `ALTER TABLE` operation on a table that resides in a shared tablespace can increase the amount of disk space occupied by the tablespace. Such operations may require as much additional space as the data in the table plus indexes. This space is not released back to the operating system as it is for file-per-table tablespaces.
- `TRUNCATE TABLE` performance is better when executed on tables that reside in file-per-table tablespaces.
- File-per-table tablespace data files can be created on separate storage devices for I/O optimization, space management, or backup purposes. See [Section 15.6.1.2, “Creating Tables Externally”](#).
- You can import a table that resides in file-per-table tablespace from another MySQL instance. See [Section 15.6.1.3, “Importing InnoDB Tables”](#).
- Tables created in file-per-table tablespaces support features associated with `DYNAMIC` and `COMPRESSED` row formats, which are not supported by the system tablespace. See [Section 15.10, “InnoDB Row Formats”](#).
- Tables stored in individual tablespace data files can save time and improve chances for a successful recovery when data corruption occurs, when backups or binary logs are unavailable, or when the MySQL server instance cannot be restarted.
- Tables created in file-per-table tablespaces can be backed up or restored quickly using MySQL Enterprise Backup, without interrupting the use of other `InnoDB` tables. This is beneficial for tables on varying backup schedules or that require backup less frequently. See [Making a Partial Backup](#) for details.

- File-per-table tablespaces permit monitoring table size on the file system by monitoring the size of the tablespace data file.
- Common Linux file systems do not permit concurrent writes to a single file such as a shared tablespace data file when `innodb_flush_method` is set to `O_DIRECT`. As a result, there are possible performance improvements when using file-per-table tablespaces in conjunction with this setting.
- Tables in a shared tablespace are limited in size by the 64TB tablespace size limit. By comparison, each file-per-table tablespace has a 64TB size limit, which provides plenty of room for individual tables to grow in size.

File-Per-Table Tablespace Disadvantages

File-per-table tablespaces have the following disadvantages compared to shared tablespaces such as the system tablespace or general tablespaces.

- With file-per-table tablespaces, each table may have unused space that can only be utilized by rows of the same table, which can lead to wasted space if not properly managed.
- `fsync` operations are performed on multiple file-per-table data files instead of a single shared tablespace data file. Because `fsync` operations are per file, write operations for multiple tables cannot be combined, which can result in a higher total number of `fsync` operations.
- `mysqld` must keep an open file handle for each file-per-table tablespace, which may impact performance if you have numerous tables in file-per-table tablespaces.
- More file descriptors are required when each table has its own data file.
- There is potential for more fragmentation, which can impede `DROP TABLE` and table scan performance. However, if fragmentation is managed, file-per-table tablespaces can improve performance for these operations.
- The buffer pool is scanned when dropping a table that resides in a file-per-table tablespace, which can take several seconds for large buffer pools. The scan is performed with a broad internal lock, which may delay other operations.
- The `innodb_autoextend_increment` variable, which defines the increment size for extending the size of an auto-extending shared tablespace file when it becomes full, does not apply to file-per-table tablespace files, which are auto-extending regardless of the `innodb_autoextend_increment` setting. Initial file-per-table tablespace extensions are by small amounts, after which extensions occur in increments of 4MB.

15.6.3.3 General Tablespaces

A general tablespace is a shared InnoDB tablespace that is created using `CREATE TABLESPACE` syntax. General tablespace capabilities and features are described under the following topics in this section:

- [General Tablespace Capabilities](#)
- [Creating a General Tablespace](#)
- [Adding Tables to a General Tablespace](#)
- [General Tablespace Row Format Support](#)
- [Moving Tables Between Tablespaces Using ALTER TABLE](#)
- [Renaming a General Tablespace](#)
- [Dropping a General Tablespace](#)

- General Tablespace Limitations

General Tablespace Capabilities

General tablespaces provide the following capabilities:

- Similar to the system tablespace, general tablespaces are shared tablespaces capable of storing data for multiple tables.
- General tablespaces have a potential memory advantage over [file-per-table tablespaces](#). The server keeps tablespace metadata in memory for the lifetime of a tablespace. Multiple tables in fewer general tablespaces consume less memory for tablespace metadata than the same number of tables in separate file-per-table tablespaces.
- General tablespace data files can be placed in a directory relative to or independent of the MySQL data directory, which provides you with many of the data file and storage management capabilities of [file-per-table tablespaces](#). As with file-per-table tablespaces, the ability to place data files outside of the MySQL data directory allows you to manage performance of critical tables separately, setup RAID or DRBD for specific tables, or bind tables to particular disks, for example.
- General tablespaces support all table row formats and associated features.
- The `TABLESPACE` option can be used with `CREATE TABLE` to create tables in a general tablespaces, file-per-table tablespace, or in the system tablespace.
- The `TABLESPACE` option can be used with `ALTER TABLE` to move tables between general tablespaces, file-per-table tablespaces, and the system tablespace.

Creating a General Tablespace

General tablespaces are created using `CREATE TABLESPACE` syntax.

```
CREATE TABLESPACE tablespace_name
  [ADD DATAFILE 'file_name']
  [FILE_BLOCK_SIZE = value]
  [ENGINE [=] engine_name]
```

A general tablespace can be created in the data directory or outside of it. To avoid conflicts with implicitly created file-per-table tablespaces, creating a general tablespace in a subdirectory under the data directory is not supported. When creating a general tablespace outside of the data directory, the directory must exist and must be known to `InnoDB` prior to creating the tablespace. To make an unknown directory known to `InnoDB`, add the directory to the `innodb_directories` argument value. `innodb_directories` is a read-only startup option. Configuring it requires restarting the server.

Examples:

Creating a general tablespace in the data directory:

```
mysql> CREATE TABLESPACE `ts1` ADD DATAFILE 'ts1.ibd' Engine=InnoDB;
```

Or

```
mysql> CREATE TABLESPACE `ts1` Engine=InnoDB;
```

The `ADD DATAFILE` clause is optional as of MySQL 8.0.14 and required before that. If the `ADD DATAFILE` clause is not specified when creating a tablespace, a tablespace data file with a unique file name is created implicitly. The unique file name is a 128 bit UUID formatted into five groups of hexadecimal numbers separated by dashes (`aaaaaaaa-bbbb-cccc-dddd-eeeeeeeeeee`). General tablespace data files include an `.ibd` file extension. In a replication environment, the data file name created on the source is not the same as the data file name created on the replica.

Creating a general tablespace in a directory outside of the data directory:

```
mysql> CREATE TABLESPACE `ts1` ADD DATAFILE '/my/tablespace/directory/ts1.ibd' Engine=InnoDB;
```

You can specify a path that is relative to the data directory as long as the tablespace directory is not under the data directory. In this example, the `my_tablespace` directory is at the same level as the data directory:

```
mysql> CREATE TABLESPACE `ts1` ADD DATAFILE '../my_tablespace/ts1.ibd' Engine=InnoDB;
```



Note

The `ENGINE = InnoDB` clause must be defined as part of the `CREATE TABLESPACE` statement, or `InnoDB` must be defined as the default storage engine (`default_storage_engine=InnoDB`).

Adding Tables to a General Tablespace

After creating a general tablespace, `CREATE TABLE tbl_name ... TABLESPACE [=] tablespace_name` or `ALTER TABLE tbl_name TABLESPACE [=] tablespace_name` statements can be used to add tables to the tablespace, as shown in the following examples:

`CREATE TABLE:`

```
mysql> CREATE TABLE t1 (c1 INT PRIMARY KEY) TABLESPACE ts1;
```

`ALTER TABLE:`

```
mysql> ALTER TABLE t2 TABLESPACE ts1;
```



Note

Support for adding table partitions to shared tablespaces was deprecated in MySQL 5.7.24 and removed in MySQL 8.0.13. Shared tablespaces include the `InnoDB` system tablespace and general tablespaces.

For detailed syntax information, see `CREATE TABLE` and `ALTER TABLE`.

General Tablespace Row Format Support

General tablespaces support all table row formats (`REDUNDANT`, `COMPACT`, `DYNAMIC`, `COMPRESSED`) with the caveat that compressed and uncompressed tables cannot coexist in the same general tablespace due to different physical page sizes.

For a general tablespace to contain compressed tables (`ROW_FORMAT=COMPRESSED`), the `FILE_BLOCK_SIZE` option must be specified, and the `FILE_BLOCK_SIZE` value must be a valid compressed page size in relation to the `innodb_page_size` value. Also, the physical page size of the compressed table (`KEY_BLOCK_SIZE`) must be equal to `FILE_BLOCK_SIZE/1024`. For example, if `innodb_page_size=16KB` and `FILE_BLOCK_SIZE=8K`, the `KEY_BLOCK_SIZE` of the table must be 8.

The following table shows permitted `innodb_page_size`, `FILE_BLOCK_SIZE`, and `KEY_BLOCK_SIZE` combinations. `FILE_BLOCK_SIZE` values may also be specified in bytes. To determine a valid `KEY_BLOCK_SIZE` value for a given `FILE_BLOCK_SIZE`, divide the `FILE_BLOCK_SIZE` value by 1024. Table compression is not supported for 32K and 64K `InnoDB` page sizes. For more information about `KEY_BLOCK_SIZE`, see `CREATE TABLE`, and [Section 15.9.1.2, “Creating Compressed Tables”](#).

Table 15.3 Permitted Page Size, FILE_BLOCK_SIZE, and KEY_BLOCK_SIZE Combinations for Compressed Tables

InnoDB Page Size (<code>innodb_page_size</code>)	Permitted FILE_BLOCK_SIZE Value	Permitted KEY_BLOCK_SIZE Value
64KB	64K (65536)	Compression is not supported
32KB	32K (32768)	Compression is not supported

InnoDB Page Size (innodb_page_size)	Permitted FILE_BLOCK_SIZE Value	Permitted KEY_BLOCK_SIZE Value
16KB	16K (16384)	None. If <code>innodb_page_size</code> is equal to <code>FILE_BLOCK_SIZE</code> , the tablespace cannot contain a compressed table.
16KB	8K (8192)	8
16KB	4K (4096)	4
16KB	2K (2048)	2
16KB	1K (1024)	1
8KB	8K (8192)	None. If <code>innodb_page_size</code> is equal to <code>FILE_BLOCK_SIZE</code> , the tablespace cannot contain a compressed table.
8KB	4K (4096)	4
8KB	2K (2048)	2
8KB	1K (1024)	1
4KB	4K (4096)	None. If <code>innodb_page_size</code> is equal to <code>FILE_BLOCK_SIZE</code> , the tablespace cannot contain a compressed table.
4KB	2K (2048)	2
4KB	1K (1024)	1

This example demonstrates creating a general tablespace and adding a compressed table. The example assumes a default `innodb_page_size` of 16KB. The `FILE_BLOCK_SIZE` of 8192 requires that the compressed table have a `KEY_BLOCK_SIZE` of 8.

```
mysql> CREATE TABLESPACE `ts2` ADD DATAFILE 'ts2.ibd' FILE_BLOCK_SIZE = 8192 Engine=InnoDB;
mysql> CREATE TABLE t4 (c1 INT PRIMARY KEY) TABLESPACE ts2 ROW_FORMAT=COMPRESSED KEY_BLOCK_SIZE=8;
```

If you do not specify `FILE_BLOCK_SIZE` when creating a general tablespace, `FILE_BLOCK_SIZE` defaults to `innodb_page_size`. When `FILE_BLOCK_SIZE` is equal to `innodb_page_size`, the tablespace may only contain tables with an uncompressed row format (`COMPACT`, `REDUNDANT`, and `DYNAMIC` row formats).

Moving Tables Between Tablespaces Using ALTER TABLE

`ALTER TABLE` with the `TABLESPACE` option can be used to move a table to an existing general tablespace, to a new file-per-table tablespace, or to the system tablespace.



Note

Support for placing table partitions in shared tablespaces was deprecated in MySQL 5.7.24 and removed MySQL 8.0.13. Shared tablespaces include the `InnoDB` system tablespace and general tablespaces.

To move a table from a file-per-table tablespace or from the system tablespace to a general tablespace, specify the name of the general tablespace. The general tablespace must exist. See `ALTER TABLESPACE` for more information.

```
ALTER TABLE tbl_name TABLESPACE [=] tablespace_name;
```

To move a table from a general tablespace or file-per-table tablespace to the system tablespace, specify `innodb_system` as the tablespace name.

```
ALTER TABLE tbl_name TABLESPACE [=] innodb_system;
```

To move a table from the system tablespace or a general tablespace to a file-per-table tablespace, specify `innodb_file_per_table` as the tablespace name.

```
ALTER TABLE tbl_name TABLESPACE [=] innodb_file_per_table;
```

`ALTER TABLE ... TABLESPACE` operations cause a full table rebuild, even if the `TABLESPACE` attribute has not changed from its previous value.

`ALTER TABLE ... TABLESPACE` syntax does not support moving a table from a temporary tablespace to a persistent tablespace.

The `DATA DIRECTORY` clause is permitted with `CREATE TABLE ... TABLESPACE=innodb_file_per_table` but is otherwise not supported for use in combination with the `TABLESPACE` option. As of MySQL 8.0.21, the directory specified in a `DATA DIRECTORY` clause must be known to InnoDB. For more information, see [Using the DATA DIRECTORY Clause](#).

Restrictions apply when moving tables from encrypted tablespaces. See [Encryption Limitations](#).

Renaming a General Tablespace

Renaming a general tablespace is supported using `ALTER TABLESPACE ... RENAME TO` syntax.

```
ALTER TABLESPACE s1 RENAME TO s2;
```

The `CREATE TABLESPACE` privilege is required to rename a general tablespace.

`RENAME TO` operations are implicitly performed in `autocommit` mode regardless of the `autocommit` setting.

A `RENAME TO` operation cannot be performed while `LOCK TABLES` or `FLUSH TABLES WITH READ LOCK` is in effect for tables that reside in the tablespace.

Exclusive `metadata locks` are taken on tables within a general tablespace while the tablespace is renamed, which prevents concurrent DDL. Concurrent DML is supported.

Dropping a General Tablespace

The `DROP TABLESPACE` statement is used to drop an InnoDB general tablespace.

All tables must be dropped from the tablespace prior to a `DROP TABLESPACE` operation. If the tablespace is not empty, `DROP TABLESPACE` returns an error.

Use a query similar to the following to identify tables in a general tablespace.

```
mysql> SELECT a.NAME AS space_name, b.NAME AS table_name FROM INFORMATION_SCHEMA.INNODB_TABLESPACES a,
    >     INFORMATION_SCHEMA.INNODB_TABLES b WHERE a.SPACE=b.SPACE AND a.NAME LIKE 'ts1';
+-----+-----+
| space_name | table_name |
+-----+-----+
| ts1        | test/t1    |
| ts1        | test/t2    |
| ts1        | test/t3    |
+-----+-----+
```

A general InnoDB tablespace is not deleted automatically when the last table in the tablespace is dropped. The tablespace must be dropped explicitly using `DROP TABLESPACE tablespace_name`.

A general tablespace does not belong to any particular database. A `DROP DATABASE` operation can drop tables that belong to a general tablespace but it cannot drop the tablespace, even if the `DROP DATABASE` operation drops all tables that belong to the tablespace.

Similar to the system tablespace, truncating or dropping tables stored in a general tablespace creates free space internally in the general tablespace .ibd data file which can only be used for new InnoDB

data. Space is not released back to the operating system as it is when a file-per-table tablespace is deleted during a `DROP TABLE` operation.

This example demonstrates how to drop an `InnoDB` general tablespace. The general tablespace `ts1` is created with a single table. The table must be dropped before dropping the tablespace.

```
mysql> CREATE TABLESPACE `ts1` ADD DATAFILE 'ts1.ibd' Engine=InnoDB;
mysql> CREATE TABLE t1 (c1 INT PRIMARY KEY) TABLESPACE ts1 Engine=InnoDB;
mysql> DROP TABLE t1;
mysql> DROP TABLESPACE ts1;
```



Note

`tablespace_name` is a case-sensitive identifier in MySQL.

General Tablespace Limitations

- A generated or existing tablespace cannot be changed to a general tablespace.
- Creation of temporary general tablespaces is not supported.
- General tablespaces do not support temporary tables.
- Similar to the system tablespace, truncating or dropping tables stored in a general tablespace creates free space internally in the general tablespace `.ibd` data file which can only be used for new `InnoDB` data. Space is not released back to the operating system as it is for `file-per-table` tablespaces.

Additionally, a table-copying `ALTER TABLE` operation on table that resides in a shared tablespace (a general tablespace or the system tablespace) can increase the amount of space used by the tablespace. Such operations require as much additional space as the data in the table plus indexes. The additional space required for the table-copying `ALTER TABLE` operation is not released back to the operating system as it is for file-per-table tablespaces.

- `ALTER TABLE ... DISCARD TABLESPACE` and `ALTER TABLE ... IMPORT TABLESPACE` are not supported for tables that belong to a general tablespace.
- Support for placing table partitions in general tablespaces was deprecated in MySQL 5.7.24 and removed in MySQL 8.0.13.
- The `ADD DATAFILE` clause is not supported in a replication environment where the source and replica reside on the same host, as it would cause the source and replica to create a tablespace of the same name in the same location, which is not supported. However, if the `ADD DATAFILE` clause is omitted, the tablespace is created in the data directory with a generated file name that is unique, which is permitted.
- As of MySQL 8.0.21, general tablespaces cannot be created in the undo tablespace directory (`innodb_undo_directory`) unless that directly is known to `InnoDB`. Known directories are those defined by the `datadir`, `innodb_data_home_dir`, and `innodb_directories` variables.

15.6.3.4 Undo Tablespaces

Undo tablespaces contain undo logs, which are collections of records containing information about how to undo the latest change by a transaction to a clustered index record.

Undo tablespaces are described under the following topics in this section:

- [Default Undo Tablespaces](#)
- [Undo Tablespace Size](#)

- [Adding Undo Tablespaces](#)
- [Dropping Undo Tablespaces](#)
- [Moving Undo Tablespaces](#)
- [Configuring the Number of Rollback Segments](#)
- [Truncating Undo Tablespaces](#)
- [Undo Tablespace Status Variables](#)

Default Undo Tablespaces

Two default undo tablespaces are created when the MySQL instance is initialized. Default undo tablespaces are created at initialization time to provide a location for rollback segments that must exist before SQL statements can be accepted. A minimum of two undo tablespaces is required to support automated truncation of undo tablespaces. See [Truncating Undo Tablespaces](#).

Default undo tablespaces are created in the location defined by the `innodb_undo_directory` variable. If the `innodb_undo_directory` variable is undefined, default undo tablespaces are created in the data directory. Default undo tablespace data files are named `undo_001` and `undo_002`. The corresponding undo tablespace names defined in the data dictionary are `innodb_undo_001` and `innodb_undo_002`.

As of MySQL 8.0.14, additional undo tablespaces can be created at runtime using SQL. See [Adding Undo Tablespaces](#).

Undo Tablespace Size

Prior to MySQL 8.0.23, the initial size of an undo tablespace depends on the `innodb_page_size` value. For the default 16KB page size, the initial undo tablespace file size is 10MiB. For 4KB, 8KB, 32KB, and 64KB page sizes, the initial undo tablespace files sizes are 7MiB, 8MiB, 20MiB, and 40MiB, respectively. As of MySQL 8.0.23, the initial undo tablespace size is normally 16MiB. The initial size may differ when a new undo tablespace is created by a truncate operation. In this case, if the file extension size is larger than 16MB, and the previous file extension occurred within the last second, the new undo tablespace is created at a quarter of the size defined by the `innodb_max_undo_log_size` variable.

Prior to MySQL 8.0.23, an undo tablespace is extended four extents at a time. From MySQL 8.0.23, an undo tablespace is extended by a minimum of 16MB. To handle aggressive growth, the file extension size is doubled if the previous file extension happened less than 0.1 seconds earlier. Doubling of the extension size can occur multiple times to a maximum of 256MB. If the previous file extension occurred more than 0.1 seconds earlier, the extension size is reduced by half, which can also occur multiple times, to a minimum of 16MB. If the `AUTOEXTEND_SIZE` option is defined for an undo tablespace, it is extended by the greater of the `AUTOEXTEND_SIZE` setting and the extension size determined by the logic described above. For information about the `AUTOEXTEND_SIZE` option, see [Section 15.6.3.9, “Tablespace AUTOEXTEND_SIZE Configuration”](#).

Adding Undo Tablespaces

Because undo logs can become large during long-running transactions, creating additional undo tablespaces can help prevent individual undo tablespaces from becoming too large. As of MySQL 8.0.14, additional undo tablespaces can be created at runtime using `CREATE UNDO TABLESPACE` syntax.

```
CREATE UNDO TABLESPACE tablespace_name ADD DATAFILE 'file_name.ibu';
```

The undo tablespace file name must have an `.ibu` extension. It is not permitted to specify a relative path when defining the undo tablespace file name. A fully qualified path is permitted, but the path must

be known to [InnoDB](#). Known paths are those defined by the `innodb_directories` variable. Unique undo tablespace file names are recommended to avoid potential file name conflicts when moving or cloning data.



Note

In a replication environment, the source and each replica must have its own undo tablespace file directory. Replicating the creation of an undo tablespace file to a common directory would cause a file name conflict.

At startup, directories defined by the `innodb_directories` variable are scanned for undo tablespace files. (The scan also traverses subdirectories.) Directories defined by the `innodb_data_home_dir`, `innodb_undo_directory`, and `datadir` variables are automatically appended to the `innodb_directories` value regardless of whether the `innodb_directories` variable is defined explicitly. An undo tablespace can therefore reside in paths defined by any of those variables.

If the undo tablespace file name does not include a path, the undo tablespace is created in the directory defined by the `innodb_undo_directory` variable. If that variable is undefined, the undo tablespace is created in the data directory.



Note

The [InnoDB](#) recovery process requires that undo tablespace files reside in known directories. Undo tablespace files must be discovered and opened before redo recovery and before other data files are opened to permit uncommitted transactions and data dictionary changes to be rolled back. An undo tablespace not found before recovery cannot be used, which can lead to database inconsistencies. An error message is reported at startup if an undo tablespace known to the data dictionary is not found. The known directory requirement also supports undo tablespace portability. See [Moving Undo Tablespaces](#).

To create undo tablespaces in a path relative to the data directory, set the `innodb_undo_directory` variable to the relative path, and specify the file name only when creating an undo tablespace.

To view undo tablespace names and paths, query `INFORMATION_SCHEMA.FILES`:

```
SELECT TABLESPACE_NAME, FILE_NAME FROM INFORMATION_SCHEMA.FILES
WHERE FILE_TYPE LIKE 'UNDO LOG';
```

A MySQL instance supports up to 127 undo tablespaces including the two default undo tablespaces created when the MySQL instance is initialized.



Note

Prior to MySQL 8.0.14, additional undo tablespaces are created by configuring the `innodb_undo_tablespaces` startup variable. This variable is deprecated and no longer configurable as of MySQL 8.0.14.

Prior to MySQL 8.0.14, increasing the `innodb_undo_tablespaces` setting creates the specified number of undo tablespaces and adds them to the list of active undo tablespaces. Decreasing the `innodb_undo_tablespaces` setting removes undo tablespaces from the list of active undo tablespaces. Undo tablespaces that are removed from the active list remain active until they are no longer used by existing transactions. The `innodb_undo_tablespaces` variable can be configured at runtime using a `SET` statement or defined in a configuration file.

Prior to MySQL 8.0.14, deactivated undo tablespaces cannot be removed. Manual removal of undo tablespace files is possible after a slow shutdown but is not recommended, as deactivated undo tablespaces may contain active undo

logs for some time after the server is restarted if open transactions were present when shutting down the server. As of MySQL 8.0.14, undo tablespaces can be dropped using `DROP UNDO TABLESPACE` syntax. See [Dropping Undo Tablespaces](#).

Dropping Undo Tablespaces

As of MySQL 8.0.14, undo tablespaces created using `CREATE UNDO TABLESPACE` syntax can be dropped at runtime using `DROP UNDO TABLESPACE` syntax.

An undo tablespace must be empty before it can be dropped. To empty an undo tablespace, the undo tablespace must first be marked as inactive using `ALTER UNDO TABLESPACE` syntax so that the tablespace is no longer used for assigning rollback segments to new transactions.

```
ALTER UNDO TABLESPACE tablespace_name SET INACTIVE;
```

After an undo tablespace is marked as inactive, transactions currently using rollback segments in the undo tablespace are permitted to finish, as are any transactions started before those transactions are completed. After transactions are completed, the purge system frees the rollback segments in the undo tablespace, and the undo tablespace is truncated to its initial size. (The same process is used when truncating undo tablespaces. See [Truncating Undo Tablespaces](#).) Once the undo tablespace is empty, it can be dropped.

```
DROP UNDO TABLESPACE tablespace_name;
```



Note

Alternatively, the undo tablespace can be left in an empty state and reactivated later, if needed, by issuing an `ALTER UNDO TABLESPACE tablespace_name SET ACTIVE` statement.

The state of an undo tablespace can be monitored by querying the Information Schema `INNODB_TABLESPACES` table.

```
SELECT NAME, STATE FROM INFORMATION_SCHEMA.INNODB_TABLESPACES
  WHERE NAME LIKE 'tablespace_name';
```

An `inactive` state indicates that rollback segments in an undo tablespace are no longer used by new transactions. An `empty` state indicates that an undo tablespace is empty and ready to be dropped, or ready to be made active again using an `ALTER UNDO TABLESPACE tablespace_name SET ACTIVE` statement. Attempting to drop an undo tablespace that is not empty returns an error.

The default undo tablespaces (`innodb_undo_001` and `innodb_undo_002`) created when the MySQL instance is initialized cannot be dropped. They can, however, be made inactive using an `ALTER UNDO TABLESPACE tablespace_name SET INACTIVE` statement. Before a default undo tablespace can be made inactive, there must be an undo tablespace to take its place. A minimum of two active undo tablespaces are required at all times to support automated truncation of undo tablespaces.

Moving Undo Tablespaces

Undo tablespaces created with `CREATE UNDO TABLESPACE` syntax can be moved while the server is offline to any known directory. Known directories are those defined by the `innodb_directories` variable. Directories defined by `innodb_data_home_dir`, `innodb_undo_directory`, and `datadir` are automatically appended to the `innodb_directories` value regardless of whether the `innodb_directories` variable is defined explicitly. Those directories and their subdirectories are scanned at startup for undo tablespaces files. An undo tablespace file moved to any of those directories is discovered at startup and assumed to be the undo tablespace that was moved.

The default undo tablespaces (`innodb_undo_001` and `innodb_undo_002`) created when the MySQL instance is initialized must reside in the directory defined by the `innodb_undo_directory`

variable. If the `innodb_undo_directory` variable is undefined, default undo tablespaces reside in the data directory. If default undo tablespaces are moved while the server is offline, the server must be started with the `innodb_undo_directory` variable configured to the new directory.

The I/O patterns for undo logs make undo tablespaces good candidates for [SSD](#) storage.

Configuring the Number of Rollback Segments

The `innodb_rollback_segments` variable defines the number of [rollback segments](#) allocated to each undo tablespace and to the global temporary tablespace. The `innodb_rollback_segments` variable can be configured at startup or while the server is running.

The default setting for `innodb_rollback_segments` is 128, which is also the maximum value. For information about the number of transactions that a rollback segment supports, see [Section 15.6.6, “Undo Logs”](#).

Truncating Undo Tablespaces

There are two methods of truncating undo tablespaces, which can be used individually or in combination to manage undo tablespace size. One method is automated, enabled using configuration variables. The other method is manual, performed using SQL statements.

The automated method does not require monitoring undo tablespace size and, once enabled, it performs deactivation, truncation, and reactivation of undo tablespaces without manual intervention. The manual truncation method may be preferable if you want to control when undo tablespaces are taken offline for truncation. For example, you may want to avoid truncating undo tablespaces during peak workload times.

Automated Truncation

Automated truncation of undo tablespaces requires a minimum of two active undo tablespaces, which ensures that one undo tablespace remains active while the other is taken offline to be truncated. By default, two undo tablespaces are created when the MySQL instance is initialized.

To have undo tablespaces automatically truncated, enable the `innodb_undo_log_truncate` variable. For example:

```
mysql> SET GLOBAL innodb_undo_log_truncate=ON;
```

When the `innodb_undo_log_truncate` variable is enabled, undo tablespaces that exceed the size limit defined by the `innodb_max_undo_log_size` variable are subject to truncation. The `innodb_max_undo_log_size` variable is dynamic and has a default value of 1073741824 bytes (1024 MiB).

```
mysql> SELECT @@innodb_max_undo_log_size;
+-----+
| @@innodb_max_undo_log_size |
+-----+
|          1073741824 |
+-----+
```

When the `innodb_undo_log_truncate` variable is enabled:

1. Default and user-defined undo tablespaces that exceed the `innodb_max_undo_log_size` setting are marked for truncation. Selection of an undo tablespace for truncation is performed in a circular fashion to avoid truncating the same undo tablespace each time.
2. Rollback segments residing in the selected undo tablespace are made inactive so that they are not assigned to new transactions. Existing transactions that are currently using rollback segments are permitted to finish.
3. The `purge` system empties rollback segments by freeing undo logs that are no longer in use.

4. After all rollback segments in the undo tablespace are freed, the truncate operation runs and truncates the undo tablespace to its initial size.

The size of an undo tablespace after a truncate operation may be larger than the initial size due to immediate use following the completion of the operation.

The `innodb_undo_directory` variable defines the location of default undo tablespace files. If the `innodb_undo_directory` variable is undefined, default undo tablespaces reside in the data directory. The location of all undo tablespace files including user-defined undo tablespaces created using `CREATE UNDO TABLESPACE` syntax can be determined by querying the Information Schema `FILES` table:

```
SELECT TABLESPACE_NAME, FILE_NAME FROM INFORMATION_SCHEMA.FILES WHERE FILE_TYPE LIKE 'UNDO LOG';
```

5. Rollback segments are reactivated so that they can be assigned to new transactions.

Manual Truncation

Manual truncation of undo tablespaces requires a minimum of three active undo tablespaces. Two active undo tablespaces are required at all times to support the possibility that automated truncation is enabled. A minimum of three undo tablespaces satisfies this requirement while permitting an undo tablespace to be taken offline manually.

To manually initiate truncation of an undo tablespace, deactivate the undo tablespace by issuing the following statement:

```
ALTER UNDO TABLESPACE tablespace_name SET INACTIVE;
```

After the undo tablespace is marked as inactive, transactions currently using rollback segments in the undo tablespace are permitted to finish, as are any transactions started before those transactions are completed. After transactions are completed, the purge system frees the rollback segments in the undo tablespace, the undo tablespace is truncated to its initial size, and the undo tablespace state changes from `inactive` to `empty`.



Note

When an `ALTER UNDO TABLESPACE tablespace_name SET INACTIVE` statement deactivates an undo tablespace, the purge thread looks for that undo tablespace at the next opportunity. Once the undo tablespace is found and marked for truncation, the purge thread returns with increased frequency to quickly empty and truncate the undo tablespace.

To check the state of an undo tablespace, query the Information Schema `INNODB_TABLESPACES` table.

```
SELECT NAME, STATE FROM INFORMATION_SCHEMA.INNODB_TABLESPACES
WHERE NAME LIKE 'tablespace_name';
```

Once the undo tablespace is in an `empty` state, it can be reactivated by issuing the following statement:

```
ALTER UNDO TABLESPACE tablespace_name SET ACTIVE;
```

An undo tablespace in an `empty` state can also be dropped. See [Dropping Undo Tablespace](#).

Expediting Automated Truncation of Undo Tablespaces

The purge thread is responsible for emptying and truncating undo tablespaces. By default, the purge thread looks for undo tablespaces to truncate once every 128 times that purge is invoked. The frequency with which the purge thread looks for undo tablespaces to truncate is controlled by the `innodb_purge_rseg_truncate_frequency` variable, which has a default setting of 128.

```
mysql> SELECT @@innodb_purge_rseg_truncate_frequency;
+-----+
| @innodb_purge_rseg_truncate_frequency |
+-----+
| 128 |
+-----+
```

To increase the frequency, decrease the `innodb_purge_rseg_truncate_frequency` setting. For example, to have the purge thread look for undo tablespaces once every 32 times that purge is invoked, set `innodb_purge_rseg_truncate_frequency` to 32.

```
mysql> SET GLOBAL innodb_purge_rseg_truncate_frequency=32;
```

Performance Impact of Truncating Undo Tablespace Files

When an undo tablespace is truncated, the rollback segments in the undo tablespace are deactivated. The active rollback segments in other undo tablespaces assume responsibility for the entire system load, which may result in a slight performance degradation. The extent to which performance is affected depends on a number of factors:

- Number of undo tablespaces
- Number of undo logs
- Undo tablespace size
- Speed of the I/O subsystem
- Existing long running transactions
- System load

The easiest way to avoid the potential performance impact is to increase the number of undo tablespaces.

Monitoring Undo Tablespace Truncation

As of MySQL 8.0.16, `undo` and `purge` subsystem counters are provided for monitoring background activities associated with undo log truncation. For counter names and descriptions, query the Information Schema `INNODB_METRICS` table.

```
SELECT NAME, SUBSYSTEM, COMMENT FROM INFORMATION_SCHEMA.INNODB_METRICS WHERE NAME LIKE '%truncate%';
```

For information about enabling counters and querying counter data, see [Section 15.15.6, “InnoDB INFORMATION_SCHEMA Metrics Table”](#).

Undo Tablespace Truncation Limit

As of MySQL 8.0.21, the number of truncate operations on the same undo tablespace between checkpoints is limited to 64. The limit prevents potential issues caused by an excessive number of undo tablespace truncate operations, which can occur if `innodb_max_undo_log_size` is set too low on a busy system, for example. If the limit is exceeded, an undo tablespace can still be made inactive, but it is not truncated until after the next checkpoint. The limit was raised from 64 to 50,000 in MySQL 8.0.22.

Undo Tablespace Truncation Recovery

An undo tablespace truncate operation creates a temporary `undo_space_number_trunc.log` file in the server log directory. That log directory is defined by `innodb_log_group_home_dir`. If a system failure occurs during the truncate operation, the temporary log file permits the startup process to identify undo tablespaces that were being truncated and to continue the operation.

Undo Tablespace Status Variables

The following status variables permit tracking the total number of undo tablespaces, implicit ([InnoDB](#)-created) undo tablespaces, explicit (user-created) undo tablespaces, and the number of active undo tablespaces:

```
mysql> SHOW STATUS LIKE 'Innodb_undo_tablespaces%';
+-----+-----+
| Variable_name          | Value |
+-----+-----+
| Innodb_undo_tablespaces_total | 2      |
| Innodb_undo_tablespaces_implicit | 2      |
| Innodb_undo_tablespaces_explicit | 0      |
| Innodb_undo_tablespaces_active | 2      |
+-----+-----+
```

For status variable descriptions, see [Section 5.1.10, “Server Status Variables”](#).

15.6.3.5 Temporary Tablespaces

[InnoDB](#) uses session temporary tablespaces and a global temporary tablespace.

Session Temporary Tablespaces

Session temporary tablespaces store user-created temporary tables and internal temporary tables created by the optimizer when [InnoDB](#) is configured as the storage engine for on-disk internal temporary tables. Beginning with MySQL 8.0.16, the storage engine used for on-disk internal temporary tables is [InnoDB](#). (Previously, the storage engine was determined by the value of `internal_tmp_disk_storage_engine`.)

Session temporary tablespaces are allocated to a session from a pool of temporary tablespaces on the first request to create an on-disk temporary table. A maximum of two tablespaces is allocated to a session, one for user-created temporary tables and the other for internal temporary tables created by the optimizer. The temporary tablespaces allocated to a session are used for all on-disk temporary tables created by the session. When a session disconnects, its temporary tablespaces are truncated and released back to the pool. A pool of 10 temporary tablespaces is created when the server is started. The size of the pool never shrinks and tablespaces are added to the pool automatically as necessary. The pool of temporary tablespaces is removed on normal shutdown or on an aborted initialization. Session temporary tablespace files are five pages in size when created and have an `.ibt` file name extension.

A range of 400 thousand space IDs is reserved for session temporary tablespaces. Because the pool of session temporary tablespaces is recreated each time the server is started, space IDs for session temporary tablespaces are not persisted when the server is shut down, and may be reused.

The `innodb_temp_tablespaces_dir` variable defines the location where session temporary tablespaces are created. The default location is the `#innodb_temp` directory in the data directory. Startup is refused if the pool of temporary tablespaces cannot be created.

```
$> cd $BASEDIR/data/#innodb_temp
$> ls
temp_10.ibt  temp_2.ibt  temp_4.ibt  temp_6.ibt  temp_8.ibt
temp_1.ibt   temp_3.ibt  temp_5.ibt  temp_7.ibt  temp_9.ibt
```

In statement based replication (SBR) mode, temporary tables created on a replica reside in a single session temporary tablespace that is truncated only when the MySQL server is shut down.

The `INNODB_SESSION_TEMP_TABLESPACES` table provides metadata about session temporary tablespaces.

The Information Schema `INNODB_TEMP_TABLE_INFO` table provides metadata about user-created temporary tables that are active in an [InnoDB](#) instance.

Global Temporary Tablespace

The global temporary tablespace (`ibtmp1`) stores rollback segments for changes made to user-created temporary tables.

The `innodb_temp_data_file_path` variable defines the relative path, name, size, and attributes for global temporary tablespace data files. If no value is specified for `innodb_temp_data_file_path`, the default behavior is to create a single auto-extending data file named `ibtmp1` in the `innodb_data_home_dir` directory. The initial file size is slightly larger than 12MB.

The global temporary tablespace is removed on normal shutdown or on an aborted initialization, and recreated each time the server is started. The global temporary tablespace receives a dynamically generated space ID when it is created. Startup is refused if the global temporary tablespace cannot be created. The global temporary tablespace is not removed if the server halts unexpectedly. In this case, a database administrator can remove the global temporary tablespace manually or restart the MySQL server. Restarting the MySQL server removes and recreates the global temporary tablespace automatically.

The global temporary tablespace cannot reside on a raw device.

The Information Schema `FILES` table provides metadata about the global temporary tablespace. Issue a query similar to this one to view global temporary tablespace metadata:

```
mysql> SELECT * FROM INFORMATION_SCHEMA.FILES WHERE TABLESPACE_NAME='innodb_temporary'\G
```

By default, the global temporary tablespace data file is autoextending and increases in size as necessary.

To determine if a global temporary tablespace data file is autoextending, check the `innodb_temp_data_file_path` setting:

```
mysql> SELECT @@innodb_temp_data_file_path;
+-----+
| @@innodb_temp_data_file_path |
+-----+
| ibtmp1:12M:autoextend      |
+-----+
```

To check the size of global temporary tablespace data files, examine the Information Schema `FILES` table using a query similar to this one:

```
mysql> SELECT FILE_NAME, TABLESPACE_NAME, ENGINE, INITIAL_SIZE, TOTAL_EXTENTS*EXTENT_SIZE
       AS TotalSizeBytes, DATA_FREE, MAXIMUM_SIZE FROM INFORMATION_SCHEMA.FILES
       WHERE TABLESPACE_NAME = 'innodb_temporary'\G
***** 1. row *****
FILE_NAME: ./ibtmp1
TABLESPACE_NAME: innodb_temporary
ENGINE: InnoDB
INITIAL_SIZE: 12582912
TotalSizeBytes: 12582912
DATA_FREE: 6291456
MAXIMUM_SIZE: NULL
```

`TotalSizeBytes` shows the current size of the global temporary tablespace data file. For information about other field values, see [Section 26.3.15, “The INFORMATION_SCHEMA FILES Table”](#).

Alternatively, check the global temporary tablespace data file size on your operating system. The global temporary tablespace data file is located in the directory defined by the `innodb_temp_data_file_path` variable.

To reclaim disk space occupied by a global temporary tablespace data file, restart the MySQL server. Restarting the server removes and recreates the global temporary tablespace data file according to the attributes defined by `innodb_temp_data_file_path`.

To limit the size of the global temporary tablespace data file, configure `innodb_temp_data_file_path` to specify a maximum file size. For example:

```
[mysqld]
innodb_temp_data_file_path=ibtmp1:12M:autoextend:max:500M
```

Configuring `innodb_temp_data_file_path` requires restarting the server.

15.6.3.6 Moving Tablespace Files While the Server is Offline

The `innodb_directories` variable, which defines directories to scan at startup for tablespace files, supports moving or restoring tablespace files to a new location while the server is offline. During startup, discovered tablespace files are used instead those referenced in the data dictionary, and the data dictionary is updated to reference the relocated files. If duplicate tablespace files are discovered by the scan, startup fails with an error indicating that multiple files were found for the same tablespace ID.

The directories defined by the `innodb_data_home_dir`, `innodb_undo_directory`, and `datadir` variables are automatically appended to the `innodb_directories` argument value. These directories are scanned at startup regardless of whether an `innodb_directories` setting is specified explicitly. The implicit addition of these directories permits moving system tablespace files, the data directory, or undo tablespace files without configuring the `innodb_directories` setting. However, settings must be updated when directories change. For example, after relocating the data directory, you must update the `--datadir` setting before restarting the server.

The `innodb_directories` variable can be specified in a startup command or MySQL option file. Quotes are used around the argument value because a semicolon (;) is interpreted as a special character by some command interpreters. (Unix shells treat it as a command terminator, for example.)

Startup command:

```
mysqld --innodb-directories="directory_path_1;directory_path_2"
```

MySQL option file:

```
[mysqld]
innodb_directories="directory_path_1;directory_path_2"
```

The following procedure is applicable to moving individual [file-per-table](#) and [general tablespace](#) files, [system tablespace](#) files, [undo tablespace](#) files, or the data directory. Before moving files or directories, review the usage notes that follow.

1. Stop the server.
2. Move the tablespace files or directories to the desired location.
3. Make the new directory known to InnoDB.
 - If moving individual [file-per-table](#) or [general tablespace](#) files, add unknown directories to the `innodb_directories` value.
 - The directories defined by the `innodb_data_home_dir`, `innodb_undo_directory`, and `datadir` variables are automatically appended to the `innodb_directories` argument value, so you need not specify these.
 - A file-per-table tablespace file can only be moved to a directory with same name as the schema. For example, if the `actor` table belongs to the `sakila` schema, then the `actor.ibd` data file can only be moved to a directory named `sakila`.
 - General tablespace files cannot be moved to the data directory or a subdirectory of the data directory.

- If moving system tablespace files, undo tablespaces, or the data directory, update the `innodb_data_home_dir`, `innodb_undo_directory`, and `datadir` settings, as necessary.
4. Restart the server.

Usage Notes

- Wildcard expressions cannot be used in the `innodb_directories` argument value.
- The `innodb_directories` scan also traverses subdirectories of specified directories. Duplicate directories and subdirectories are discarded from the list of directories to be scanned.
- `innodb_directories` supports moving InnoDB tablespace files. Moving files that belong to a storage engine other than InnoDB is not supported. This restriction also applies when moving the entire data directory.
- `innodb_directories` supports renaming of tablespace files when moving files to a scanned directory. It also supports moving tablespaces files to other supported operating systems.
- When moving tablespace files to a different operating system, ensure that tablespace file names do not include prohibited characters or characters with a special meaning on the destination system.
- When moving a data directory from a Windows operating system to a Linux operating system, modify the binary log file paths in the binary log index file to use backward slashes instead of forward slashes. By default, the binary log index file has the same base name as the binary log file, with the extension '`.index`'. The location of the binary log index file is defined by `--log-bin`. The default location is the data directory.
- If moving tablespace files to a different operating system introduces cross-platform replication, it is the database administrator's responsibility to ensure proper replication of DDL statements that contain platform-specific directories. Statements that permit specifying directories include `CREATE TABLE ... DATA DIRECTORY` and `CREATE TABLESPACE ... ADD DATAFILE`.
- Add the directories of file-per-table and general tablespaces created with an absolute path or in a location outside of the data directory to the `innodb_directories` setting. Otherwise, InnoDB is not able to locate the files during recovery. For related information, see [Tablespace Discovery During Crash Recovery](#).

To view tablespace file locations, query the Information Schema `FILES` table:

```
mysql> SELECT TABLESPACE_NAME, FILE_NAME FROM INFORMATION_SCHEMA.FILES \G
```

15.6.3.7 Disabling Tablespace Path Validation

At startup, InnoDB scans directories defined by the `innodb_directories` variable for tablespace files. The paths of discovered tablespace files are validated against the paths recorded in the data dictionary. If the paths do not match, the paths in the data dictionary are updated.

The `innodb_validate_tablespace_paths` variable, introduced in MySQL 8.0.21, permits disabling tablespace path validation. This feature is intended for environments where tablespaces files are not moved. Disabling path validation improves startup time on systems with a large number of tablespace files. If `log_error_verbosity` is set to 3, the following message is printed at startup when tablespace path validation is disabled:

```
[InnoDB] Skipping InnoDB tablespace path validation.  
Manually moved tablespace files will not be detected!
```



Warning

Starting the server with tablespace path validation disabled after moving tablespace files can lead to undefined behavior.

15.6.3.8 Optimizing Tablespace Space Allocation on Linux

As of MySQL 8.0.22, you can optimize how `InnoDB` allocates space to file-per-table and general tablespaces on Linux. By default, when additional space is required, `InnoDB` allocates pages to the tablespace and physically writes NULLs to those pages. This behavior can affect performance if new pages are allocated frequently. As of MySQL 8.0.22, you can disable `innodb_extend_and_initialize` on Linux systems to avoid physically writing NULLs to newly allocated tablespace pages. When `innodb_extend_and_initialize` is disabled, space is allocated to tablespace files using `posix_fallocate()` calls, which reserve space without physically writing NULLs.

When pages are allocated using `posix_fallocate()` calls, the extension size is small by default and pages are often allocated only a few at a time, which can cause fragmentation and increase random I/O. To avoid this issue, increase the tablespace extension size when enabling `posix_fallocate()` calls. Tablespace extension size can be increased up to 4GB using the `AUTOEXTEND_SIZE` option. For more information, see [Section 15.6.3.9, “Tablespace AUTOEXTEND_SIZE Configuration”](#).

`InnoDB` writes a redo log record before allocating a new tablespace page. If a page allocation operation is interrupted, the operation is replayed from the redo log record during recovery. (A page allocation operation replayed from a redo log record physically writes NULLs to the newly allocated page.) A redo log record is written before allocating a page regardless of the `innodb_extend_and_initialize` setting.

On non-Linux systems and Windows, `InnoDB` allocates new pages to the tablespace and physically writes NULLs to those pages, which is the default behavior. Attempting to disable `innodb_extend_and_initialize` on those systems returns the following error:

```
Changing innodb_extend_and_initialize not supported on this platform.  
Falling back to the default.
```

15.6.3.9 Tablespace AUTOEXTEND_SIZE Configuration

By default, when a file-per-table or general tablespace requires additional space, the tablespace is extended incrementally according to the following rules:

- If the tablespace is less than an extent in size, it is extended one page at a time.
- If the tablespace is greater than 1 extent but smaller than 32 extents in size, it is extended one extent at a time.
- If the tablespace is more than 32 extents in size, it is extended four extents at a time.

For information about extent size, see [Section 15.11.2, “File Space Management”](#).

From MySQL 8.0.23, the amount by which a file-per-table or general tablespace is extended is configurable by specifying the `AUTOEXTEND_SIZE` option. Configuring a larger extension size can help avoid fragmentation and facilitate ingestion of large amounts of data.

To configure the extension size for a file-per-table tablespace, specify the `AUTOEXTEND_SIZE` size in a `CREATE TABLE` or `ALTER TABLE` statement:

```
CREATE TABLE t1 (c1 INT) AUTOEXTEND_SIZE = 4M;  
ALTER TABLE t1 AUTOEXTEND_SIZE = 8M;
```

To configure the extension size for a general tablespace, specify the `AUTOEXTEND_SIZE` size in a `CREATE TABLESPACE` or `ALTER TABLESPACE` statement:

```
CREATE TABLESPACE ts1 AUTOEXTEND_SIZE = 4M;  
ALTER TABLESPACE ts1 AUTOEXTEND_SIZE = 8M;
```

**Note**

The `AUTOEXTEND_SIZE` option can also be used when creating an undo tablespace, but the extension behavior for undo tablespaces differs. For more information, see [Section 15.6.3.4, “Undo Tablespaces”](#).

The `AUTOEXTEND_SIZE` setting must be a multiple of 4M. Specifying an `AUTOEXTEND_SIZE` setting that is not a multiple of 4M returns an error.

The `AUTOEXTEND_SIZE` default setting is 0, which causes the tablespace to be extended according to the default behavior described above.

The maximum `AUTOEXTEND_SIZE` setting is 64M in MySQL 8.0.23. From MySQL 8.0.24, the maximum setting is 4GB.

The minimum `AUTOEXTEND_SIZE` setting depends on the `InnoDB` page size, as shown in the following table:

InnoDB Page Size	Minimum AUTOEXTEND_SIZE
4K	4M
8K	4M
16K	4M
32K	8M
64K	16M

The default `InnoDB` page size is 16K (16384 bytes). To determine the `InnoDB` page size for your MySQL instance, query the `innodb_page_size` setting:

```
mysql> SELECT @@GLOBAL.innodb_page_size;
+-----+
| @@GLOBAL.innodb_page_size |
+-----+
|          16384 |
+-----+
```

When the `AUTOEXTEND_SIZE` setting for a tablespace is altered, the first extension that occurs afterward increases the tablespace size to a multiple of the `AUTOEXTEND_SIZE` setting. Subsequent extensions are of the configured size.

When a file-per-table or general tablespace is created with a non-zero `AUTOEXTEND_SIZE` setting, the tablespace is initialized at the specified `AUTOEXTEND_SIZE` size.

`ALTER TABLESPACE` cannot be used to configure the `AUTOEXTEND_SIZE` of a file-per-table tablespace. `ALTER TABLE` must be used.

For tables created in file-per-table tablespaces, `SHOW CREATE TABLE` shows the `AUTOEXTEND_SIZE` option only when it is configured to a non-zero value.

To determine the `AUTOEXTEND_SIZE` for any `InnoDB` tablespace, query the Information Schema `INNODB_TABLESPACES` table. For example:

```
mysql> SELECT NAME, AUTOEXTEND_SIZE FROM INFORMATION_SCHEMA.INNODB_TABLESPACES
      WHERE NAME LIKE 'test/t1';
+-----+-----+
| NAME   | AUTOEXTEND_SIZE |
+-----+-----+
| test/t1 |        4194304 |
+-----+-----+

mysql> SELECT NAME, AUTOEXTEND_SIZE FROM INFORMATION_SCHEMA.INNODB_TABLESPACES
      WHERE NAME LIKE 'ts1';
+-----+-----+
```

NAME	AUTOEXTEND_SIZE
ts1	4194304

**Note**

An `AUTOEXTEND_SIZE` of 0, which is the default setting, means that the tablespace is extended according to the default tablespace extension behavior described above.

15.6.4 Doublewrite Buffer

The doublewrite buffer is a storage area where `InnoDB` writes pages flushed from the buffer pool before writing the pages to their proper positions in the `InnoDB` data files. If there is an operating system, storage subsystem, or unexpected `mysqld` process exit in the middle of a page write, `InnoDB` can find a good copy of the page from the doublewrite buffer during crash recovery.

Although data is 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 doublewrite buffer in a large sequential chunk, with a single `fsync()` call to the operating system (except in the case that `innodb_flush_method` is set to `O_DIRECT_NO_FSYNC`).

Prior to MySQL 8.0.20, the doublewrite buffer storage area is located in the `InnoDB` system tablespace. As of MySQL 8.0.20, the doublewrite buffer storage area is located in doublewrite files.

The following variables are provided for doublewrite buffer configuration:

- `innodb_doublewrite`

The `innodb_doublewrite` variable controls whether the doublewrite buffer is enabled. It is enabled by default in most cases. To disable the doublewrite buffer, set `innodb_doublewrite` to `OFF`. Consider disabling the doublewrite buffer if you are more concerned with performance than data integrity, as may be the case when performing benchmarks, for example.

From MySQL 8.0.30, `innodb_doublewrite` supports `DETECT_AND_RECOVER` and `DETECT_ONLY` settings.

The `DETECT_AND_RECOVER` setting is the same as the `ON` setting. With this setting, the doublewrite buffer is fully enabled, with database page content written to the doublewrite buffer where it is accessed during recovery to fix incomplete page writes.

With the `DETECT_ONLY` setting, only metadata is written to the doublewrite buffer. Database page content is not written to the doublewrite buffer, and recovery does not use the doublewrite buffer to fix incomplete page writes. This lightweight setting is intended for detecting incomplete page writes only.

MySQL 8.0.30 onwards supports dynamic changes to the `innodb_doublewrite` setting that enables the doublewrite buffer, between `ON`, `DETECT_AND_RECOVER`, and `DETECT_ONLY`. MySQL does not support dynamic changes between a setting that enables the doublewrite buffer and `OFF` or vice versa.

If the doublewrite buffer is located on a Fusion-io device that supports atomic writes, the doublewrite buffer is automatically disabled and data file writes are performed using Fusion-io atomic writes instead. However, be aware that the `innodb_doublewrite` setting is global. When the doublewrite buffer is disabled, it is disabled for all data files including those that do not reside on Fusion-io hardware. This feature is only supported on Fusion-io hardware and is only enabled for Fusion-io NVMFS on Linux. To take full advantage of this feature, an `innodb_flush_method` setting of `O_DIRECT` is recommended.

- `innodb_doublewrite_dir`

The `innodb_doublewrite_dir` variable (introduced in MySQL 8.0.20) defines the directory where InnoDB creates doublewrite files. If no directory is specified, doublewrite files are created in the `innodb_data_home_dir` directory, which defaults to the data directory if unspecified.

A hash symbol '#' is automatically prefixed to the specified directory name to avoid conflicts with schema names. However, if a '.', '#'. or '/' prefix is specified explicitly in the directory name, the hash symbol '#' is not prefixed to the directory name.

Ideally, the doublewrite directory should be placed on the fastest storage media available.

- `innodb_doublewrite_files`

The `innodb_doublewrite_files` variable defines the number of doublewrite files. By default, two doublewrite files are created for each buffer pool instance: A flush list doublewrite file and an LRU list doublewrite file.

The flush list doublewrite file is for pages flushed from the buffer pool flush list. The default size of a flush list doublewrite file is the InnoDB page size * doublewrite page bytes.

The LRU list doublewrite file is for pages flushed from the buffer pool LRU list. It also contains slots for single page flushes. The default size of an LRU list doublewrite file is the InnoDB page size * (doublewrite pages + (512 / the number of buffer pool instances)) where 512 is the total number of slots reserved for single page flushes.

At a minimum, there are two doublewrite files. The maximum number of doublewrite files is two times the number of buffer pool instances. (The number of buffer pool instances is controlled by the `innodb_buffer_pool_instances` variable.)

Doublewrite file names have the following format: `#ib_page_size_file_number dblwr` (or `.dblwr` with the `DETECT_ONLY` setting). For example, the following doublewrite files are created for a MySQL instance with an InnoDB pages size of 16KB and a single buffer pool:

```
#ib_16384_0.dblwr  
#ib_16384_1.dblwr
```

The `innodb_doublewrite_files` variable is intended for advanced performance tuning. The default setting should be suitable for most users.

- `innodb_doublewrite_pages`

The `innodb_doublewrite_pages` variable (introduced in MySQL 8.0.20) controls the maximum number of doublewrite pages per thread. If no value is specified, `innodb_doublewrite_pages` is set to the `innodb_write_io_threads` value. This variable is intended for advanced performance tuning. The default value should be suitable for most users.

- `innodb_doublewrite_batch_size`

The `innodb_doublewrite_batch_size` variable (introduced in MySQL 8.0.20) controls the number of doublewrite pages to write in a batch. This variable is intended for advanced performance tuning. The default value should be suitable for most users.

As of MySQL 8.0.23, InnoDB automatically encrypts doublewrite file pages that belong to encrypted tablespaces (see [Section 15.13, “InnoDB Data-at-Rest Encryption”](#)). Likewise, doublewrite file pages belonging to page-compressed tablespaces are compressed. As a result, doublewrite files can contain different page types including unencrypted and uncompressed pages, encrypted pages, compressed pages, and pages that are both encrypted and compressed.

15.6.5 Redo Log

The redo log is a disk-based data structure used during crash recovery to correct data written by incomplete transactions. During normal operations, the redo log encodes requests to change table

data that result from SQL statements or low-level API calls. Modifications that did not finish updating data files before an unexpected shutdown are replayed automatically during initialization and before connections are accepted. For information about the role of the redo log in crash recovery, see [Section 15.18.2, “InnoDB Recovery”](#).

The redo log is physically represented on disk by redo log files. Data that is written to redo log files is encoded in terms of records affected, and this data is collectively referred to as redo. The passage of data through redo log files is represented by an ever-increasing [LSN](#) value. Redo log data is appended as data modifications occur, and the oldest data is truncated as the checkpoint progresses.

Information and procedures related to redo logs are described under the following topics in the section:

- [Configuring Redo Log Capacity \(MySQL 8.0.30 or Higher\)](#)
- [Configuring Redo Log Capacity \(Before MySQL 8.0.30\)](#)
- [Automatic Redo Log Capacity Configuration](#)
- [Redo Log Archiving](#)
- [Disabling Redo Logging](#)
- [Related Topics](#)

Configuring Redo Log Capacity (MySQL 8.0.30 or Higher)

From MySQL 8.0.30, the `innodb_redo_log_capacity` system variable controls the amount of disk space occupied by redo log files. You can set this variable in an option file at startup or at runtime using a `SET GLOBAL` statement; for example, the following statement sets the redo log capacity to 8GB:

```
SET GLOBAL innodb_redo_log_capacity = 8589934592;
```

When set at runtime, the configuration change occurs immediately but it may take some time for the new limit to be fully implemented. If the redo log files occupy less space than the specified value, dirty pages are flushed from the buffer pool to tablespace data files less aggressively, eventually increasing the disk space occupied by the redo log files. If the redo log files occupy more space than the specified value, dirty pages are flushed more aggressively, eventually decreasing the disk space occupied by redo log files.

The `innodb_redo_log_capacity` variable supersedes the `innodb_log_files_in_group` and `innodb_log_file_size` variables, which are deprecated. When the `innodb_redo_log_capacity` setting is defined, the `innodb_log_files_in_group` and `innodb_log_file_size` settings are ignored; otherwise, these settings are used to compute the `innodb_redo_log_capacity` setting (`innodb_log_files_in_group * innodb_log_file_size = innodb_redo_log_capacity`). If none of those variables are set, redo log capacity is set to the `innodb_redo_log_capacity` default value, which is 104857600 bytes (100MB). The maximum redo log capacity is 128GB.

Redo log files reside in the `#innodb_redo` directory in the data directory unless a different directory was specified by the `innodb_log_group_home_dir` variable. If `innodb_log_group_home_dir` was defined, the redo log files reside in the `#innodb_redo` directory in that directory. There are two types of redo log files, ordinary and spare. Ordinary redo log files are those being used. Spare redo log files are those waiting to be used. InnoDB tries to maintain 32 redo log files in total, with each file equal in size to $1/32 * \text{innodb_redo_log_capacity}$; however, file sizes may differ for a time after modifying the `innodb_redo_log_capacity` setting.

Redo log files use an `#ib_redoN` naming convention, where N is the redo log file number. Spare redo log files are denoted by a `_tmp` suffix. The following example shows the redo log files in an `#innodb_redo` directory, where there are 21 active redo log files and 11 spare redo log files, numbered sequentially.

```
'#ib_redo582' '#ib_redo590' '#ib_redo598' '#ib_redo606_tmp'
```

```
'#ib_redo583' '#ib_redo591' '#ib_redo599' '#ib_redo607_tmp'
'#ib_redo584' '#ib_redo592' '#ib_redo600' '#ib_redo608_tmp'
'#ib_redo585' '#ib_redo593' '#ib_redo601' '#ib_redo609_tmp'
'#ib_redo586' '#ib_redo594' '#ib_redo602' '#ib_redo610_tmp'
'#ib_redo587' '#ib_redo595' '#ib_redo603_tmp' '#ib_redo611_tmp'
'#ib_redo588' '#ib_redo596' '#ib_redo604_tmp' '#ib_redo612_tmp'
'#ib_redo589' '#ib_redo597' '#ib_redo605_tmp' '#ib_redo613_tmp'
```

Each ordinary redo log file is associated with a particular range of LSN values; for example, the following query shows the `START_LSN` and `END_LSN` values for the active redo log files listed in the previous example:

```
mysql> SELECT FILE_NAME, START_LSN, END_LSN FROM performance_schema.innodb_redo_log_files;
+-----+-----+-----+
| FILE_NAME | START_LSN | END_LSN |
+-----+-----+-----+
| ./innodb_redo/#ib_redo582 | 117654982144 | 117658256896 |
| ./innodb_redo/#ib_redo583 | 117658256896 | 117661531648 |
| ./innodb_redo/#ib_redo584 | 117661531648 | 117664806400 |
| ./innodb_redo/#ib_redo585 | 117664806400 | 117668081152 |
| ./innodb_redo/#ib_redo586 | 117668081152 | 117671355904 |
| ./innodb_redo/#ib_redo587 | 117671355904 | 117674630656 |
| ./innodb_redo/#ib_redo588 | 117674630656 | 117677905408 |
| ./innodb_redo/#ib_redo589 | 117677905408 | 117681180160 |
| ./innodb_redo/#ib_redo590 | 117681180160 | 117684454912 |
| ./innodb_redo/#ib_redo591 | 117684454912 | 117687729664 |
| ./innodb_redo/#ib_redo592 | 117687729664 | 117691004416 |
| ./innodb_redo/#ib_redo593 | 117691004416 | 117694279168 |
| ./innodb_redo/#ib_redo594 | 117694279168 | 117697553920 |
| ./innodb_redo/#ib_redo595 | 117697553920 | 117700828672 |
| ./innodb_redo/#ib_redo596 | 117700828672 | 117704103424 |
| ./innodb_redo/#ib_redo597 | 117704103424 | 117707378176 |
| ./innodb_redo/#ib_redo598 | 117707378176 | 117710652928 |
| ./innodb_redo/#ib_redo599 | 117710652928 | 117713927680 |
| ./innodb_redo/#ib_redo600 | 117713927680 | 117717202432 |
| ./innodb_redo/#ib_redo601 | 117717202432 | 117720477184 |
| ./innodb_redo/#ib_redo602 | 117720477184 | 117723751936 |
+-----+-----+-----+
```

When doing a checkpoint, InnoDB stores the checkpoint LSN in the header of the file which contains this LSN. During recovery, all redo log files are checked and recovery starts at the latest checkpoint LSN.

Several status variables are provided for monitoring the redo log and redo log capacity resize operations; for example, you can query `Innodb_redo_log_resize_status` to view the status of a resize operation:

```
mysql> SHOW STATUS LIKE 'Innodb_redo_log_resize_status';
+-----+-----+
| Variable_name | Value |
+-----+-----+
| Innodb_redo_log_resize_status | OK |
+-----+-----+
```

The `Innodb_redo_log_capacity_resized` status variable shows the current redo log capacity limit:

```
mysql> SHOW STATUS LIKE 'Innodb_redo_log_capacity_resized';
+-----+-----+
| Variable_name | Value |
+-----+-----+
| Innodb_redo_log_capacity_resized | 104857600 |
+-----+-----+
```

Other applicable status variables include:

- `Innodb_redo_log_checkpoint_lsn`
- `Innodb_redo_log_current_lsn`

- `Innodb_redo_log_flushed_to_disk_lsn`
- `Innodb_redo_log_logical_size`
- `Innodb_redo_log_physical_size`
- `Innodb_redo_log_read_only`
- `Innodb_redo_log_uuid`

Refer to the status variable descriptions for more information.

You can view information about active redo log files by querying the `innodb_redo_log_files` Performance Schema table. The following query retrieves data from all of the table's columns:

```
SELECT FILE_ID, START_LSN, END_LSN, SIZE_IN_BYTES, IS_FULL, CONSUMER_LEVEL  
FROM performance_schema.innodb_redo_log_files;
```

Configuring Redo Log Capacity (Before MySQL 8.0.30)

Prior to MySQL 8.0.30, `InnoDB` creates two redo log files in the data directory by default, named `ib_logfile0` and `ib_logfile1`, and writes to these files in a circular fashion.

Modifying redo log capacity requires changing the number or the size of `redo log` files, or both.

1. Stop the MySQL server and make sure that it shuts down without errors.
2. Edit `my.cnf` to change the redo log file configuration. To change the redo log file size, configure `innodb_log_file_size`. To increase the number of redo log files, configure `innodb_log_files_in_group`.
3. Start the MySQL server again.

If `InnoDB` detects that the `innodb_log_file_size` differs from the redo log file size, it writes a log checkpoint, closes and removes the old log files, creates new log files at the requested size, and opens the new log files.

Automatic Redo Log Capacity Configuration

When `innodb_dedicated_server` is enabled, `InnoDB` automatically configures certain `InnoDB` parameters, including redo log capacity. Automated configuration is intended for MySQL instances that reside on a server dedicated to MySQL, where the MySQL server can use all available system resources. For more information, see [Section 15.8.12, “Enabling Automatic Configuration for a Dedicated MySQL Server”](#).

Redo Log Archiving

Backup utilities that copy redo log records may sometimes fail to keep pace with redo log generation while a backup operation is in progress, resulting in lost redo log records due to those records being overwritten. This issue most often occurs when there is significant MySQL server activity during the backup operation, and the redo log file storage media operates at a faster speed than the backup storage media. The redo log archiving feature, introduced in MySQL 8.0.17, addresses this issue by sequentially writing redo log records to an archive file in addition to the redo log files. Backup utilities can copy redo log records from the archive file as necessary, thereby avoiding the potential loss of data.

If redo log archiving is configured on the server, [MySQL Enterprise Backup](#), available with the [MySQL Enterprise Edition](#), uses the redo log archiving feature when backing up a MySQL server.

Enabling redo log archiving on the server requires setting a value for the `innodb_redo_log_archive_dirs` system variable. The value is specified as a semicolon-separated list of labeled redo log archive directories. The `label:directory` pair is separated by a colon (`:`). For example:

```
mysql> SET GLOBAL innodb_redo_log_archive_dirs='label1:directory_path1[;label2:directory_path2;...]';
```

The `label` is an arbitrary identifier for the archive directory. It can be any string of characters, with the exception of colons (:), which are not permitted. An empty label is also permitted, but the colon (:) is still required in this case. A `directory_path` must be specified. The directory selected for the redo log archive file must exist when redo log archiving is activated, or an error is returned. The path can contain colons (':'), but semicolons (;) are not permitted.

The `innodb_redo_log_archive_dirs` variable must be configured before redo log archiving can be activated. The default value is `NULL`, which does not permit activating redo log archiving.



Notes

The archive directories that you specify must satisfy the following requirements. (The requirements are enforced when redo log archiving is activated.):

- Directories must exist. Directories are not created by the redo log archive process. Otherwise, the following error is returned:

```
ERROR 3844 (HY000): Redo log archive directory
'directory_path1' does not exist or is not a directory
```

- Directories must not be world-accessible. This is to prevent the redo log data from being exposed to unauthorized users on the system. Otherwise, the following error is returned:

```
ERROR 3846 (HY000): Redo log archive directory
'directory_path1' is accessible to all OS users
```

- Directories cannot be those defined by `datadir`, `innodb_data_home_dir`, `innodb_directories`, `innodb_log_group_home_dir`, `innodb_temp_tablespaces_dir`, `innodb_tmpdir`, `innodb_undo_directory`, or `secure_file_priv`, nor can they be parent directories or subdirectories of those directories. Otherwise, an error similar to the following is returned:

```
ERROR 3845 (HY000): Redo log archive directory
'directory_path1' is in, under, or over server directory
'datadir' - '/path/to/data_directory'
```

When a backup utility that supports redo log archiving initiates a backup, the backup utility activates redo log archiving by invoking the `innodb_redo_log_archive_start()` function.

If you are not using a backup utility that supports redo log archiving, redo log archiving can also be activated manually, as shown:

```
mysql> SELECT innodb_redo_log_archive_start('label', 'subdir');
+-----+
| innodb_redo_log_archive_start('label') |
+-----+
| 0                                         |
+-----+
```

Or:

```
mysql> DO innodb_redo_log_archive_start('label', 'subdir');
Query OK, 0 rows affected (0.09 sec)
```



Note

The MySQL session that activates redo log archiving (using `innodb_redo_log_archive_start()`) must remain open for the duration of the archiving. The same session must deactivate redo log archiving (using

`innodb_redo_log_archive_stop()`). If the session is terminated before the redo log archiving is explicitly deactivated, the server deactivates redo log archiving implicitly and removes the redo log archive file.

where `label` is a label defined by `innodb_redo_log_archive_dirs`; `subdir` is an optional argument for specifying a subdirectory of the directory identified by `label` for saving the archive file; it must be a simple directory name (no slash (/), backslash (\), or colon (:) is permitted). `subdir` can be empty, null, or it can be left out.

Only users with the `INNODB_REDO_LOG_ARCHIVE` privilege can activate redo log archiving by invoking `innodb_redo_log_archive_start()`, or deactivate it using `innodb_redo_log_archive_stop()`. The MySQL user running the backup utility or the MySQL user activating and deactivating redo log archiving manually must have this privilege.

The redo log archive file path is `directory_identified_by_label/[subdir/]archive.serverUUID.000001.log`, where `directory_identified_by_label` is the archive directory identified by the `label` argument for `innodb_redo_log_archive_start()`. `subdir` is the optional argument used for `innodb_redo_log_archive_start()`.

For example, the full path and name for a redo log archive file appears similar to the following:

```
/directory_path/subdirectory/archive.e71a47dc-61f8-11e9-a3cb-080027154b4d.000001.log
```

After the backup utility finishes copying InnoDB data files, it deactivates redo log archiving by calling the `innodb_redo_log_archive_stop()` function.

If you are not using a backup utility that supports redo log archiving, redo log archiving can also be deactivated manually, as shown:

```
mysql> SELECT innodb_redo_log_archive_stop();
+-----+
| innodb_redo_log_archive_stop() |
+-----+
| 0                                |
+-----+
```

Or:

```
mysql> DO innodb_redo_log_archive_stop();
Query OK, 0 rows affected (0.01 sec)
```

After the stop function completes successfully, the backup utility looks for the relevant section of redo log data from the archive file and copies it into the backup.

After the backup utility finishes copying the redo log data and no longer needs the redo log archive file, it deletes the archive file.

Removal of the archive file is the responsibility of the backup utility in normal situations. However, if the redo log archiving operation quits unexpectedly before `innodb_redo_log_archive_stop()` is called, the MySQL server removes the file.

Performance Considerations

Activating redo log archiving typically has a minor performance cost due to the additional write activity.

On Unix and Unix-like operating systems, the performance impact is typically minor, assuming there is not a sustained high rate of updates. On Windows, the performance impact is typically a bit higher, assuming the same.

If there is a sustained high rate of updates and the redo log archive file is on the same storage media as the redo log files, the performance impact may be more significant due to compounded write activity.

If there is a sustained high rate of updates and the redo log archive file is on slower storage media than the redo log files, performance is impacted arbitrarily.

Writing to the redo log archive file does not impede normal transactional logging except in the case that the redo log archive file storage media operates at a much slower rate than the redo log file storage media, and there is a large backlog of persisted redo log blocks waiting to be written to the redo log archive file. In this case, the transactional logging rate is reduced to a level that can be managed by the slower storage media where the redo log archive file resides.

Disabling Redo Logging

As of MySQL 8.0.21, you can disable redo logging using the `ALTER INSTANCE DISABLE INNODB REDO_LOG` statement. This functionality is intended for loading data into a new MySQL instance. Disabling redo logging speeds up data loading by avoiding redo log writes and doublewrite buffering.



Warning

This feature is intended only for loading data into a new MySQL instance. *Do not disable redo logging on a production system.* It is permitted to shutdown and restart the server while redo logging is disabled, but an unexpected server stoppage while redo logging is disabled can cause data loss and instance corruption.

Attempting to restart the server after an unexpected server stoppage while redo logging is disabled is refused with the following error:

```
[ERROR] [MY-013598] [InnoDB] Server was killed when Innodb Redo
logging was disabled. Data files could be corrupt. You can try
to restart the database with innodb_force_recovery=6
```

In this case, initialize a new MySQL instance and start the data loading procedure again.

The `INNODB_REDOLG_ENABLE` privilege is required to enable and disable redo logging.

The `Innodb_redo_log_enabled` status variable permits monitoring redo logging status.

Cloning operations and redo log archiving are not permitted while redo logging is disabled and vice versa.

An `ALTER INSTANCE [ENABLE|DISABLE] INNODB REDO_LOG` operation requires an exclusive backup metadata lock, which prevents other `ALTER INSTANCE` operations from executing concurrently. Other `ALTER INSTANCE` operations must wait for the lock to be released before executing.

The following procedure demonstrates how to disable redo logging when loading data into a new MySQL instance.

- On the new MySQL instance, grant the `INNODB_REDOLG_ENABLE` privilege to the user account responsible for disabling redo logging.

```
mysql> GRANT INNODB_REDOLG_ENABLE ON *.* to 'data_load_admin';
```

- As the `data_load_admin` user, disable redo logging:

```
mysql> ALTER INSTANCE DISABLE INNODB REDO_LOG;
```

- Check the `Innodb_redo_log_enabled` status variable to ensure that redo logging is disabled.

```
mysql> SHOW GLOBAL STATUS LIKE 'Innodb_redo_log_enabled';
+-----+-----+
| Variable_name      | Value   |
+-----+-----+
| Innodb_redo_log_enabled | OFF    |
+-----+-----+
```

- Run the data load operation.

5. As the `data_load_admin` user, enable redo logging after the data load operation finishes:

```
mysql> ALTER INSTANCE ENABLE INNODB REDO_LOG;
```

6. Check the `Innodb_redo_log_enabled` status variable to ensure that redo logging is enabled.

```
mysql> SHOW GLOBAL STATUS LIKE 'Innodb_redo_log_enabled';
+-----+-----+
| Variable_name      | Value   |
+-----+-----+
| Innodb_redo_log_enabled | ON     |
+-----+-----+
```

Related Topics

- [Redo Log Configuration](#)
- [Section 8.5.4, “Optimizing InnoDB Redo Logging”](#)
- [Redo Log Encryption](#)

15.6.6 Undo Logs

An undo log is a collection of undo log records associated with a single read-write transaction. An undo log record contains information about how to undo the latest change by a transaction to a [clustered index](#) record. If another transaction needs to see the original data as part of a consistent read operation, the unmodified data is retrieved from undo log records. Undo logs exist within [undo log segments](#), which are contained within [rollback segments](#). Rollback segments reside in [undo tablespaces](#) and in the [global temporary tablespace](#).

Undo logs that reside in the global temporary tablespace are used for transactions that modify data in user-defined temporary tables. These undo logs are not redo-logged, as they are not required for crash recovery. They are used only for rollback while the server is running. This type of undo log benefits performance by avoiding redo logging I/O.

For information about data-at-rest encryption for undo logs, see [Undo Log Encryption](#).

Each undo tablespace and the global temporary tablespace individually support a maximum of 128 rollback segments. The `innodb_rollback_segments` variable defines the number of rollback segments.

The number of transactions that a rollback segment supports depends on the number of undo slots in the rollback segment and the number of undo logs required by each transaction. The number of undo slots in a rollback segment differs according to [InnoDB](#) page size.

InnoDB Page Size	Number of Undo Slots in a Rollback Segment (InnoDB Page Size / 16)
4096 (4KB)	256
8192 (8KB)	512
16384 (16KB)	1024
32768 (32KB)	2048
65536 (64KB)	4096

A transaction is assigned up to four undo logs, one for each of the following operation types:

1. `INSERT` operations on user-defined tables
2. `UPDATE` and `DELETE` operations on user-defined tables
3. `INSERT` operations on user-defined temporary tables

4. `UPDATE` and `DELETE` operations on user-defined temporary tables

Undo logs are assigned as needed. For example, a transaction that performs `INSERT`, `UPDATE`, and `DELETE` operations on regular and temporary tables requires a full assignment of four undo logs. A transaction that performs only `INSERT` operations on regular tables requires a single undo log.

A transaction that performs operations on regular tables is assigned undo logs from an assigned undo tablespace rollback segment. A transaction that performs operations on temporary tables is assigned undo logs from an assigned global temporary tablespace rollback segment.

An undo log assigned to a transaction remains attached to the transaction for its duration. For example, an undo log assigned to a transaction for an `INSERT` operation on a regular table is used for all `INSERT` operations on regular tables performed by that transaction.

Given the factors described above, the following formulas can be used to estimate the number of concurrent read-write transactions that `InnoDB` is capable of supporting.



Note

It is possible to encounter a concurrent transaction limit error before reaching the number of concurrent read-write transactions that `InnoDB` is capable of supporting. This occurs when a rollback segment assigned to a transaction runs out of undo slots. In such cases, try rerunning the transaction.

When transactions perform operations on temporary tables, the number of concurrent read-write transactions that `InnoDB` is capable of supporting is constrained by the number of rollback segments allocated to the global temporary tablespace, which is 128 by default.

- If each transaction performs either an `INSERT` or an `UPDATE` or `DELETE` operation, the number of concurrent read-write transactions that `InnoDB` is capable of supporting is:

```
(innodb_page_size / 16) * innodb_rollback_segments * number of undo tablespaces
```

- If each transaction performs an `INSERT` and an `UPDATE` or `DELETE` operation, the number of concurrent read-write transactions that `InnoDB` is capable of supporting is:

```
(innodb_page_size / 16 / 2) * innodb_rollback_segments * number of undo tablespaces
```

- If each transaction performs an `INSERT` operation on a temporary table, the number of concurrent read-write transactions that `InnoDB` is capable of supporting is:

```
(innodb_page_size / 16) * innodb_rollback_segments
```

- If each transaction performs an `INSERT` and an `UPDATE` or `DELETE` operation on a temporary table, the number of concurrent read-write transactions that `InnoDB` is capable of supporting is:

```
(innodb_page_size / 16 / 2) * innodb_rollback_segments
```

15.7 InnoDB Locking and Transaction Model

To implement a large-scale, busy, or highly reliable database application, to port substantial code from a different database system, or to tune MySQL performance, it is important to understand `InnoDB` locking and the `InnoDB` transaction model.

This section discusses several topics related to `InnoDB` locking and the `InnoDB` transaction model with which you should be familiar.

- [Section 15.7.1, “InnoDB Locking”](#) describes lock types used by `InnoDB`.
- [Section 15.7.2, “InnoDB Transaction Model”](#) describes transaction isolation levels and the locking strategies used by each. It also discusses the use of `autocommit`, consistent non-locking reads, and locking reads.

- [Section 15.7.3, “Locks Set by Different SQL Statements in InnoDB”](#) discusses specific types of locks set in [InnoDB](#) for various statements.
- [Section 15.7.4, “Phantom Rows”](#) describes how [InnoDB](#) uses next-key locking to avoid phantom rows.
- [Section 15.7.5, “Deadlocks in InnoDB”](#) provides a deadlock example, discusses deadlock detection, and provides tips for minimizing and handling deadlocks in [InnoDB](#).

15.7.1 InnoDB Locking

This section describes lock types used by [InnoDB](#).

- [Shared and Exclusive Locks](#)
- [Intention Locks](#)
- [Record Locks](#)
- [Gap Locks](#)
- [Next-Key Locks](#)
- [Insert Intention Locks](#)
- [AUTO-INC Locks](#)
- [Predicate Locks for Spatial Indexes](#)

Shared and Exclusive Locks

[InnoDB](#) implements standard row-level locking where there are two types of locks, [shared \(S\)](#) locks and [exclusive \(X\)](#) locks.

- A [shared \(S\) lock](#) permits the transaction that holds the lock to read a row.
- An [exclusive \(X\) lock](#) permits the transaction that holds the lock to update or delete a row.

If transaction [T1](#) holds a shared ([S](#)) lock on row [r](#), then requests from some distinct transaction [T2](#) for a lock on row [r](#) are handled as follows:

- A request by [T2](#) for an [S](#) lock can be granted immediately. As a result, both [T1](#) and [T2](#) hold an [S](#) lock on [r](#).
- A request by [T2](#) for an [X](#) lock cannot be granted immediately.

If a transaction [T1](#) holds an exclusive ([X](#)) lock on row [r](#), a request from some distinct transaction [T2](#) for a lock of either type on [r](#) cannot be granted immediately. Instead, transaction [T2](#) has to wait for transaction [T1](#) to release its lock on row [r](#).

Intention Locks

[InnoDB](#) supports *multiple granularity locking* which permits coexistence of row locks and table locks. For example, a statement such as `LOCK TABLES ... WRITE` takes an exclusive lock (an [X](#) lock) on the specified table. To make locking at multiple granularity levels practical, [InnoDB](#) uses [intention locks](#). Intention locks are table-level locks that indicate which type of lock (shared or exclusive) a transaction requires later for a row in a table. There are two types of intention locks:

- An [intention shared lock \(IS\)](#) indicates that a transaction intends to set a [shared](#) lock on individual rows in a table.
- An [intention exclusive lock \(IX\)](#) indicates that a transaction intends to set an [exclusive](#) lock on individual rows in a table.

For example, `SELECT ... FOR SHARE` sets an `IS` lock, and `SELECT ... FOR UPDATE` sets an `IX` lock.

The intention locking protocol is as follows:

- Before a transaction can acquire a shared lock on a row in a table, it must first acquire an `IS` lock or stronger on the table.
- Before a transaction can acquire an exclusive lock on a row in a table, it must first acquire an `IX` lock on the table.

Table-level lock type compatibility is summarized in the following matrix.

	<code>X</code>	<code>IX</code>	<code>S</code>	<code>IS</code>
<code>X</code>	Conflict	Conflict	Conflict	Conflict
<code>IX</code>	Conflict	Compatible	Conflict	Compatible
<code>S</code>	Conflict	Conflict	Compatible	Compatible
<code>IS</code>	Conflict	Compatible	Compatible	Compatible

A lock is granted to a requesting transaction if it is compatible with existing locks, but not if it conflicts with existing locks. A transaction waits until the conflicting existing lock is released. If a lock request conflicts with an existing lock and cannot be granted because it would cause `deadlock`, an error occurs.

Intention locks do not block anything except full table requests (for example, `LOCK TABLES ... WRITE`). The main purpose of intention locks is to show that someone is locking a row, or going to lock a row in the table.

Transaction data for an intention lock appears similar to the following in `SHOW ENGINE INNODB STATUS` and `InnoDB monitor` output:

```
TABLE LOCK table `test`.`t` trx id 10080 lock mode IX
```

Record Locks

A record lock is a lock on an index record. For example, `SELECT c1 FROM t WHERE c1 = 10 FOR UPDATE;` prevents any other transaction from inserting, updating, or deleting rows where the value of `t.c1` is `10`.

Record locks always lock index records, even if a table is defined with no indexes. For such cases, InnoDB creates a hidden clustered index and uses this index for record locking. See [Section 15.6.2.1, "Clustered and Secondary Indexes"](#).

Transaction data for a record lock appears similar to the following in `SHOW ENGINE INNODB STATUS` and `InnoDB monitor` output:

```
RECORD LOCKS space id 58 page no 3 n bits 72 index `PRIMARY` of table `test`.`t`
trx id 10078 lock_mode X locks rec but not gap
Record lock, heap no 2 PHYSICAL RECORD: n_fields 3; compact format; info bits 0
 0: len 4; hex 8000000a; asc     ;;
 1: len 6; hex 00000000274f; asc     'O';
 2: len 7; hex b60000019d0110; asc     ;;
```

Gap Locks

A gap lock is 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 FROM t WHERE c1 BETWEEN 10 and 20 FOR UPDATE;` prevents other transactions from inserting a value of `15` into 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.

A gap might span a single index value, multiple index values, or even be empty.

Gap locks are part of the tradeoff between performance and concurrency, and are used in some transaction isolation levels and not others.

Gap locking is not needed for statements that lock rows using a unique index to search for a unique row. (This does not include the case that the search condition includes only some columns of a multiple-column unique index; in that case, gap locking does occur.) For example, if the `id` column has a unique index, the following statement uses only an index-record lock for the row having `id` value 100 and it does not matter whether other sessions insert rows in the preceding gap:

```
SELECT * FROM child WHERE id = 100;
```

If `id` is not indexed or has a nonunique index, the statement does lock the preceding gap.

It is also worth noting here that conflicting locks can be held on a gap by different transactions. For example, transaction A can hold a shared gap lock (gap S-lock) on a gap while transaction B holds an exclusive gap lock (gap X-lock) on the same gap. The reason conflicting gap locks are allowed is that if a record is purged from an index, the gap locks held on the record by different transactions must be merged.

Gap locks in InnoDB are “purely inhibitive”, which means that their only purpose is to prevent other transactions from inserting to the gap. Gap locks can co-exist. A gap lock taken by one transaction does not prevent another transaction from taking a gap lock on the same gap. There is no difference between shared and exclusive gap locks. They do not conflict with each other, and they perform the same function.

Gap locking can be disabled explicitly. This occurs if you change the transaction isolation level to `READ COMMITTED`. In this case, gap locking is disabled for searches and index scans and is used only for foreign-key constraint checking and duplicate-key checking.

There are also other effects of using the `READ COMMITTED` isolation level. Record locks for nonmatching rows are released after MySQL has evaluated the `WHERE` condition. For `UPDATE` statements, InnoDB does a “semi-consistent” read, such that it returns the latest committed version to MySQL so that MySQL can determine whether the row matches the `WHERE` condition of the `UPDATE`.

Next-Key Locks

A next-key lock is a combination of a record lock on the index record and a gap lock on the gap before the index record.

InnoDB performs row-level locking in such a way that when it searches or scans a table index, it sets shared or exclusive locks on the index records it encounters. Thus, the row-level locks are actually index-record locks. A next-key lock on an index record also affects the “gap” before that index record. That is, a next-key lock is an index-record lock plus a gap lock on the gap preceding the index record. If one session has a shared or exclusive lock on record `R` in an index, another session cannot insert a new index record in the gap immediately before `R` in the index order.

Suppose that an index contains the values 10, 11, 13, and 20. The possible next-key locks for this index cover the following intervals, where a round bracket denotes exclusion of the interval endpoint and a square bracket denotes inclusion of the endpoint:

```
(negative infinity, 10]
(10, 11]
(11, 13]
(13, 20]
(20, positive infinity)
```

For the last interval, the next-key lock locks the gap above the largest value in the index and the “supremum” pseudo-record having a value higher than any value actually in the index. The supremum is not a real index record, so, in effect, this next-key lock locks only the gap following the largest index value.

By default, InnoDB operates in REPEATABLE READ transaction isolation level. In this case, InnoDB uses next-key locks for searches and index scans, which prevents phantom rows (see [Section 15.7.4, “Phantom Rows”](#)).

Transaction data for a next-key lock appears similar to the following in `SHOW ENGINE INNODB STATUS` and [InnoDB monitor](#) output:

```
RECORD LOCKS space id 58 page no 3 n bits 72 index `PRIMARY` of table `test`.`t`
trx id 10080 lock_mode X
Record lock, heap no 1 PHYSICAL RECORD: n_fields 1; compact format; info bits 0
 0: len 8; hex 73757072656d756d; asc supremum;;
Record lock, heap no 2 PHYSICAL RECORD: n_fields 3; compact format; info bits 0
 0: len 4; hex 8000000a; asc      ;;
 1: len 6; hex 00000000274f; asc      '0;;
 2: len 7; hex b60000019d0110; asc      ;;
```

Insert Intention Locks

An insert intention lock is a type of gap lock set by `INSERT` operations prior to row insertion. This lock signals the intent to insert in such a way that multiple transactions inserting into the same index gap need not wait for each other if they are not inserting at the same position within the gap. Suppose that there are index records with values of 4 and 7. Separate transactions that attempt to insert values of 5 and 6, respectively, each lock the gap between 4 and 7 with insert intention locks prior to obtaining the exclusive lock on the inserted row, but do not block each other because the rows are nonconflicting.

The following example demonstrates a transaction taking an insert intention lock prior to obtaining an exclusive lock on the inserted record. The example involves two clients, A and B.

Client A creates a table containing two index records (90 and 102) and then starts a transaction that places an exclusive lock on index records with an ID greater than 100. The exclusive lock includes a gap lock before record 102:

```
mysql> CREATE TABLE child (id int(11) NOT NULL, PRIMARY KEY(id)) ENGINE=InnoDB;
mysql> INSERT INTO child (id) values (90),(102);

mysql> START TRANSACTION;
mysql> SELECT * FROM child WHERE id > 100 FOR UPDATE;
+----+
| id |
+----+
| 102 |
+----+
```

Client B begins a transaction to insert a record into the gap. The transaction takes an insert intention lock while it waits to obtain an exclusive lock.

```
mysql> START TRANSACTION;
mysql> INSERT INTO child (id) VALUES (101);
```

Transaction data for an insert intention lock appears similar to the following in `SHOW ENGINE INNODB STATUS` and [InnoDB monitor](#) output:

```
RECORD LOCKS space id 31 page no 3 n bits 72 index `PRIMARY` of table `test`.`child`
trx id 8731 lock_mode X locks gap before rec insert intention waiting
Record lock, heap no 3 PHYSICAL RECORD: n_fields 3; compact format; info bits 0
 0: len 4; hex 80000066; asc   f;;
 1: len 6; hex 000000002215; asc      " ;;
 2: len 7; hex 9000000172011c; asc      r  ;...;
```

AUTO-INC Locks

An AUTO-INC lock is a special table-level lock taken by transactions inserting into tables with `AUTO_INCREMENT` columns. 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.

The `innodb_autoinc_lock_mode` variable controls the algorithm used for auto-increment locking. It allows you to choose how to trade off between predictable sequences of auto-increment values and maximum concurrency for insert operations.

For more information, see [Section 15.6.1.6, “AUTO_INCREMENT Handling in InnoDB”](#).

Predicate Locks for Spatial Indexes

InnoDB supports `SPATIAL` indexing of columns containing spatial data (see [Section 11.4.9, “Optimizing Spatial Analysis”](#)).

To handle locking for operations involving `SPATIAL` indexes, next-key locking does not work well to support `REPEATABLE READ` or `SERIALIZABLE` transaction isolation levels. There is no absolute ordering concept in multidimensional data, so it is not clear which is the “next” key.

To enable support of isolation levels for tables with `SPATIAL` indexes, InnoDB uses predicate locks. A `SPATIAL` index contains minimum bounding rectangle (MBR) values, so InnoDB enforces consistent read on the index by setting a predicate lock on the MBR value used for a query. Other transactions cannot insert or modify a row that would match the query condition.

15.7.2 InnoDB Transaction Model

The InnoDB transaction model aims to combine the best properties of a `multi-versioning` database with traditional two-phase locking. InnoDB performs locking at the row level and runs queries as nonlocking `consistent reads` by default, in the style of Oracle. The lock information in InnoDB is stored space-efficiently so that lock escalation is not needed. Typically, several users are permitted to lock every row in InnoDB tables, or any random subset of the rows, without causing InnoDB memory exhaustion.

15.7.2.1 Transaction Isolation Levels

Transaction isolation is 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.

InnoDB offers all four transaction isolation levels described by the SQL:1992 standard: `READ UNCOMMITTED`, `READ COMMITTED`, `REPEATABLE READ`, and `SERIALIZABLE`. The default isolation level for InnoDB is `REPEATABLE READ`.

A user can change the isolation level for a single session or for all subsequent connections with the `SET TRANSACTION` statement. To set the server's default isolation level for all connections, use the `--transaction-isolation` option on the command line or in an option file. For detailed information about isolation levels and level-setting syntax, see [Section 13.3.7, “SET TRANSACTION Statement”](#).

InnoDB supports each of the transaction isolation levels described here using different `locking` strategies. You can enforce a high degree of consistency with the default `REPEATABLE READ` level, for operations on crucial data where `ACID` compliance is important. Or you can relax the consistency rules with `READ COMMITTED` or even `READ UNCOMMITTED`, in situations such as bulk reporting where precise consistency and repeatable results are less important than minimizing the amount of overhead for locking. `SERIALIZABLE` enforces even stricter rules than `REPEATABLE READ`, and is used mainly in specialized situations, such as with `X` transactions and for troubleshooting issues with concurrency and `deadlocks`.

The following list describes how MySQL supports the different transaction levels. The list goes from the most commonly used level to the least used.

- `REPEATABLE READ`

This is the default isolation level for InnoDB. `Consistent reads` within the same transaction read the `snapshot` established by the first read. This means that if you issue several plain (nonlocking) `SELECT` statements within the same transaction, these `SELECT` statements are consistent also with respect to each other. See [Section 15.7.2.3, “Consistent Nonlocking Reads”](#).

For locking reads (`SELECT` with `FOR UPDATE` or `FOR SHARE`), `UPDATE`, and `DELETE` statements, locking depends on whether the statement uses a unique index with a unique search condition, or a range-type search condition.

- For a unique index with a unique search condition, InnoDB locks only the index record found, not the gap before it.
- For other search conditions, InnoDB locks the index range scanned, using **gap locks** or **next-key locks** to block insertions by other sessions into the gaps covered by the range. For information about gap locks and next-key locks, see [Section 15.7.1, “InnoDB Locking”](#).
- `READ COMMITTED`

Each consistent read, even within the same transaction, sets and reads its own fresh snapshot. For information about consistent reads, see [Section 15.7.2.3, “Consistent Nonlocking Reads”](#).

For locking reads (`SELECT` with `FOR UPDATE` or `FOR SHARE`), `UPDATE` statements, and `DELETE` statements, InnoDB locks only index records, not the gaps before them, and thus permits the free insertion of new records next to locked records. Gap locking is only used for foreign-key constraint checking and duplicate-key checking.

Because gap locking is disabled, phantom row problems may occur, as other sessions can insert new rows into the gaps. For information about phantom rows, see [Section 15.7.4, “Phantom Rows”](#).

Only row-based binary logging is supported with the `READ COMMITTED` isolation level. If you use `READ COMMITTED` with `binlog_format=MIXED`, the server automatically uses row-based logging.

Using `READ COMMITTED` has additional effects:

- For `UPDATE` or `DELETE` statements, InnoDB holds locks only for rows that it updates or deletes. Record locks for nonmatching rows are released after MySQL has evaluated the `WHERE` condition. This greatly reduces the probability of deadlocks, but they can still happen.
- For `UPDATE` statements, if a row is already locked, InnoDB performs a “semi-consistent” read, returning 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.

Consider the following example, beginning with this table:

```
CREATE TABLE t (a INT NOT NULL, b INT) ENGINE = InnoDB;
INSERT INTO t VALUES (1,2),(2,3),(3,2),(4,3),(5,2);
COMMIT;
```

In this case, the table has no indexes, so searches and index scans use the hidden clustered index for record locking (see [Section 15.6.2.1, “Clustered and Secondary Indexes”](#)) rather than indexed columns.

Suppose that one session performs an `UPDATE` using these statements:

```
# Session A
START TRANSACTION;
UPDATE t SET b = 5 WHERE b = 3;
```

Suppose also that a second session performs an `UPDATE` by executing these statements following those of the first session:

```
# Session B
UPDATE t SET b = 4 WHERE b = 2;
```

As InnoDB executes each `UPDATE`, it first acquires an exclusive lock for each row, and then determines whether to modify it. If InnoDB does not modify the row, it releases the lock. Otherwise,

InnoDB retains the lock until the end of the transaction. This affects transaction processing as follows.

When using the default REPEATABLE READ isolation level, the first UPDATE acquires an x-lock on each row that it reads and does not release any of them:

```
x-lock(1,2); retain x-lock
x-lock(2,3); update(2,3) to (2,5); retain x-lock
x-lock(3,2); retain x-lock
x-lock(4,3); update(4,3) to (4,5); retain x-lock
x-lock(5,2); retain x-lock
```

The second UPDATE blocks as soon as it tries to acquire any locks (because first update has retained locks on all rows), and does not proceed until the first UPDATE commits or rolls back:

```
x-lock(1,2); block and wait for first UPDATE to commit or roll back
```

If READ COMMITTED is used instead, the first UPDATE acquires an x-lock on each row that it reads and releases those for rows that it does not modify:

```
x-lock(1,2); unlock(1,2)
x-lock(2,3); update(2,3) to (2,5); retain x-lock
x-lock(3,2); unlock(3,2)
x-lock(4,3); update(4,3) to (4,5); retain x-lock
x-lock(5,2); unlock(5,2)
```

For the second UPDATE, InnoDB does a “semi-consistent” read, returning the latest committed version of each row that it reads to MySQL so that MySQL can determine whether the row matches the WHERE condition of the UPDATE:

```
x-lock(1,2); update(1,2) to (1,4); retain x-lock
x-lock(2,3); unlock(2,3)
x-lock(3,2); update(3,2) to (3,4); retain x-lock
x-lock(4,3); unlock(4,3)
x-lock(5,2); update(5,2) to (5,4); retain x-lock
```

However, if the WHERE condition includes an indexed column, and InnoDB uses the index, only the indexed column is considered when taking and retaining record locks. In the following example, the first UPDATE takes and retains an x-lock on each row where b = 2. The second UPDATE blocks when it tries to acquire x-locks on the same records, as it also uses the index defined on column b.

```
CREATE TABLE t (a INT NOT NULL, b INT, c INT, INDEX (b)) ENGINE = InnoDB;
INSERT INTO t VALUES (1,2,3),(2,2,4);
COMMIT;

# Session A
START TRANSACTION;
UPDATE t SET b = 3 WHERE b = 2 AND c = 3;

# Session B
UPDATE t SET b = 4 WHERE b = 2 AND c = 4;
```

The READ COMMITTED isolation level can be set at startup or changed at runtime. At runtime, it can be set globally for all sessions, or individually per session.

- [READ UNCOMMITTED](#)

SELECT statements are performed in a nonlocking fashion, but a possible earlier version of a row might be used. Thus, using this isolation level, such reads are not consistent. This is also called a **dirty read**. Otherwise, this isolation level works like READ COMMITTED.

- [SERIALIZABLE](#)

This level is like REPEATABLE READ, but InnoDB implicitly converts all plain SELECT statements to SELECT ... FOR SHARE if autocommit is disabled. If autocommit is enabled, the SELECT is its own transaction. It therefore is known to be read only and can be serialized if performed as a

consistent (nonlocking) read and need not block for other transactions. (To force a plain `SELECT` to block if other transactions have modified the selected rows, disable `autocommit`.)

**Note**

As of MySQL 8.0.22, DML operations that read data from MySQL grant tables (through a join list or subquery) but do not modify them do not acquire read locks on the MySQL grant tables, regardless of the isolation level. For more information, see [Grant Table Concurrency](#).

15.7.2.2 autocommit, Commit, and Rollback

In `InnoDB`, all user activity occurs inside a transaction. If `autocommit` mode is enabled, each SQL statement forms a single transaction on its own. By default, MySQL starts the session for each new connection with `autocommit` enabled, so MySQL does a commit after each SQL statement if that statement did not return an error. If a statement returns an error, the commit or rollback behavior depends on the error. See [Section 15.21.5, “InnoDB Error Handling”](#).

A session that has `autocommit` enabled can perform a multiple-statement transaction by starting it with an explicit `START TRANSACTION` or `BEGIN` statement and ending it with a `COMMIT` or `ROLLBACK` statement. See [Section 13.3.1, “START TRANSACTION, COMMIT, and ROLLBACK Statements”](#).

If `autocommit` mode is disabled within a session with `SET autocommit = 0`, the session always has a transaction open. A `COMMIT` or `ROLLBACK` statement ends the current transaction and a new one starts.

If a session that has `autocommit` disabled ends without explicitly committing the final transaction, MySQL rolls back that transaction.

Some statements implicitly end a transaction, as if you had done a `COMMIT` before executing the statement. For details, see [Section 13.3.3, “Statements That Cause an Implicit Commit”](#).

A `COMMIT` means that the changes made in the current transaction are made permanent and become visible to other sessions. A `ROLLBACK` statement, on the other hand, cancels all modifications made by the current transaction. Both `COMMIT` and `ROLLBACK` release all `InnoDB` locks that were set during the current transaction.

Grouping DML Operations with Transactions

By default, connection to the MySQL server begins with `autocommit` mode enabled, which automatically commits every SQL statement as you execute it. This mode of operation might be unfamiliar if you have experience with other database systems, where it is standard practice to issue a sequence of `DML` statements and commit them or roll them back all together.

To use multiple-statement `transactions`, switch autocommit off with the SQL statement `SET autocommit = 0` and end each transaction with `COMMIT` or `ROLLBACK` as appropriate. To leave autocommit on, begin each transaction with `START TRANSACTION` and end it with `COMMIT` or `ROLLBACK`. The following example shows two transactions. The first is committed; the second is rolled back.

```
$> mysql test
mysql> CREATE TABLE customer (a INT, b CHAR (20), INDEX (a));
Query OK, 0 rows affected (0.00 sec)
mysql> -- Do a transaction with autocommit turned on.
mysql> START TRANSACTION;
Query OK, 0 rows affected (0.00 sec)
mysql> INSERT INTO customer VALUES (10, 'Heikki');
Query OK, 1 row affected (0.00 sec)
mysql> COMMIT;
Query OK, 0 rows affected (0.00 sec)
mysql> -- Do another transaction with autocommit turned off.
```

```
mysql> SET autocommit=0;
Query OK, 0 rows affected (0.00 sec)
mysql> INSERT INTO customer VALUES (15, 'John');
Query OK, 1 row affected (0.00 sec)
mysql> INSERT INTO customer VALUES (20, 'Paul');
Query OK, 1 row affected (0.00 sec)
mysql> DELETE FROM customer WHERE b = 'Heikki';
Query OK, 1 row affected (0.00 sec)
mysql> -- Now we undo those last 2 inserts and the delete.
mysql> ROLLBACK;
Query OK, 0 rows affected (0.00 sec)
mysql> SELECT * FROM customer;
+----+----+
| a | b |
+----+----+
| 10 | Heikki |
+----+----+
1 row in set (0.00 sec)
mysql>
```

Transactions in Client-Side Languages

In APIs such as PHP, Perl DBI, JDBC, ODBC, or the standard C call interface of MySQL, you can send transaction control statements such as `COMMIT` to the MySQL server as strings just like any other SQL statements such as `SELECT` or `INSERT`. Some APIs also offer separate special transaction commit and rollback functions or methods.

15.7.2.3 Consistent Nonlocking Reads

A **consistent read** means that InnoDB uses multi-versioning to present to a query a snapshot of the database at a point in time. The query sees the changes made by transactions that committed before that point in time, and no changes made by later or uncommitted transactions. The exception to this rule is that the query sees the changes made by earlier statements within the same transaction. This exception causes the following anomaly: If you update some rows in a table, a `SELECT` sees the latest version of the updated rows, but it might also see older versions of any rows. If other sessions simultaneously update the same table, the anomaly means that you might see the table in a state that never existed in the database.

If the transaction **isolation level** is `REPEATABLE READ` (the default level), all consistent reads within the same transaction read the snapshot established by the first such read in that transaction. You can get a fresher snapshot for your queries by committing the current transaction and after that issuing new queries.

With `READ COMMITTED` isolation level, each consistent read within a transaction sets and reads its own fresh snapshot.

Consistent read is the default mode in which InnoDB processes `SELECT` statements in `READ COMMITTED` and `REPEATABLE READ` isolation levels. A consistent read does not set any locks on the tables it accesses, and therefore other sessions are free to modify those tables at the same time a consistent read is being performed on the table.

Suppose that you are running in the default `REPEATABLE READ` isolation level. When you issue a consistent read (that is, an ordinary `SELECT` statement), InnoDB gives your transaction a timepoint according to which your query sees the database. If another transaction deletes a row and commits after your timepoint was assigned, you do not see the row as having been deleted. Inserts and updates are treated similarly.



Note

The snapshot of the database state applies to `SELECT` statements within a transaction, not necessarily to `DML` statements. If you insert or modify some rows and then commit that transaction, a `DELETE` or `UPDATE` statement issued from another concurrent `REPEATABLE READ` transaction could affect those just-

committed rows, even though the session could not query them. If a transaction does update or delete rows committed by a different transaction, those changes do become visible to the current transaction. For example, you might encounter a situation like the following:

```
SELECT COUNT(c1) FROM t1 WHERE c1 = 'xyz';
-- Returns 0: no rows match.
DELETE FROM t1 WHERE c1 = 'xyz';
-- Deletes several rows recently committed by other transaction.

SELECT COUNT(c2) FROM t1 WHERE c2 = 'abc';
-- Returns 0: no rows match.
UPDATE t1 SET c2 = 'cba' WHERE c2 = 'abc';
-- Affects 10 rows: another txn just committed 10 rows with 'abc' values.
SELECT COUNT(c2) FROM t1 WHERE c2 = 'cba';
-- Returns 10: this txn can now see the rows it just updated.
```

You can advance your timepoint by committing your transaction and then doing another `SELECT` or `START TRANSACTION WITH CONSISTENT SNAPSHOT`.

This is called *multi-versioned concurrency control*.

In the following example, session A sees the row inserted by B only when B has committed the insert and A has committed as well, so that the timepoint is advanced past the commit of B.

	Session A	Session B
time	SET autocommit=0;	SET autocommit=0;
	SELECT * FROM t;	
	empty set	INSERT INTO t VALUES (1, 2);
v	SELECT * FROM t;	
	empty set	COMMIT;
	SELECT * FROM t;	
	empty set	
	COMMIT;	
	SELECT * FROM t;	

	1 2	

If you want to see the “freshest” state of the database, use either the `READ COMMITTED` isolation level or a `locking read`:

```
SELECT * FROM t FOR SHARE;
```

With `READ COMMITTED` isolation level, each consistent read within a transaction sets and reads its own fresh snapshot. With `FOR SHARE`, a locking read occurs instead: A `SELECT` blocks until the transaction containing the freshest rows ends (see [Section 15.7.2.4, “Locking Reads”](#)).

Consistent read does not work over certain DDL statements:

- Consistent read does not work over `DROP TABLE`, because MySQL cannot use a table that has been dropped and InnoDB destroys the table.
- Consistent read does not work over `ALTER TABLE` operations that make a temporary copy of the original table and delete the original table when the temporary copy is built. When you reissue a consistent read within a transaction, rows in the new table are not visible because those rows did not exist when the transaction’s snapshot was taken. In this case, the transaction returns an error: `ER_TABLE_DEF_CHANGED`, “Table definition has changed, please retry transaction”.

The type of read varies for selects in clauses like `INSERT INTO ... SELECT, UPDATE ... (SELECT)`, and `CREATE TABLE ... SELECT` that do not specify `FOR UPDATE` or `FOR SHARE`:

- By default, InnoDB uses stronger locks for those statements and the `SELECT` part acts like `READ COMMITTED`, where each consistent read, even within the same transaction, sets and reads its own fresh snapshot.
- To perform a nonlocking read in such cases, set the isolation level of the transaction to `READ UNCOMMITTED` or `READ COMMITTED` to avoid setting locks on rows read from the selected table.

15.7.2.4 Locking Reads

If you query data and then insert or update related data within the same transaction, the regular `SELECT` statement does not give enough protection. Other transactions can update or delete the same rows you just queried. InnoDB supports two types of **locking reads** that offer extra safety:

- `SELECT ... FOR SHARE`

Sets a shared mode lock on any rows that are read. Other sessions can read the rows, but cannot modify them until your transaction commits. If any of these rows were changed by another transaction that has not yet committed, your query waits until that transaction ends and then uses the latest values.



Note

`SELECT ... FOR SHARE` is a replacement for `SELECT ... LOCK IN SHARE MODE`, but `LOCK IN SHARE MODE` remains available for backward compatibility. The statements are equivalent. However, `FOR SHARE` supports `OF table_name`, `NOWAIT`, and `SKIP LOCKED` options. See [Locking Read Concurrency with NOWAIT and SKIP LOCKED](#).

Prior to MySQL 8.0.22, `SELECT ... FOR SHARE` requires the `SELECT` privilege and at least one of the `DELETE`, `LOCK TABLES`, or `UPDATE` privileges. From MySQL 8.0.22, only the `SELECT` privilege is required.

From MySQL 8.0.22, `SELECT ... FOR SHARE` statements do not acquire read locks on MySQL grant tables. For more information, see [Grant Table Concurrency](#).

- `SELECT ... FOR UPDATE`

For index records the search encounters, locks the rows and any associated index entries, the same as if you issued an `UPDATE` statement for those rows. Other transactions are blocked from updating those rows, from doing `SELECT ... FOR SHARE`, or from reading the data in certain transaction isolation levels. Consistent reads ignore any locks set on the records that exist in the read view. (Old versions of a record cannot be locked; they are reconstructed by applying `undo logs` on an in-memory copy of the record.)

`SELECT ... FOR UPDATE` requires the `SELECT` privilege and at least one of the `DELETE`, `LOCK TABLES`, or `UPDATE` privileges.

These clauses are primarily useful when dealing with tree-structured or graph-structured data, either in a single table or split across multiple tables. You traverse edges or tree branches from one place to another, while reserving the right to come back and change any of these “pointer” values.

All locks set by `FOR SHARE` and `FOR UPDATE` queries are released when the transaction is committed or rolled back.



Note

Locking reads are only possible when autocommit is disabled (either by beginning transaction with `START TRANSACTION` or by setting `autocommit` to 0).

A locking read clause in an outer statement does not lock the rows of a table in a nested subquery unless a locking read clause is also specified in the subquery. For example, the following statement does not lock rows in table `t2`.

```
SELECT * FROM t1 WHERE c1 = (SELECT c1 FROM t2) FOR UPDATE;
```

To lock rows in table `t2`, add a locking read clause to the subquery:

```
SELECT * FROM t1 WHERE c1 = (SELECT c1 FROM t2 FOR UPDATE) FOR UPDATE;
```

Locking Read Examples

Suppose that you want to insert a new row into a table `child`, and make sure that the child row has a parent row in table `parent`. Your application code can ensure referential integrity throughout this sequence of operations.

First, use a consistent read to query the table `PARENT` and verify that the parent row exists. Can you safely insert the child row to table `CHILD`? No, because some other session could delete the parent row in the moment between your `SELECT` and your `INSERT`, without you being aware of it.

To avoid this potential issue, perform the `SELECT` using `FOR SHARE`:

```
SELECT * FROM parent WHERE NAME = 'Jones' FOR SHARE;
```

After the `FOR SHARE` query returns the parent '`Jones`', you can safely add the child record to the `CHILD` table and commit the transaction. Any transaction that tries to acquire an exclusive lock in the applicable row in the `PARENT` table waits until you are finished, that is, until the data in all tables is in a consistent state.

For another example, consider an integer counter field in a table `CHILD_CODES`, used to assign a unique identifier to each child added to table `CHILD`. Do not use either consistent read or a shared mode read to read the present value of the counter, because two users of the database could see the same value for the counter, and a duplicate-key error occurs if two transactions attempt to add rows with the same identifier to the `CHILD` table.

Here, `FOR SHARE` is not a good solution because if two users read the counter at the same time, at least one of them ends up in deadlock when it attempts to update the counter.

To implement reading and incrementing the counter, first perform a locking read of the counter using `FOR UPDATE`, and then increment the counter. For example:

```
SELECT counter_field FROM child_codes FOR UPDATE;
UPDATE child_codes SET counter_field = counter_field + 1;
```

A `SELECT ... FOR UPDATE` reads the latest available data, setting exclusive locks on each row it reads. Thus, it sets the same locks a searched SQL `UPDATE` would set on the rows.

The preceding description is merely an example of how `SELECT ... FOR UPDATE` works. In MySQL, the specific task of generating a unique identifier actually can be accomplished using only a single access to the table:

```
UPDATE child_codes SET counter_field = LAST_INSERT_ID(counter_field + 1);
SELECT LAST_INSERT_ID();
```

The `SELECT` statement merely retrieves the identifier information (specific to the current connection). It does not access any table.

Locking Read Concurrency with NOWAIT and SKIP LOCKED

If a row is locked by a transaction, a `SELECT ... FOR UPDATE` or `SELECT ... FOR SHARE` transaction that requests the same locked row must wait until the blocking transaction releases the row

lock. This behavior prevents transactions from updating or deleting rows that are queried for updates by other transactions. However, waiting for a row lock to be released is not necessary if you want the query to return immediately when a requested row is locked, or if excluding locked rows from the result set is acceptable.

To avoid waiting for other transactions to release row locks, `NOWAIT` and `SKIP LOCKED` options may be used with `SELECT ... FOR UPDATE` or `SELECT ... FOR SHARE` locking read statements.

- `NOWAIT`

A locking read that uses `NOWAIT` never waits to acquire a row lock. The query executes immediately, failing with an error if a requested row is locked.

- `SKIP LOCKED`

A locking read that uses `SKIP LOCKED` never waits to acquire a row lock. The query executes immediately, removing locked rows from the result set.



Note

Queries that skip locked rows return an inconsistent view of the data. `SKIP LOCKED` is therefore not suitable for general transactional work. However, it may be used to avoid lock contention when multiple sessions access the same queue-like table.

`NOWAIT` and `SKIP LOCKED` only apply to row-level locks.

Statements that use `NOWAIT` or `SKIP LOCKED` are unsafe for statement based replication.

The following example demonstrates `NOWAIT` and `SKIP LOCKED`. Session 1 starts a transaction that takes a row lock on a single record. Session 2 attempts a locking read on the same record using the `NOWAIT` option. Because the requested row is locked by Session 1, the locking read returns immediately with an error. In Session 3, the locking read with `SKIP LOCKED` returns the requested rows except for the row that is locked by Session 1.

```
# Session 1:

mysql> CREATE TABLE t (i INT, PRIMARY KEY (i)) ENGINE = InnoDB;
mysql> INSERT INTO t (i) VALUES(1),(2),(3);
mysql> START TRANSACTION;

mysql> SELECT * FROM t WHERE i = 2 FOR UPDATE;
+---+
| i |
+---+
| 2 |
+---+

# Session 2:

mysql> START TRANSACTION;

mysql> SELECT * FROM t WHERE i = 2 FOR UPDATE NOWAIT;
ERROR 3572 (HY000): Do not wait for lock.

# Session 3:

mysql> START TRANSACTION;

mysql> SELECT * FROM t FOR UPDATE SKIP LOCKED;
+---+
| i |
+---+
| 1 |
+---+
```

	3	
+	---	+

15.7.3 Locks Set by Different SQL Statements in InnoDB

A [locking read](#), an [UPDATE](#), or a [DELETE](#) generally set record locks on every index record that is scanned in the processing of an SQL statement. It does not matter whether there are [WHERE](#) conditions in the statement that would exclude the row. InnoDB does not remember the exact [WHERE](#) condition, but only knows which index ranges were scanned. The locks are normally [next-key locks](#) that also block inserts into the “gap” immediately before the record. However, [gap locking](#) can be disabled explicitly, which causes next-key locking not to be used. For more information, see [Section 15.7.1, “InnoDB Locking”](#). The transaction isolation level can also affect which locks are set; see [Section 15.7.2.1, “Transaction Isolation Levels”](#).

If a secondary index is used in a search and the index record locks to be set are exclusive, InnoDB also retrieves the corresponding clustered index records and sets locks on them.

If you have no indexes suitable for your statement and MySQL must scan the entire table to process the statement, every row of the table becomes locked, which in turn blocks all inserts by other users to the table. It is important to create good indexes so that your queries do not scan more rows than necessary.

InnoDB sets specific types of locks as follows.

- `SELECT ... FROM` is a consistent read, reading a snapshot of the database and setting no locks unless the transaction isolation level is set to [SERIALIZABLE](#). For [SERIALIZABLE](#) level, the search sets shared next-key locks on the index records it encounters. However, only an index record lock is required for statements that lock rows using a unique index to search for a unique row.
- `SELECT ... FOR UPDATE` and `SELECT ... FOR SHARE` statements that use a unique index acquire locks for scanned rows, and release the locks for rows that do not qualify for inclusion in the result set (for example, if they do not meet the criteria given in the [WHERE](#) clause). However, in some cases, rows might not be unlocked immediately because the relationship between a result row and its original source is lost during query execution. For example, in a [UNION](#), scanned (and locked) rows from a table might be inserted into a temporary table before evaluating whether they qualify for the result set. In this circumstance, the relationship of the rows in the temporary table to the rows in the original table is lost and the latter rows are not unlocked until the end of query execution.
- For [locking reads](#) (`SELECT` with `FOR UPDATE` or `FOR SHARE`), `UPDATE`, and `DELETE` statements, the locks that are taken depend on whether the statement uses a unique index with a unique search condition or a range-type search condition.
 - For a unique index with a unique search condition, InnoDB locks only the index record found, not the [gap](#) before it.
 - For other search conditions, and for non-unique indexes, InnoDB locks the index range scanned, using [gap locks](#) or [next-key locks](#) to block insertions by other sessions into the gaps covered by the range. For information about gap locks and next-key locks, see [Section 15.7.1, “InnoDB Locking”](#).
- For index records the search encounters, `SELECT ... FOR UPDATE` blocks other sessions from doing `SELECT ... FOR SHARE` or from reading in certain transaction isolation levels. Consistent reads ignore any locks set on the records that exist in the read view.
- `UPDATE ... WHERE ...` sets an exclusive next-key lock on every record the search encounters. However, only an index record lock is required for statements that lock rows using a unique index to search for a unique row.
- When `UPDATE` modifies a clustered index record, implicit locks are taken on affected secondary index records. The `UPDATE` operation also takes shared locks on affected secondary index records

when performing duplicate check scans prior to inserting new secondary index records, and when inserting new secondary index records.

- `DELETE FROM ... WHERE ...` sets an exclusive next-key lock on every record the search encounters. However, only an index record lock is required for statements that lock rows using a unique index to search for a unique row.
- `INSERT` sets an exclusive lock on the inserted row. This lock is an index-record lock, not a next-key lock (that is, there is no gap lock) and does not prevent other sessions from inserting into the gap before the inserted row.

Prior to inserting the row, a type of gap lock called an insert intention gap lock is set. This lock signals the intent to insert in such a way that multiple transactions inserting into the same index gap need not wait for each other if they are not inserting at the same position within the gap. Suppose that there are index records with values of 4 and 7. Separate transactions that attempt to insert values of 5 and 6 each lock the gap between 4 and 7 with insert intention locks prior to obtaining the exclusive lock on the inserted row, but do not block each other because the rows are nonconflicting.

If a duplicate-key error occurs, a shared lock on the duplicate index record is set. This use of a shared lock can result in deadlock should there be multiple sessions trying to insert the same row if another session already has an exclusive lock. This can occur if another session deletes the row. Suppose that an `InnoDB` table `t1` has the following structure:

```
CREATE TABLE t1 (i INT, PRIMARY KEY (i)) ENGINE = InnoDB;
```

Now suppose that three sessions perform the following operations in order:

Session 1:

```
START TRANSACTION;
INSERT INTO t1 VALUES(1);
```

Session 2:

```
START TRANSACTION;
INSERT INTO t1 VALUES(1);
```

Session 3:

```
START TRANSACTION;
INSERT INTO t1 VALUES(1);
```

Session 1:

```
ROLLBACK;
```

The first operation by session 1 acquires an exclusive lock for the row. The operations by sessions 2 and 3 both result in a duplicate-key error and they both request a shared lock for the row. When session 1 rolls back, it releases its exclusive lock on the row and the queued shared lock requests for sessions 2 and 3 are granted. At this point, sessions 2 and 3 deadlock: Neither can acquire an exclusive lock for the row because of the shared lock held by the other.

A similar situation occurs if the table already contains a row with key value 1 and three sessions perform the following operations in order:

Session 1:

```
START TRANSACTION;
DELETE FROM t1 WHERE i = 1;
```

Session 2:

```
START TRANSACTION;
INSERT INTO t1 VALUES(1);
```

Session 3:

```
START TRANSACTION;
INSERT INTO t1 VALUES(1);
```

Session 1:

```
COMMIT;
```

The first operation by session 1 acquires an exclusive lock for the row. The operations by sessions 2 and 3 both result in a duplicate-key error and they both request a shared lock for the row. When session 1 commits, it releases its exclusive lock on the row and the queued shared lock requests for sessions 2 and 3 are granted. At this point, sessions 2 and 3 deadlock: Neither can acquire an exclusive lock for the row because of the shared lock held by the other.

- `INSERT ... ON DUPLICATE KEY UPDATE` differs from a simple `INSERT` in that an exclusive lock rather than a shared lock is placed on the row to be updated when a duplicate-key error occurs. An exclusive index-record lock is taken for a duplicate primary key value. An exclusive next-key lock is taken for a duplicate unique key value.
- `REPLACE` is done like an `INSERT` if there is no collision on a unique key. Otherwise, an exclusive next-key lock is placed on the row to be replaced.
- `INSERT INTO T SELECT ... FROM S WHERE ...` sets an exclusive index record lock (without a gap lock) on each row inserted into `T`. If the transaction isolation level is `READ COMMITTED`, InnoDB does the search on `S` as a consistent read (no locks). Otherwise, InnoDB sets shared next-key locks on rows from `S`. InnoDB has to set locks in the latter case: During roll-forward recovery using a statement-based binary log, every SQL statement must be executed in exactly the same way it was done originally.

`CREATE TABLE ... SELECT ...` performs the `SELECT` with shared next-key locks or as a consistent read, as for `INSERT ... SELECT`.

When a `SELECT` is used in the constructs `REPLACE INTO t SELECT ... FROM s WHERE ...` or `UPDATE t ... WHERE col IN (SELECT ... FROM s ...)`, InnoDB sets shared next-key locks on rows from table `s`.

- InnoDB sets an exclusive lock on the end of the index associated with the `AUTO_INCREMENT` column while initializing a previously specified `AUTO_INCREMENT` column on a table.

With `innodb_autoinc_lock_mode=0`, InnoDB uses a special `AUTO-INC` table lock mode where the lock is obtained and held to the end of the current SQL statement (not to the end of the entire transaction) while accessing the auto-increment counter. Other clients cannot insert into the table while the `AUTO-INC` table lock is held. The same behavior occurs for “bulk inserts” with `innodb_autoinc_lock_mode=1`. Table-level `AUTO-INC` locks are not used with `innodb_autoinc_lock_mode=2`. For more information, See [Section 15.6.1.6, “AUTO_INCREMENT Handling in InnoDB”](#).

InnoDB fetches the value of a previously initialized `AUTO_INCREMENT` column without setting any locks.

- If a `FOREIGN KEY` constraint is defined on a table, any insert, update, or delete that requires the constraint condition to be checked sets shared record-level locks on the records that it looks at to check the constraint. InnoDB also sets these locks in the case where the constraint fails.
- `LOCK TABLES` sets table locks, but it is the higher MySQL layer above the InnoDB layer that sets these locks. InnoDB is aware of table locks if `innodb_table_locks = 1` (the default) and `autocommit = 0`, and the MySQL layer above InnoDB knows about row-level locks.

Otherwise, InnoDB's automatic deadlock detection cannot detect deadlocks where such table locks are involved. Also, because in this case the higher MySQL layer does not know about row-level

locks, it is possible to get a table lock on a table where another session currently has row-level locks. However, this does not endanger transaction integrity, as discussed in [Section 15.7.5.2, “Deadlock Detection”](#).

- `LOCK TABLES` acquires two locks on each table if `innodb_table_locks=1` (the default). In addition to a table lock on the MySQL layer, it also acquires an `InnoDB` table lock. To avoid acquiring `InnoDB` table locks, set `innodb_table_locks=0`. If no `InnoDB` table lock is acquired, `LOCK TABLES` completes even if some records of the tables are being locked by other transactions.

In MySQL 8.0, `innodb_table_locks=0` has no effect for tables locked explicitly with `LOCK TABLES ... WRITE`. It does have an effect for tables locked for read or write by `LOCK TABLES ... WRITE` implicitly (for example, through triggers) or by `LOCK TABLES ... READ`.

- All `InnoDB` locks held by a transaction are released when the transaction is committed or aborted. Thus, it does not make much sense to invoke `LOCK TABLES` on `InnoDB` tables in `autocommit=1` mode because the acquired `InnoDB` table locks would be released immediately.
- You cannot lock additional tables in the middle of a transaction because `LOCK TABLES` performs an implicit `COMMIT` and `UNLOCK TABLES`.

15.7.4 Phantom Rows

The so-called *phantom* problem occurs within a transaction when the same query produces different sets of rows at different times. For example, if a `SELECT` is executed twice, but returns a row the second time that was not returned the first time, the row is a “phantom” row.

Suppose that there is an index on the `id` column of the `child` table and that you want to read and lock all rows from the table having an identifier value larger than 100, with the intention of updating some column in the selected rows later:

```
SELECT * FROM child WHERE id > 100 FOR UPDATE;
```

The query scans the index starting from the first record where `id` is bigger than 100. Let the table contain rows having `id` values of 90 and 102. If the locks set on the index records in the scanned range do not lock out inserts made in the gaps (in this case, the gap between 90 and 102), another session can insert a new row into the table with an `id` of 101. If you were to execute the same `SELECT` within the same transaction, you would see a new row with an `id` of 101 (a “phantom”) in the result set returned by the query. If we regard a set of rows as a data item, the new phantom child would violate the isolation principle of transactions that a transaction should be able to run so that the data it has read does not change during the transaction.

To prevent phantoms, `InnoDB` uses an algorithm called *next-key locking* that combines index-row locking with gap locking. `InnoDB` performs row-level locking in such a way that when it searches or scans a table index, it sets shared or exclusive locks on the index records it encounters. Thus, the row-level locks are actually index-record locks. In addition, a next-key lock on an index record also affects the “gap” before the index record. That is, a next-key lock is an index-record lock plus a gap lock on the gap preceding the index record. If one session has a shared or exclusive lock on record `R` in an index, another session cannot insert a new index record in the gap immediately before `R` in the index order.

When `InnoDB` scans an index, it can also lock the gap after the last record in the index. Just that happens in the preceding example: To prevent any insert into the table where `id` would be bigger than 100, the locks set by `InnoDB` include a lock on the gap following `id` value 102.

You can use next-key locking to implement a uniqueness check in your application: If you read your data in share mode and do not see a duplicate for a row you are going to insert, then you can safely insert your row and know that the next-key lock set on the successor of your row during the read prevents anyone meanwhile inserting a duplicate for your row. Thus, the next-key locking enables you to “lock” the nonexistence of something in your table.

Gap locking can be disabled as discussed in [Section 15.7.1, “InnoDB Locking”](#). This may cause phantom problems because other sessions can insert new rows into the gaps when gap locking is disabled.

15.7.5 Deadlocks in InnoDB

A deadlock is 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 ever release the locks it holds.

A deadlock can occur when 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. For a deadlock example, see [Section 15.7.5.1, “An InnoDB Deadlock Example”](#).

To reduce the possibility of deadlocks, use transactions rather than `LOCK TABLES` 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. For more information about avoiding and recovering from deadlock conditions, see [Section 15.7.5.3, “How to Minimize and Handle Deadlocks”](#).

When deadlock detection is enabled (the default) and a deadlock does occur, InnoDB detects the condition and rolls back one of the transactions (the victim). If deadlock detection is disabled using the `innodb_deadlock_detect` variable, InnoDB relies on the `innodb_lock_wait_timeout` setting to roll back transactions in case of a deadlock. Thus, even if your application logic is correct, you must still handle the case where a transaction must be retried. To view the last deadlock in an InnoDB user transaction, use `SHOW ENGINE INNODB STATUS`. If frequent deadlocks highlight a problem with transaction structure or application error handling, enable `innodb_print_all_deadlocks` to print information about all deadlocks to the `mysqld` error log. For more information about how deadlocks are automatically detected and handled, see [Section 15.7.5.2, “Deadlock Detection”](#).

15.7.5.1 An InnoDB Deadlock Example

The following example illustrates how an error can occur when a lock request causes a deadlock. The example involves two clients, A and B.

InnoDB status contains details of the last deadlock. For frequent deadlocks, enable global variable `innodb_print_all_deadlocks`. This adds deadlock information to the error log.

Client A enables `innodb_print_all_deadlocks`, creates two tables, 'Animals' and 'Birds', and inserts data into each. Client A begins a transaction, and selects a row in Animals in share mode:

```
mysql> SET GLOBAL innodb_print_all_deadlocks = ON;
Query OK, 0 rows affected (0.00 sec)

mysql> CREATE TABLE Animals (name VARCHAR(10) PRIMARY KEY, value INT) ENGINE = InnoDB;
Query OK, 0 rows affected (0.01 sec)

mysql> CREATE TABLE Birds (name VARCHAR(10) PRIMARY KEY, value INT) ENGINE = InnoDB;
Query OK, 0 rows affected (0.01 sec)

mysql> INSERT INTO Animals (name,value) VALUES ("Aardvark",10);
Query OK, 1 row affected (0.00 sec)

mysql> INSERT INTO Birds (name,value) VALUES ("Buzzard",20);
Query OK, 1 row affected (0.00 sec)

mysql> START TRANSACTION;
Query OK, 0 rows affected (0.00 sec)

mysql> SELECT value FROM Animals WHERE name='Aardvark' FOR SHARE;
+-----+
```

```
| value |
+-----+
|    10 |
+-----+
1 row in set (0.00 sec)
```

Next, client B begins a transaction, and selects a row in Birds in share mode:

```
mysql> START TRANSACTION;
Query OK, 0 rows affected (0.00 sec)

mysql> SELECT value FROM Birds WHERE name='Buzzard' FOR SHARE;
+-----+
| value |
+-----+
|    20 |
+-----+
1 row in set (0.00 sec)
```

The Performance Schema shows the locks after the two select statements:

```
mysql> SELECT ENGINE_TRANSACTION_ID as Trx_Id,
          OBJECT_NAME as `Table`,
          INDEX_NAME as `Index`,
          LOCK_DATA as Data,
          LOCK_MODE as Mode,
          LOCK_STATUS as Status,
          LOCK_TYPE as Type
      FROM performance_schema.data_locks;
+-----+-----+-----+-----+-----+-----+-----+-----+
| Trx_Id | Table | Index | Data | Mode | Status | Type |
+-----+-----+-----+-----+-----+-----+-----+-----+
| 421291106147544 | Animals | NULL | NULL | IS | GRANTED | TABLE |
| 421291106147544 | Animals | PRIMARY | 'Aardvark' | S,REC_NOT_GAP | GRANTED | RECORD |
| 421291106148352 | Birds | NULL | NULL | IS | GRANTED | TABLE |
| 421291106148352 | Birds | PRIMARY | 'Buzzard' | S,REC_NOT_GAP | GRANTED | RECORD |
+-----+-----+-----+-----+-----+-----+-----+-----+
4 rows in set (0.00 sec)
```

Client B then updates a row in Animals:

```
mysql> UPDATE Animals SET value=30 WHERE name='Aardvark';
```

Client B has to wait. The Performance Schema shows the wait for a lock:

```
mysql> SELECT REQUESTING_ENGINE_LOCK_ID as Req_Lock_Id,
          REQUESTING_ENGINE_TRANSACTION_ID as Req_Trx_Id,
          BLOCKING_ENGINE_LOCK_ID as Blk_Lock_Id,
          BLOCKING_ENGINE_TRANSACTION_ID as Blk_Trx_Id
      FROM performance_schema.data_lock_waits;
+-----+-----+-----+-----+
| Req_Lock_Id | Req_Trx_Id | Blk_Lock_Id | Blk_Trx_Id |
+-----+-----+-----+-----+
| 139816129437696:27:4:2:139816016601240 | 43260 | 139816129436888:27:4:2:139816016594720 | 4212911061
+-----+-----+-----+-----+
1 row in set (0.00 sec)

mysql> SELECT ENGINE_LOCK_ID as Lock_Id,
          ENGINE_TRANSACTION_ID as Trx_id,
          OBJECT_NAME as `Table`,
          INDEX_NAME as `Index`,
          LOCK_DATA as Data,
          LOCK_MODE as Mode,
          LOCK_STATUS as Status,
          LOCK_TYPE as Type
      FROM performance_schema.data_locks;
+-----+-----+-----+-----+-----+-----+
| Lock_Id | Trx_Id | Table | Index | Data | Mode |
+-----+-----+-----+-----+-----+-----+
| 139816129437696:1187:139816016603896 | 43260 | Animals | NULL | NULL | IX |
| 139816129437696:1188:139816016603808 | 43260 | Birds | NULL | NULL | IS |
+-----+-----+-----+-----+-----+-----+
```

```

| 139816129437696:28:4:2:139816016600896 |           43260 | Birds   | PRIMARY | 'Buzzard' | S,REC_NOT
| 139816129437696:27:4:2:139816016601240 |           43260 | Animals | PRIMARY | 'Aardvark' | X,REC_NOT
| 139816129436888:1187:139816016597712 | 421291106147544 | Animals | NULL    | NULL     | IS
| 139816129436888:27:4:2:139816016594720 | 421291106147544 | Animals | PRIMARY | 'Aardvark' | S,REC_NOT
+-----+-----+-----+-----+-----+
6 rows in set (0.00 sec)

```

InnoDB only uses sequential transaction ids when a transaction attempts to modify the database. Thererfore, the previous read-only transaction id changes from 421291106148352 to 43260.

If client A attempts to update a row in Birds at the same time, this will lead to a deadlock:

```

mysql> UPDATE Birds SET value=40 WHERE name='Buzzard';
ERROR 1213 (40001): Deadlock found when trying to get lock; try restarting transaction

```

InnoDB rolls back the transaction that caused the deadlock. The first update, from Client B, can now proceed.

The Information Schema contains the number of deadlocks:

```

mysql> SELECT `count` FROM INFORMATION_SCHEMA.INNODB_METRICS
      WHERE NAME="lock_deadlocks";
+-----+
| count |
+-----+
|     1 |
+-----+
1 row in set (0.00 sec)

```

The InnoDB status contains the following information about the deadlock and transactions. It also shows that the read-only transaction id 421291106147544 changes to sequential transaction id 43261.

```

mysql> SHOW ENGINE INNODB STATUS;
-----
LATEST DETECTED DEADLOCK
-----
2022-11-25 15:58:22 139815661168384
*** (1) TRANSACTION:
TRANSACTION 43260, ACTIVE 186 sec starting index read
mysql tables in use 1, locked 1
LOCK WAIT 4 lock struct(s), heap size 1128, 2 row lock(s)
MySQL thread id 19, OS thread handle 139815619204864, query id 143 localhost u2 updating
UPDATE Animals SET value=30 WHERE name='Aardvark'

*** (1) HOLDS THE LOCK(S):
RECORD LOCKS space id 28 page no 4 n bits 72 index PRIMARY of table `test`.`Birds` trx id 43260 lock mode
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 7; hex 42757a7a617264; asc Buzzard;;
 1: len 6; hex 00000000a8fb; asc      ;;
 2: len 7; hex 82000000e40110; asc      ;;
 3: len 4; hex 80000014; asc      ;;

*** (1) WAITING FOR THIS LOCK TO BE GRANTED:
RECORD LOCKS space id 27 page no 4 n bits 72 index PRIMARY of table `test`.`Animals` trx id 43260 lock mode
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 8; hex 416172647661726b; asc Aardvark;;
 1: len 6; hex 00000000a8f9; asc      ;;
 2: len 7; hex 82000000e20110; asc      ;;
 3: len 4; hex 8000000a; asc      ;;

*** (2) TRANSACTION:
TRANSACTION 43261, ACTIVE 209 sec starting index read
mysql tables in use 1, locked 1
LOCK WAIT 4 lock struct(s), heap size 1128, 2 row lock(s)
MySQL thread id 18, OS thread handle 139815618148096, query id 146 localhost u1 updating
UPDATE Birds SET value=40 WHERE name='Buzzard'

*** (2) HOLDS THE LOCK(S):
RECORD LOCKS space id 27 page no 4 n bits 72 index PRIMARY of table `test`.`Animals` trx id 43261 lock mode
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 7; hex 42757a7a617264; asc Buzzard;;
 1: len 6; hex 00000000a8fb; asc      ;;
 2: len 7; hex 82000000e40110; asc      ;;
 3: len 4; hex 80000014; asc      ;;


```

```

Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 8; hex 416172647661726b; asc Aardvark;;
 1: len 6; hex 00000000a8f9; asc           ;;
 2: len 7; hex 82000000e20110; asc           ;;
 3: len 4; hex 8000000a; asc           ;;

*** (2) WAITING FOR THIS LOCK TO BE GRANTED:
RECORD LOCKS space id 28 page no 4 n bits 72 index PRIMARY of table `test`.`Birds` trx id 43261 lock_mode X
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 7; hex 42757a7a617264; asc Buzzard;;
 1: len 6; hex 00000000a8fb; asc           ;;
 2: len 7; hex 82000000e40110; asc           ;;
 3: len 4; hex 80000014; asc           ;;

*** WE ROLL BACK TRANSACTION (2)
-----
TRANSACTIONS
-----
Trx id counter 43262
Purge done for trx's n:o < 43256 undo n:o < 0 state: running but idle
History list length 0
LIST OF TRANSACTIONS FOR EACH SESSION:
---TRANSACTION 421291106147544, not started
0 lock struct(s), heap size 1128, 0 row lock(s)
---TRANSACTION 421291106146736, not started
0 lock struct(s), heap size 1128, 0 row lock(s)
---TRANSACTION 421291106145928, not started
0 lock struct(s), heap size 1128, 0 row lock(s)
---TRANSACTION 43260, ACTIVE 219 sec
4 lock struct(s), heap size 1128, 2 row lock(s), undo log entries 1
MySQL thread id 19, OS thread handle 139815619204864, query id 143 localhost u2

```

The error log contains this information about transactions and locks:

```

mysql> SELECT @@log_error;
+-----+
| @@log_error      |
+-----+
| /var/log/mysqld.log |
+-----+
1 row in set (0.00 sec)

TRANSACTION 43260, ACTIVE 186 sec starting index read
mysql tables in use 1, locked 1
LOCK WAIT 4 lock struct(s), heap size 1128, 2 row lock(s)
MySQL thread id 19, OS thread handle 139815619204864, query id 143 localhost u2 updating
UPDATE Animals SET value=30 WHERE name='Aardvark'
RECORD LOCKS space id 28 page no 4 n bits 72 index PRIMARY of table `test`.`Birds` trx id 43260 lock mode S
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 7; hex 42757a7a617264; asc Buzzard;;
 1: len 6; hex 00000000a8fb; asc           ;;
 2: len 7; hex 82000000e40110; asc           ;;
 3: len 4; hex 80000014; asc           ;;

RECORD LOCKS space id 27 page no 4 n bits 72 index PRIMARY of table `test`.`Animals` trx id 43260 lock mode S
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 8; hex 416172647661726b; asc Aardvark;;
 1: len 6; hex 00000000a8f9; asc           ;;
 2: len 7; hex 82000000e20110; asc           ;;
 3: len 4; hex 8000000a; asc           ;;

TRANSACTION 43261, ACTIVE 209 sec starting index read
mysql tables in use 1, locked 1
LOCK WAIT 4 lock struct(s), heap size 1128, 2 row lock(s)
MySQL thread id 18, OS thread handle 139815618148096, query id 146 localhost u1 updating
UPDATE Birds SET value=40 WHERE name='Buzzard'
RECORD LOCKS space id 27 page no 4 n bits 72 index PRIMARY of table `test`.`Animals` trx id 43261 lock mode S
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 8; hex 416172647661726b; asc Aardvark;;
 1: len 6; hex 00000000a8f9; asc           ;;
 2: len 7; hex 82000000e20110; asc           ;;
```

```

3: len 4; hex 8000000a; asc    ;;
RECORD LOCKS space id 28 page no 4 n bits 72 index PRIMARY of table `test`.`Birds` trx id 43261 lock_m
Record lock, heap no 2 PHYSICAL RECORD: n_fields 4; compact format; info bits 0
 0: len 7; hex 42757a7a617264; asc Buzzard;;
 1: len 6; hex 00000000a8fb; asc      ;;
 2: len 7; hex 82000000e40110; asc      ;;
 3: len 4; hex 80000014; asc    ;;

```

15.7.5.2 Deadlock Detection

When [deadlock detection](#) is enabled (the default), [InnoDB](#) automatically detects transaction [deadlocks](#) and rolls back a transaction or transactions to break the deadlock. [InnoDB](#) tries to pick small transactions to roll back, where the size of a transaction is determined by the number of rows inserted, updated, or deleted.

[InnoDB](#) is aware of table locks if `innodb_table_locks = 1` (the default) and `autocommit = 0`, and the MySQL layer above it knows about row-level locks. Otherwise, [InnoDB](#) cannot detect deadlocks where a table lock set by a MySQL `LOCK TABLES` statement or a lock set by a storage engine other than [InnoDB](#) is involved. Resolve these situations by setting the value of the `innodb_lock_wait_timeout` system variable.

If the `LATEST DETECTED DEADLOCK` section of [InnoDB](#) Monitor output includes a message stating `TOO DEEP OR LONG SEARCH IN THE LOCK TABLE WAITS-FOR GRAPH, WE WILL ROLL BACK FOLLOWING TRANSACTION`, this indicates that the number of transactions on the wait-for list has reached a limit of 200. A wait-for list that exceeds 200 transactions is treated as a deadlock and the transaction attempting to check the wait-for list is rolled back. The same error may also occur if the locking thread must look at more than 1,000,000 locks owned by transactions on the wait-for list.

For techniques to organize database operations to avoid deadlocks, see [Section 15.7.5, “Deadlocks in InnoDB”](#).

Disabling Deadlock Detection

On high concurrency systems, deadlock detection can cause a slowdown when numerous threads wait for the same lock. At times, it may be more efficient to disable deadlock detection and rely on the `innodb_lock_wait_timeout` setting for transaction rollback when a deadlock occurs. Deadlock detection can be disabled using the `innodb_deadlock_detect` variable.

15.7.5.3 How to Minimize and Handle Deadlocks

This section builds on the conceptual information about deadlocks in [Section 15.7.5.2, “Deadlock Detection”](#). It explains how to organize database operations to minimize deadlocks and the subsequent error handling required in applications.

[Deadlocks](#) are a classic problem in transactional databases, but they are not dangerous unless they are so frequent that you cannot run certain transactions at all. Normally, you must write your applications so that they are always prepared to re-issue a transaction if it gets rolled back because of a deadlock.

[InnoDB](#) uses automatic row-level locking. You can get deadlocks even in the case of transactions that just insert or delete a single row. That is because these operations are not really “atomic”; they automatically set locks on the (possibly several) index records of the row inserted or deleted.

You can cope with deadlocks and reduce the likelihood of their occurrence with the following techniques:

- At any time, issue `SHOW ENGINE INNODB STATUS` to determine the cause of the most recent deadlock. That can help you to tune your application to avoid deadlocks.
- If frequent deadlock warnings cause concern, collect more extensive debugging information by enabling the `innodb_print_all_deadlocks` variable. Information about each deadlock, not

just the latest one, is recorded in the MySQL [error log](#). Disable this option when you are finished debugging.

- Always be prepared to re-issue a transaction if it fails due to deadlock. Deadlocks are not dangerous. Just try again.
- Keep transactions small and short in duration to make them less prone to collision.
- Commit transactions immediately after making a set of related changes to make them less prone to collision. In particular, do not leave an interactive `mysql` session open for a long time with an uncommitted transaction.
- If you use [locking reads](#) (`SELECT ... FOR UPDATE` or `SELECT ... FOR SHARE`), try using a lower isolation level such as `READ COMMITTED`.
- When modifying multiple tables within a transaction, or different sets of rows in the same table, do those operations in a consistent order each time. Then transactions form well-defined queues and do not deadlock. For example, organize database operations into functions within your application, or call stored routines, rather than coding multiple similar sequences of `INSERT`, `UPDATE`, and `DELETE` statements in different places.
- Add well-chosen indexes to your tables so that your queries scan fewer index records and set fewer locks. Use `EXPLAIN SELECT` to determine which indexes the MySQL server regards as the most appropriate for your queries.
- Use less locking. If you can afford to permit a `SELECT` to return data from an old snapshot, do not add a `FOR UPDATE` or `FOR SHARE` clause to it. Using the `READ COMMITTED` isolation level is good here, because each consistent read within the same transaction reads from its own fresh snapshot.
- If nothing else helps, serialize your transactions with table-level locks. The correct way to use `LOCK TABLES` with transactional tables, such as `InnoDB` tables, is to begin a transaction with `SET autocommit = 0` (not `START TRANSACTION`) followed by `LOCK TABLES`, and to not call `UNLOCK TABLES` until you commit the transaction explicitly. For example, if you need to write to table `t1` and read from table `t2`, you can do this:

```
SET autocommit=0;
LOCK TABLES t1 WRITE, t2 READ, ...;
... do something with tables t1 and t2 here ...
COMMIT;
UNLOCK TABLES;
```

Table-level locks prevent concurrent updates to the table, avoiding deadlocks at the expense of less responsiveness for a busy system.

- Another way to serialize transactions is to create an auxiliary “semaphore” table that contains just a single row. Have each transaction update that row before accessing other tables. In that way, all transactions happen in a serial fashion. Note that the `InnoDB` instant deadlock detection algorithm also works in this case, because the serializing lock is a row-level lock. With MySQL table-level locks, the timeout method must be used to resolve deadlocks.

15.7.6 Transaction Scheduling

`InnoDB` uses the Contention-Aware Transaction Scheduling (CATS) algorithm to prioritize transactions that are waiting for locks. When multiple transactions are waiting for a lock on the same object, the CATS algorithm determines which transaction receives the lock first.

The CATS algorithm prioritizes waiting transactions by assigning a scheduling weight, which is computed based on the number of transactions that a transaction blocks. For example, if two transactions are waiting for a lock on the same object, the transaction that blocks the most transactions is assigned a greater scheduling weight. If weights are equal, priority is given to the longest waiting transaction.

**Note**

Prior to MySQL 8.0.20, InnoDB also uses a First In First Out (FIFO) algorithm to schedule transactions, and the CATS algorithm is used under heavy lock contention only. CATS algorithm enhancements in MySQL 8.0.20 rendered the FIFO algorithm redundant, permitting its removal. Transaction scheduling previously performed by the FIFO algorithm is performed by the CATS algorithm as of MySQL 8.0.20. In some cases, this change may affect the order in which transactions are granted locks.

You can view transaction scheduling weights by querying the `TRX_SCHEDULE_WEIGHT` column in the Information Schema `INNODB_TRX` table. Weights are computed for waiting transactions only. Waiting transactions are those in a `LOCK_WAIT` transaction execution state, as reported by the `TRX_STATE` column. A transaction that is not waiting for a lock reports a NULL `TRX_SCHEDULE_WEIGHT` value.

`INNODB_METRICS` counters are provided for monitoring of code-level transaction scheduling events. For information about using `INNODB_METRICS` counters, see [Section 15.15.6, “InnoDB INFORMATION_SCHEMA Metrics Table”](#).

- `lock_rec_release_attempts`

The number of attempts to release record locks. A single attempt may lead to zero or more record locks being released, as there may be zero or more record locks in a single structure.

- `lock_rec_grant_attempts`

The number of attempts to grant record locks. A single attempt may result in zero or more record locks being granted.

- `lock_schedule_refreshes`

The number of times the wait-for graph was analyzed to update the scheduled transaction weights.

15.8 InnoDB Configuration

This section provides configuration information and procedures for InnoDB initialization, startup, and various components and features of the InnoDB storage engine. For information about optimizing database operations for InnoDB tables, see [Section 8.5, “Optimizing for InnoDB Tables”](#).

15.8.1 InnoDB Startup Configuration

The first decisions to make about InnoDB configuration involve the configuration of data files, log files, page size, and memory buffers, which should be configured before initializing InnoDB. Modifying the configuration after InnoDB is initialized may involve non-trivial procedures.

This section provides information about specifying InnoDB settings in a configuration file, viewing InnoDB initialization information, and important storage considerations.

- [Specifying Options in a MySQL Option File](#)
- [Viewing InnoDB Initialization Information](#)
- [Important Storage Considerations](#)
- [System Tablespace Data File Configuration](#)
- [InnoDB Doublewrite Buffer File Configuration](#)
- [Redo Log Configuration](#)
- [Undo Tablespace Configuration](#)

- [Global Temporary Tablespace Configuration](#)
- [Session Temporary Tablespace Configuration](#)
- [Page Size Configuration](#)
- [Memory Configuration](#)

Specifying Options in a MySQL Option File

Because MySQL uses data file, log file, and page size settings to initialize [InnoDB](#), it is recommended that you define these settings in an option file that MySQL reads at startup, prior to initializing [InnoDB](#). Normally, [InnoDB](#) is initialized when the MySQL server is started for the first time.

You can place [InnoDB](#) options in the `[mysqld]` group of any option file that your server reads when it starts. The locations of MySQL option files are described in [Section 4.2.2.2, “Using Option Files”](#).

To make sure that `mysqld` reads options only from a specific file (and `mysqld-auto.cnf`), use the `--defaults-file` option as the first option on the command line when starting the server:

```
mysqld --defaults-file=path_to_option_file
```

Viewing InnoDB Initialization Information

To view [InnoDB](#) initialization information during startup, start `mysqld` from a command prompt, which prints initialization information to the console.

For example, on Windows, if `mysqld` is located in `C:\Program Files\MySQL\MySQL Server 8.0\bin`, start the MySQL server like this:

```
C:\> "C:\Program Files\MySQL\MySQL Server 8.0\bin\mysqld" --console
```

On Unix-like systems, `mysqld` is located in the `bin` directory of your MySQL installation:

```
$> bin/mysqld --user=mysql &
```

If you do not send server output to the console, check the error log after startup to see the initialization information [InnoDB](#) printed during the startup process.

For information about starting MySQL using other methods, see [Section 2.9.5, “Starting and Stopping MySQL Automatically”](#).



Note

[InnoDB](#) does not open all user tables and associated data files at startup. However, [InnoDB](#) does check for the existence of tablespace files referenced in the data dictionary. If a tablespace file is not found, [InnoDB](#) logs an error and continues the startup sequence. Tablespace files referenced in the redo log may be opened during crash recovery for redo application.

Important Storage Considerations

Review the following storage-related considerations before proceeding with your startup configuration.

- In some cases, you can improve database performance by placing data and log files on separate physical disks. You can also use raw disk partitions (raw devices) for [InnoDB](#) data files, which may speed up I/O. See [Using Raw Disk Partitions for the System Tablespace](#).
- [InnoDB](#) is a transaction-safe (ACID compliant) storage engine with commit, rollback, and crash-recovery capabilities to protect user data. **However, it cannot do so** if the underlying operating system or hardware does not work as advertised. Many operating systems or disk subsystems may delay or reorder write operations to improve performance. On some operating systems, the very `fsync()` system call that should wait until all unwritten data for a file has been flushed might actually return before the data has been flushed to stable storage. Because of this, an operating system

crash or a power outage may destroy recently committed data, or in the worst case, even corrupt the database because write operations have been reordered. If data integrity is important to you, perform “pull-the-plug” tests before using anything in production. On macOS, InnoDB uses a special `fcntl()` file flush method. Under Linux, it is advisable to **disable the write-back cache**.

On ATA/SATA disk drives, a command such `hdparm -W0 /dev/hda` may work to disable the write-back cache. **Beware that some drives or disk controllers may be unable to disable the write-back cache.**

- With regard to InnoDB recovery capabilities that protect user data, InnoDB uses a file flush technique involving a structure called the **doublewrite buffer**, which is enabled by default (`innodb_doublewrite=ON`). The doublewrite buffer adds safety to recovery following an unexpected exit or power outage, and improves performance on most varieties of Unix by reducing the need for `fsync()` operations. It is recommended that the `innodb_doublewrite` option remains enabled if you are concerned with data integrity or possible failures. For information about the doublewrite buffer, see [Section 15.11.1, “InnoDB Disk I/O”](#).
- Before using NFS with InnoDB, review potential issues outlined in [Using NFS with MySQL](#).

System Tablespace Data File Configuration

The `innodb_data_file_path` option defines the name, size, and attributes of InnoDB system tablespace data files. If you do not configure this option prior to initializing the MySQL server, the default behavior is to create a single auto-extending data file, slightly larger than 12MB, named `ibdata1`:

```
mysql> SHOW VARIABLES LIKE 'innodb_data_file_path';
+-----+-----+
| Variable_name      | Value           |
+-----+-----+
| innodb_data_file_path | ibdata1:12M:autoextend |
+-----+-----+
```

The full data file specification syntax includes the file name, file size, `autoextend` attribute, and `max` attribute:

```
file_name:file_size[:autoextend[:max:max_file_size]]
```

File sizes are specified in kilobytes, megabytes, or gigabytes by appending `K`, `M` or `G` to the size value. If specifying the data file size in kilobytes, do so in multiples of 1024. Otherwise, kilobyte values are rounded to nearest megabyte (MB) boundary. The sum of file sizes must be, at a minimum, slightly larger than 12MB.

You can specify more than one data file using a semicolon-separated list. For example:

```
[mysqld]
innodb_data_file_path=ibdata1:50M;ibdata2:50M:autoextend
```

The `autoextend` and `max` attributes can be used only for the data file that is specified last.

When the `autoextend` attribute is specified, the data file automatically increases in size by 64MB increments as space is required. The `innodb_autoextend_increment` variable controls the increment size.

To specify a maximum size for an auto-extending data file, use the `max` attribute following the `autoextend` attribute. Use the `max` attribute only in cases where constraining disk usage is of critical importance. The following configuration permits `ibdata1` to grow to a limit of 500MB:

```
[mysqld]
innodb_data_file_path=ibdata1:12M:autoextend:max:500M
```

A minimum file size is enforced for the *first* system tablespace data file to ensure that there is enough space for doublewrite buffer pages. The following table shows minimum file sizes for each InnoDB page size. The default InnoDB page size is 16384 (16KB).

Page Size (innodb_page_size)	Minimum File Size
16384 (16KB) or less	3MB
32768 (32KB)	6MB
65536 (64KB)	12MB

If your disk becomes full, you can add a data file on another disk. For instructions, see [Resizing the System Tablespace](#).

The size limit for individual files is determined by your operating system. You can set the file size to more than 4GB on operating systems that support large files. You can also use raw disk partitions as data files. See [Using Raw Disk Partitions for the System Tablespace](#).

InnoDB is not aware of the file system maximum file size, so be cautious on file systems where the maximum file size is a small value such as 2GB.

System tablespace files are created in the data directory by default (`datadir`). To specify an alternate location, use the `innodb_data_home_dir` option. For example, to create a system tablespace data file in a directory named `myibdata`, use this configuration:

```
[mysqld]
innodb_data_home_dir = /myibdata/
innodb_data_file_path=ibdata1:50M:autoextend
```

A trailing slash is required when specifying a value for `innodb_data_home_dir`. InnoDB does not create directories, so ensure that the specified directory exists before you start the server. Also, ensure sure that the MySQL server has the proper access rights to create files in the directory.

InnoDB forms the directory path for each data file by textually concatenating the value of `innodb_data_home_dir` to the data file name. If `innodb_data_home_dir` is not defined, the default value is “.”, which is the data directory. (The MySQL server changes its current working directory to the data directory when it begins executing.)

Alternatively, you can specify an absolute path for system tablespace data files. The following configuration is equivalent to the preceding one:

```
[mysqld]
innodb_data_file_path=/myibdata/ibdata1:50M:autoextend
```

When you specify an absolute path for `innodb_data_file_path`, the setting is not concatenated with the `innodb_data_home_dir` setting. System tablespace files are created in the specified absolute path. The specified directory must exist before you start the server.

InnoDB Doublewrite Buffer File Configuration

As of MySQL 8.0.20, the doublewrite buffer storage area resides in doublewrite files, which provides flexibility with respect to the storage location of doublewrite pages. In previous releases, the doublewrite buffer storage area resided in the system tablespace. The `innodb_doublewrite_dir` variable defines the directory where InnoDB creates doublewrite files at startup. If no directory is specified, doublewrite files are created in the `innodb_data_home_dir` directory, which defaults to the data directory if unspecified.

To have doublewrite files created in a location other than the `innodb_data_home_dir` directory, configure `innodb_doublewrite_dir` variable. For example:

```
innodb_doublewrite_dir=/path/to/doublewrite_directory
```

Other doublewrite buffer variables permit defining the number of doublewrite files, the number of pages per thread, and the doublewrite batch size. For more information about doublewrite buffer configuration, see [Section 15.6.4, “Doublewrite Buffer”](#).

Redo Log Configuration

From MySQL 8.0.30, the amount of disk space occupied by redo log files is controlled by the `innodb_redo_log_capacity` variable, which can be set at startup or runtime; for example, to set the variable to 8GB in an option file, add the following entry:

```
[mysqld]
innodb_redo_log_capacity = 8589934592
```

For information about configuring redo log capacity at runtime, see [Configuring Redo Log Capacity \(MySQL 8.0.30 or Higher\)](#).

The `innodb_redo_log_capacity` variable supersedes the `innodb_log_file_size` and `innodb_log_files_in_group` variables, which are deprecated. When the `innodb_redo_log_capacity` setting is defined, the `innodb_log_file_size` and `innodb_log_files_in_group` settings are ignored; otherwise, these settings are used to compute the `innodb_redo_log_capacity` setting (`innodb_log_files_in_group * innodb_log_file_size = innodb_redo_log_capacity`). If none of those variables are set, `innodb_redo_log_capacity` is set to the default value, which is 104857600 bytes (100MB). The maximum setting is 128GB.

From MySQL 8.0.30, `InnoDB` attempts to maintain 32 redo log files, with each file equal to 1/32 * `innodb_redo_log_capacity`. The redo log files reside in the `#innodb_redo` directory in the data directory unless a different directory was specified by the `innodb_log_group_home_dir` variable. If `innodb_log_group_home_dir` was defined, the redo log files reside in the `#innodb_redo` directory in that directory. For more information, see [Section 15.6.5, “Redo Log”](#).

Before MySQL 8.0.30, `InnoDB` creates two 5MB redo log files named `ib_logfile0` and `ib_logfile1` in the data directory by default. You can define a different number of redo log files and different redo log file size when initializing the MySQL Server instance by configuring the `innodb_log_files_in_group` and `innodb_log_file_size` variables.

- `innodb_log_files_in_group` defines the number of log files in the log group. The default and recommended value is 2.
- `innodb_log_file_size` defines the size in bytes of each log file in the log group. The combined log file size (`innodb_log_file_size * innodb_log_files_in_group`) cannot exceed the maximum value, which is slightly less than 512GB. A pair of 255 GB log files, for example, approaches the limit but does not exceed it. The default log file size is 48MB. Generally, the combined size of the log files should be large enough that the server can smooth out peaks and troughs in workload activity, which often means that there is enough redo log space to handle more than an hour of write activity. A larger log file size means less checkpoint flush activity in the buffer pool, which reduces disk I/O. For additional information, see [Section 8.5.4, “Optimizing InnoDB Redo Logging”](#).

The `innodb_log_group_home_dir` defines directory path to the `InnoDB` log files. You might use this option to place `InnoDB` redo log files in a different physical storage location than `InnoDB` data files to avoid potential I/O resource conflicts; for example:

```
[mysqld]
innodb_log_group_home_dir = /dr3/iblogs
```



Note

`InnoDB` does not create directories, so make sure that the log directory exists before you start the server. Use the Unix or DOS `mkdir` command to create any necessary directories.

Make sure that the MySQL server has the proper access rights to create files in the log directory. More generally, the server must have access rights in any directory where it needs to create files.

Undo Tablespace Configuration

Undo logs, by default, reside in two undo tablespaces created when the MySQL instance is initialized.

The `innodb_undo_directory` variable defines the path where InnoDB creates default undo tablespaces. If that variable is undefined, default undo tablespaces are created in the data directory. The `innodb_undo_directory` variable is not dynamic. Configuring it requires restarting the server.

The I/O patterns for undo logs make undo tablespaces good candidates for SSD storage.

For information about configuring additional undo tablespaces, see [Section 15.6.3.4, “Undo Tablespaces”](#).

Global Temporary Tablespace Configuration

The global temporary tablespace stores rollback segments for changes made to user-created temporary tables.

A single auto-extending global temporary tablespace data file named `ibtmp1` in the `innodb_data_home_dir` directory by default. The initial file size is slightly larger than 12MB.

The `innodb_temp_data_file_path` option specifies the path, file name, and file size for global temporary tablespace data files. File size is specified in KB, MB, or GB by appending K, M, or G to the size value. The file size or combined file size must be slightly larger than 12MB.

To specify an alternate location for global temporary tablespace data files, configure the `innodb_temp_data_file_path` option at startup.

Session Temporary Tablespace Configuration

In MySQL 8.0.15 and earlier, session temporary tablespaces store user-created temporary tables and internal temporary tables created by the optimizer when InnoDB is configured as the on-disk storage engine for internal temporary tables (`internal_tmp_disk_storage_engine=InnoDB`). From MySQL 8.0.16, InnoDB is always used as the on-disk storage engine for internal temporary tables.

The `innodb_temp_tablespaces_dir` variable defines the location where InnoDB creates session temporary tablespaces. The default location is the `#innodb_temp` directory in the data directory.

To specify an alternate location for session temporary tablespaces, configure the `innodb_temp_tablespaces_dir` variable at startup. A fully qualified path or path relative to the data directory is permitted.

Page Size Configuration

The `innodb_page_size` option specifies the page size for all InnoDB tablespaces in a MySQL instance. This value is set when the instance is created and remains constant afterward. Valid values are 64KB, 32KB, 16KB (the default), 8KB, and 4KB. Alternatively, you can specify page size in bytes (65536, 32768, 16384, 8192, 4096).

The default 16KB page size is appropriate for a wide range of workloads, particularly for queries involving table scans and DML operations involving bulk updates. Smaller page sizes might be more efficient for OLTP workloads involving many small writes, where contention can be an issue when a single page contains many rows. Smaller pages can also be more efficient for SSD storage devices, which typically use small block sizes. Keeping the InnoDB page size close to the storage device block size minimizes the amount of unchanged data that is rewritten to disk.

Memory Configuration

MySQL allocates memory to various caches and buffers to improve performance of database operations. When allocating memory for InnoDB, always consider memory required by the operating system, memory allocated to other applications, and memory allocated for other MySQL buffers and caches. For example, if you use MyISAM tables, consider the amount of memory allocated for the key

buffer (`key_buffer_size`). For an overview of MySQL buffers and caches, see [Section 8.12.3.1, “How MySQL Uses Memory”](#).

Buffers specific to [InnoDB](#) are configured using the following parameters:

- `innodb_buffer_pool_size` defines size of the buffer pool, which is the memory area that holds cached data for [InnoDB](#) tables, indexes, and other auxiliary buffers. The size of the buffer pool is important for system performance, and it is typically recommended that `innodb_buffer_pool_size` is configured to 50 to 75 percent of system memory. The default buffer pool size is 128MB. For additional guidance, see [Section 8.12.3.1, “How MySQL Uses Memory”](#). For information about how to configure [InnoDB](#) buffer pool size, see [Section 15.8.3.1, “Configuring InnoDB Buffer Pool Size”](#). Buffer pool size can be configured at startup or dynamically.

On systems with a large amount of memory, you can improve concurrency by dividing the buffer pool into multiple buffer pool instances. The number of buffer pool instances is controlled by the `innodb_buffer_pool_instances` option. By default, [InnoDB](#) creates one buffer pool instance. The number of buffer pool instances can be configured at startup. For more information, see [Section 15.8.3.2, “Configuring Multiple Buffer Pool Instances”](#).

- `innodb_log_buffer_size` defines the size of the buffer that [InnoDB](#) uses to write to the log files on disk. The default size is 16MB. A large log buffer enables large transactions to run without writing the log to disk before the transactions commit. If you have transactions that update, insert, or delete many rows, you might consider increasing the size of the log buffer to save disk I/O. `innodb_log_buffer_size` can be configured at startup. For related information, see [Section 8.5.4, “Optimizing InnoDB Redo Logging”](#).



Warning

On 32-bit GNU/Linux x86, if memory usage is set too high, `glibc` may permit the process heap to grow over the thread stacks, causing a server failure. It is a risk if the memory allocated to the `mysqld` process for global and per-thread buffers and caches is close to or exceeds 2GB.

A formula similar to the following that calculates global and per-thread memory allocation for MySQL can be used to estimate MySQL memory usage. You may need to modify the formula to account for buffers and caches in your MySQL version and configuration. For an overview of MySQL buffers and caches, see [Section 8.12.3.1, “How MySQL Uses Memory”](#).

```
innodb_buffer_pool_size
+ key_buffer_size
+ max_connections*(sort_buffer_size+read_buffer_size+binlog_cache_size)
+ max_connections*2MB
```

Each thread uses a stack (often 2MB, but only 256KB in MySQL binaries provided by Oracle Corporation.) and in the worst case also uses `sort_buffer_size + read_buffer_size` additional memory.

On Linux, if the kernel is enabled for large page support, [InnoDB](#) can use large pages to allocate memory for its buffer pool. See [Section 8.12.3.3, “Enabling Large Page Support”](#).

15.8.2 Configuring InnoDB for Read-Only Operation

You can query [InnoDB](#) tables where the MySQL data directory is on read-only media by enabling the `--innodb-read-only` configuration option at server startup.

How to Enable

To prepare an instance for read-only operation, make sure all the necessary information is [flushed](#) to the data files before storing it on the read-only medium. Run the server with change buffering disabled (`innodb_change_buffering=0`) and do a [slow shutdown](#).

To enable read-only mode for an entire MySQL instance, specify the following configuration options at server startup:

- `--innodb-read-only=1`
- If the instance is on read-only media such as a DVD or CD, or the `/var` directory is not writeable by all: `--pid-file=path_on_writeable_media` and `--event-scheduler=disabled`
- `--innodb-temp-data-file-path`. This option specifies the path, file name, and file size for InnoDB temporary tablespace data files. The default setting is `ibtmp1:12M:autoextend`, which creates the `ibtmp1` temporary tablespace data file in the data directory. To prepare an instance for read-only operation, set `innodb_temp_data_file_path` to a location outside of the data directory. The path must be relative to the data directory. For example:

```
--innodb-temp-data-file-path=../../../../tmp/ibtmp1:12M:autoextend
```

As of MySQL 8.0, enabling `innodb_read_only` prevents table creation and drop operations for all storage engines. These operations modify data dictionary tables in the `mysql` system database, but those tables use the InnoDB storage engine and cannot be modified when `innodb_read_only` is enabled. The same restriction applies to any operation that modifies data dictionary tables, such as `ANALYZE TABLE` and `ALTER TABLE tbl_name ENGINE=engine_name`.

In addition, other tables in the `mysql` system database use the InnoDB storage engine in MySQL 8.0. Making those tables read only results in restrictions on operations that modify them. For example, `CREATE USER`, `GRANT`, `REVOKE`, and `INSTALL PLUGIN` operations are not permitted in read-only mode.

Usage Scenarios

This mode of operation is appropriate in situations such as:

- Distributing a MySQL application, or a set of MySQL data, on a read-only storage medium such as a DVD or CD.
- Multiple MySQL instances querying the same data directory simultaneously, typically in a data warehousing configuration. You might use this technique to avoid `bottlenecks` that can occur with a heavily loaded MySQL instance, or you might use different configuration options for the various instances to tune each one for particular kinds of queries.
- Querying data that has been put into a read-only state for security or data integrity reasons, such as archived backup data.



Note

This feature is mainly intended for flexibility in distribution and deployment, rather than raw performance based on the read-only aspect. See [Section 8.5.3, “Optimizing InnoDB Read-Only Transactions”](#) for ways to tune the performance of read-only queries, which do not require making the entire server read-only.

How It Works

When the server is run in read-only mode through the `--innodb-read-only` option, certain InnoDB features and components are reduced or turned off entirely:

- No `change buffering` is done, in particular no merges from the change buffer. To make sure the change buffer is empty when you prepare the instance for read-only operation, disable change buffering (`innodb_change_buffering=0`) and do a `slow shutdown` first.
- There is no `crash recovery` phase at startup. The instance must have performed a `slow shutdown` before being put into the read-only state.
- Because the `redo log` is not used in read-only operation, you can set `innodb_log_file_size` to the smallest size possible (1 MB) before making the instance read-only.

- Most background threads are turned off. I/O read threads remain, as well as I/O write threads and a page flush coordinator thread for writes to temporary files, which are permitted in read-only mode. A buffer pool resize thread also remains active to enable online resizing of the buffer pool.
- Information about deadlocks, monitor output, and so on is not written to temporary files. As a consequence, `SHOW ENGINE INNODB STATUS` does not produce any output.
- Changes to configuration option settings that would normally change the behavior of write operations, have no effect when the server is in read-only mode.
- The MVCC processing to enforce isolation levels is turned off. All queries read the latest version of a record, because update and deletes are not possible.
- The undo log is not used. Disable any settings for the `innodb_undo_tablespaces` and `innodb_undo_directory` configuration options.

15.8.3 InnoDB Buffer Pool Configuration

This section provides configuration and tuning information for the InnoDB buffer pool.

15.8.3.1 Configuring InnoDB Buffer Pool Size

You can configure InnoDB buffer pool size offline or while the server is running. Behavior described in this section applies to both methods. For additional information about configuring buffer pool size online, see [Configuring InnoDB Buffer Pool Size Online](#).

When increasing or decreasing `innodb_buffer_pool_size`, the operation is performed in chunks. Chunk size is defined by the `innodb_buffer_pool_chunk_size` configuration option, which has a default of `128M`. For more information, see [Configuring InnoDB Buffer Pool Chunk Size](#).

Buffer pool size must always be equal to or a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`. If you configure `innodb_buffer_pool_size` to a value that is not equal to or a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`, buffer pool size is automatically adjusted to a value that is equal to or a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`.

In the following example, `innodb_buffer_pool_size` is set to `8G`, and `innodb_buffer_pool_instances` is set to `16`. `innodb_buffer_pool_chunk_size` is `128M`, which is the default value.

`8G` is a valid `innodb_buffer_pool_size` value because `8G` is a multiple of `innodb_buffer_pool_instances=16 * innodb_buffer_pool_chunk_size=128M`, which is `2G`.

```
$> mysqld --innodb-buffer-pool-size=8G --innodb-buffer-pool-instances=16
mysql> SELECT @@innodb_buffer_pool_size/1024/1024/1024;
+-----+
| @@innodb_buffer_pool_size/1024/1024/1024 |
+-----+
| 8.000000000000000 |
+-----+
```

In this example, `innodb_buffer_pool_size` is set to `9G`, and `innodb_buffer_pool_instances` is set to `16`. `innodb_buffer_pool_chunk_size` is `128M`, which is the default value. In this case, `9G` is not a multiple of `innodb_buffer_pool_instances=16 * innodb_buffer_pool_chunk_size=128M`, so `innodb_buffer_pool_size` is adjusted to `10G`, which is a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`.

```
$> mysqld --innodb-buffer-pool-size=9G --innodb-buffer-pool-instances=16
mysql> SELECT @@innodb_buffer_pool_size/1024/1024/1024;
+-----+
| @@innodb_buffer_pool_size/1024/1024/1024 |
+-----+
```

	10.000000000000
--	-----------------

Configuring InnoDB Buffer Pool Chunk Size

`innodb_buffer_pool_chunk_size` can be increased or decreased in 1MB (1048576 byte) units but can only be modified at startup, in a command line string or in a MySQL configuration file.

Command line:

```
$> mysqld --innodb-buffer-pool-chunk-size=134217728
```

Configuration file:

```
[mysqld]
innodb_buffer_pool_chunk_size=134217728
```

The following conditions apply when altering `innodb_buffer_pool_chunk_size`:

- If the new `innodb_buffer_pool_chunk_size` value * `innodb_buffer_pool_instances` is larger than the current buffer pool size when the buffer pool is initialized, `innodb_buffer_pool_chunk_size` is truncated to `innodb_buffer_pool_size / innodb_buffer_pool_instances`.

For example, if the buffer pool is initialized with a size of **2GB** (2147483648 bytes), **4** buffer pool instances, and a chunk size of **1GB** (1073741824 bytes), chunk size is truncated to a value equal to `innodb_buffer_pool_size / innodb_buffer_pool_instances`, as shown below:

```
$> mysqld --innodb-buffer-pool-size=2147483648 --innodb-buffer-pool-instances=4
--innodb-buffer-pool-chunk-size=1073741824;
```

```
mysql> SELECT @@innodb_buffer_pool_size;
+-----+
| @@innodb_buffer_pool_size |
+-----+
| 2147483648 |
+-----+

mysql> SELECT @@innodb_buffer_pool_instances;
+-----+
| @@innodb_buffer_pool_instances |
+-----+
| 4 |
+-----+

# Chunk size was set to 1GB (1073741824 bytes) on startup but was
# truncated to innodb_buffer_pool_size / innodb_buffer_pool_instances

mysql> SELECT @@innodb_buffer_pool_chunk_size;
+-----+
| @@innodb_buffer_pool_chunk_size |
+-----+
| 536870912 |
+-----+
```

- Buffer pool size must always be equal to or a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`. If you alter `innodb_buffer_pool_chunk_size`, `innodb_buffer_pool_size` is automatically adjusted to a value that is equal to or a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`. The adjustment occurs when the buffer pool is initialized. This behavior is demonstrated in the following example:

```
# The buffer pool has a default size of 128MB (134217728 bytes)

mysql> SELECT @@innodb_buffer_pool_size;
+-----+
| @@innodb_buffer_pool_size |
+-----+
| 134217728 |
+-----+
```

```
# The chunk size is also 128MB (134217728 bytes)

mysql> SELECT @@innodb_buffer_pool_chunk_size;
+-----+
| @@innodb_buffer_pool_chunk_size |
+-----+
| 134217728 |
+-----+

# There is a single buffer pool instance

mysql> SELECT @@innodb_buffer_pool_instances;
+-----+
| @@innodb_buffer_pool_instances |
+-----+
| 1 |
+-----+

# Chunk size is decreased by 1MB (1048576 bytes) at startup
# (134217728 - 1048576 = 133169152):

$> mysqld --innodb-buffer-pool-chunk-size=133169152

mysql> SELECT @@innodb_buffer_pool_chunk_size;
+-----+
| @@innodb_buffer_pool_chunk_size |
+-----+
| 133169152 |
+-----+

# Buffer pool size increases from 134217728 to 266338304
# Buffer pool size is automatically adjusted to a value that is equal to
# or a multiple of innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances

mysql> SELECT @@innodb_buffer_pool_size;
+-----+
| @@innodb_buffer_pool_size |
+-----+
| 266338304 |
+-----+
```

This example demonstrates the same behavior but with multiple buffer pool instances:

```
# The buffer pool has a default size of 2GB (2147483648 bytes)

mysql> SELECT @@innodb_buffer_pool_size;
+-----+
| @@innodb_buffer_pool_size |
+-----+
| 2147483648 |
+-----+

# The chunk size is .5 GB (536870912 bytes)

mysql> SELECT @@innodb_buffer_pool_chunk_size;
+-----+
| @@innodb_buffer_pool_chunk_size |
+-----+
| 536870912 |
+-----+

# There are 4 buffer pool instances

mysql> SELECT @@innodb_buffer_pool_instances;
+-----+
| @@innodb_buffer_pool_instances |
+-----+
| 4 |
+-----+

# Chunk size is decreased by 1MB (1048576 bytes) at startup
```

```
# (536870912 - 1048576 = 535822336):

$> mysqld --innodb-buffer-pool-chunk-size=535822336

mysql> SELECT @@innodb_buffer_pool_chunk_size;
+-----+
| @@innodb_buffer_pool_chunk_size |
+-----+
| 535822336 |
+-----+

# Buffer pool size increases from 2147483648 to 4286578688
# Buffer pool size is automatically adjusted to a value that is equal to
# or a multiple of innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances

mysql> SELECT @@innodb_buffer_pool_size;
+-----+
| @@innodb_buffer_pool_size |
+-----+
| 4286578688 |
+-----+
```

Care should be taken when changing `innodb_buffer_pool_chunk_size`, as changing this value can increase the size of the buffer pool, as shown in the examples above. Before you change `innodb_buffer_pool_chunk_size`, calculate the effect on `innodb_buffer_pool_size` to ensure that the resulting buffer pool size is acceptable.



Note

To avoid potential performance issues, the number of chunks (`innodb_buffer_pool_size / innodb_buffer_pool_chunk_size`) should not exceed 1000.

Configuring InnoDB Buffer Pool Size Online

The `innodb_buffer_pool_size` configuration option can be set dynamically using a `SET` statement, allowing you to resize the buffer pool without restarting the server. For example:

```
mysql> SET GLOBAL innodb_buffer_pool_size=402653184;
```



Note

The buffer pool size must be equal to or a multiple of `innodb_buffer_pool_chunk_size * innodb_buffer_pool_instances`. Changing those variable settings requires restarting the server.

Active transactions and operations performed through `InnoDB` APIs should be completed before resizing the buffer pool. When initiating a resizing operation, the operation does not start until all active transactions are completed. Once the resizing operation is in progress, new transactions and operations that require access to the buffer pool must wait until the resizing operation finishes. The exception to the rule is that concurrent access to the buffer pool is permitted while the buffer pool is defragmented and pages are withdrawn when buffer pool size is decreased. A drawback of allowing concurrent access is that it could result in a temporary shortage of available pages while pages are being withdrawn.



Note

Nested transactions could fail if initiated after the buffer pool resizing operation begins.

Monitoring Online Buffer Pool Resizing Progress

The `Innodb_buffer_pool_resize_status` variable reports a string value indicating buffer pool resizing progress; for example:

```
mysql> SHOW STATUS WHERE Variable_name='InnoDB_buffer_pool_resize_status';
+-----+-----+
| Variable_name | Value |
+-----+-----+
| Innodb_buffer_pool_resize_status | Resizing also other hash tables. |
+-----+-----+
```

From MySQL 8.0.31, you can also monitor an online buffer pool resizing operation using the [Innodb_buffer_pool_resize_status_code](#) and [Innodb_buffer_pool_resize_status_progress](#) status variables, which report numeric values, preferable for programmatic monitoring.

The [Innodb_buffer_pool_resize_status_code](#) status variable reports a status code indicating the stage of an online buffer pool resizing operation. Status codes include:

- 0: No Resize operation in progress
- 1: Starting Resize
- 2: Disabling AHI (Adaptive Hash Index)
- 3: Withdrawing Blocks
- 4: Acquiring Global Lock
- 5: Resizing Pool
- 6: Resizing Hash
- 7: Resizing Failed

The [Innodb_buffer_pool_resize_status_progress](#) status variable reports a percentage value indicating the progress of each stage. The percentage value is updated after each buffer pool instance is processed. As the status (reported by [Innodb_buffer_pool_resize_status_code](#)) changes from one status to another, the percentage value is reset to 0.

The following query returns a string value indicating the buffer pool resizing progress, a code indicating the current stage of the operation, and the current progress of that stage, expressed as a percentage value:

```
SELECT variable_name, variable_value
  FROM performance_schema.global_status
 WHERE LOWER(variable_name) LIKE "innodb_buffer_pool_resize%";
```

Buffer pool resizing progress is also visible in the server error log. This example shows notes that are logged when increasing the size of the buffer pool:

```
[Note] InnoDB: Resizing buffer pool from 134217728 to 4294967296. (unit=134217728)
[Note] InnoDB: disabled adaptive hash index.
[Note] InnoDB: buffer pool 0 : 31 chunks (253952 blocks) was added.
[Note] InnoDB: buffer pool 0 : hash tables were resized.
[Note] InnoDB: Resized hash tables at lock_sys, adaptive hash index, dictionary.
[Note] InnoDB: completed to resize buffer pool from 134217728 to 4294967296.
[Note] InnoDB: re-enabled adaptive hash index.
```

This example shows notes that are logged when decreasing the size of the buffer pool:

```
[Note] InnoDB: Resizing buffer pool from 4294967296 to 134217728. (unit=134217728)
[Note] InnoDB: disabled adaptive hash index.
[Note] InnoDB: buffer pool 0 : start to withdraw the last 253952 blocks.
[Note] InnoDB: buffer pool 0 : withdrew 253952 blocks from free list. tried to relocate
0 pages. (253952/253952)
[Note] InnoDB: buffer pool 0 : withdrawn target 253952 blocks.
[Note] InnoDB: buffer pool 0 : 31 chunks (253952 blocks) was freed.
[Note] InnoDB: buffer pool 0 : hash tables were resized.
[Note] InnoDB: Resized hash tables at lock_sys, adaptive hash index, dictionary.
[Note] InnoDB: completed to resize buffer pool from 4294967296 to 134217728.
[Note] InnoDB: re-enabled adaptive hash index.
```

From MySQL 8.0.31, starting the server with `--log-error-verbosity=3` logs additional information to the error log during an online buffer pool resizing operation. Additional information includes the status codes reported by `Innodb_buffer_pool_resize_status_code` and the percentage progress value reported by `Innodb_buffer_pool_resize_status_progress`.

```
[Note] [MY-012398] [InnoDB] Requested to resize buffer pool. (new size: 1073741824 bytes)
[Note] [MY-013954] [InnoDB] Status code 1: Resizing buffer pool from 134217728 to 1073741824
(unit=134217728).
[Note] [MY-013953] [InnoDB] Status code 1: 100% complete
[Note] [MY-013952] [InnoDB] Status code 1: Completed
[Note] [MY-013954] [InnoDB] Status code 2: Disabling adaptive hash index.
[Note] [MY-011885] [InnoDB] disabled adaptive hash index.
[Note] [MY-013953] [InnoDB] Status code 2: 100% complete
[Note] [MY-013952] [InnoDB] Status code 2: Completed
[Note] [MY-013954] [InnoDB] Status code 3: Withdrawing blocks to be shrunken.
[Note] [MY-013953] [InnoDB] Status code 3: 100% complete
[Note] [MY-013952] [InnoDB] Status code 3: Completed
[Note] [MY-013954] [InnoDB] Status code 4: Latching whole of buffer pool.
[Note] [MY-013953] [InnoDB] Status code 4: 14% complete
[Note] [MY-013953] [InnoDB] Status code 4: 28% complete
[Note] [MY-013953] [InnoDB] Status code 4: 42% complete
[Note] [MY-013953] [InnoDB] Status code 4: 57% complete
[Note] [MY-013953] [InnoDB] Status code 4: 71% complete
[Note] [MY-013953] [InnoDB] Status code 4: 85% complete
[Note] [MY-013953] [InnoDB] Status code 4: 100% complete
[Note] [MY-013952] [InnoDB] Status code 4: Completed
[Note] [MY-013954] [InnoDB] Status code 5: Starting pool resize
[Note] [MY-013954] [InnoDB] Status code 5: buffer pool 0 : resizing with chunks 1 to 8.
[Note] [MY-011891] [InnoDB] buffer pool 0 : 7 chunks (57339 blocks) were added.
[Note] [MY-013953] [InnoDB] Status code 5: 100% complete
[Note] [MY-013952] [InnoDB] Status code 5: Completed
[Note] [MY-013954] [InnoDB] Status code 6: Resizing hash tables.
[Note] [MY-011892] [InnoDB] buffer pool 0 : hash tables were resized.
[Note] [MY-013953] [InnoDB] Status code 6: 100% complete
[Note] [MY-013954] [InnoDB] Status code 6: Resizing also other hash tables.
[Note] [MY-011893] [InnoDB] Resized hash tables at lock_sys, adaptive hash index, dictionary.
[Note] [MY-011894] [InnoDB] Completed to resize buffer pool from 134217728 to 1073741824.
[Note] [MY-011895] [InnoDB] Re-enabled adaptive hash index.
[Note] [MY-013952] [InnoDB] Status code 6: Completed
[Note] [MY-013954] [InnoDB] Status code 0: Completed resizing buffer pool at 220826 6:25:46.
[Note] [MY-013953] [InnoDB] Status code 0: 100% complete
```

Online Buffer Pool Resizing Internals

The resizing operation is performed by a background thread. When increasing the size of the buffer pool, the resizing operation:

- Adds pages in `chunks` (chunk size is defined by `innodb_buffer_pool_chunk_size`)
- Converts hash tables, lists, and pointers to use new addresses in memory
- Adds new pages to the free list

While these operations are in progress, other threads are blocked from accessing the buffer pool.

When decreasing the size of the buffer pool, the resizing operation:

- Defragments the buffer pool and withdraws (frees) pages
- Removes pages in `chunks` (chunk size is defined by `innodb_buffer_pool_chunk_size`)
- Converts hash tables, lists, and pointers to use new addresses in memory

Of these operations, only defragmenting the buffer pool and withdrawing pages allow other threads to access to the buffer pool concurrently.

15.8.3.2 Configuring Multiple Buffer Pool Instances

For systems with buffer pools in the multi-gigabyte range, dividing the buffer pool into separate instances can improve concurrency, by reducing contention as different threads read and write to cached pages. This feature is typically intended for systems with a [buffer pool](#) size in the multi-gigabyte range. Multiple buffer pool instances are configured using the [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. 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, LRU, and all other data structures connected to a buffer pool. Prior to MySQL 8.0, each buffer pool was protected by its own buffer pool mutex. In MySQL 8.0 and later, the buffer pool mutex was replaced by several list and hash protecting mutexes, to reduce contention.

To enable multiple buffer pool instances, 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 [innodb_buffer_pool_size](#) to a size of 1GB 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 1GB.

For information about modifying [InnoDB](#) buffer pool size, see [Section 15.8.3.1, “Configuring InnoDB Buffer Pool Size”](#).

15.8.3.3 Making the Buffer Pool Scan Resistant

Rather than using a strict [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 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 specifics about the LRU algorithm, see [Section 15.5.1, “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 [1000](#). 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](#) can be specified in the MySQL option file ([my.cnf](#) or [my.ini](#)) or changed at runtime with the [SET GLOBAL](#) statement. Changing the value at runtime requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#).

To help you gauge the effect of setting these parameters, the `SHOW ENGINE INNODB STATUS` command reports buffer pool statistics. For details, see [Monitoring the Buffer Pool Using the InnoDB Standard Monitor](#).

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.

15.8.3.4 Configuring InnoDB Buffer Pool Prefetching (Read-Ahead)

A [read-ahead](#) request is an I/O request to prefetch multiple pages in the [buffer pool](#) asynchronously, in anticipation of impending need for these pages. The requests bring in all the pages in one [extent](#). [InnoDB](#) uses 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 the last page of the current extent.

The configuration parameter `innodb_read_ahead_threshold` controls how sensitive [InnoDB](#) is in detecting patterns of sequential page access. If the number of pages read sequentially from an extent is greater than or equal to `innodb_read_ahead_threshold`, [InnoDB](#) initiates an asynchronous read-ahead operation of the entire following extent. `innodb_read_ahead_threshold` can be set to any value from 0-64. The default value is 56. 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](#) triggers an asynchronous read-ahead even if as few as 8 pages in the extent are accessed sequentially. You can set the value of this parameter in the MySQL [configuration file](#), or change it dynamically with the `SET GLOBAL` statement, which requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#).

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. To enable this feature, set the configuration variable `innodb_random_read_ahead` to `ON`.

The `SHOW ENGINE INNODB STATUS` command displays statistics to help you evaluate the effectiveness of the read-ahead algorithm. Statistics include counter information for the following global status variables:

- `Innodb_buffer_pool_read_ahead`
- `Innodb_buffer_pool_read_ahead_evicted`
- `Innodb_buffer_pool_read_ahead_rnd`

This information can be useful when fine-tuning the `innodb_random_read_ahead` setting.

For more information about I/O performance, see [Section 8.5.8, “Optimizing InnoDB Disk I/O”](#) and [Section 8.12.1, “Optimizing Disk I/O”](#).

15.8.3.5 Configuring Buffer Pool Flushing

InnoDB performs certain tasks in the background, including flushing of dirty pages from the buffer pool. Dirty pages are those that have been modified but are not yet written to the data files on disk.

In MySQL 8.0, buffer pool flushing is performed by page cleaner threads. The number of page cleaner threads is controlled by the `innodb_page_cleaners` variable, which has a default value of 4. However, if the number of page cleaner threads exceeds the number of buffer pool instances, `innodb_page_cleaners` is automatically set to the same value as `innodb_buffer_pool_instances`.

Buffer pool flushing is initiated when the percentage of dirty pages reaches the low water mark value defined by the `innodb_max_dirty_pages_pct_lwm` variable. The default low water mark is 10% of buffer pool pages. A `innodb_max_dirty_pages_pct_lwm` value of 0 disables this early flushing behaviour.

The purpose of the `innodb_max_dirty_pages_pct_lwm` threshold is to control the percentage dirty pages in the buffer pool and to prevent the amount of dirty pages from reaching the threshold defined by the `innodb_max_dirty_pages_pct` variable, which has a default value of 90. InnoDB aggressively flushes buffer pool pages if the percentage of dirty pages in the buffer pool reaches the `innodb_max_dirty_pages_pct` threshold.

When configuring `innodb_max_dirty_pages_pct_lwm`, the value should always be lower than the `innodb_max_dirty_pages_pct` value.

Additional variables permit fine-tuning of buffer pool flushing behavior:

- The `innodb_flush_neighbors` variable defines whether flushing a page from the buffer pool also flushes other dirty pages in the same extent.
 - The default setting of 0 disables `innodb_flush_neighbors`. Dirty pages in the same extent are not flushed. This setting is recommended for non-rotational storage (SSD) devices where seek time is not a significant factor.
 - A setting of 1 flushes contiguous dirty pages in the same extent.
 - A setting of 2 flushes dirty pages in the same extent.

When table data is stored on a traditional HDD storage device, flushing neighbor pages in one operation reduces I/O overhead (primarily for disk seek operations) compared to flushing individual pages at different times. For table data stored on SSD, seek time is not a significant factor and you can disable this setting to spread out write operations.

- The `innodb_lru_scan_depth` variable specifies, per buffer pool instance, how far down the buffer pool LRU list the page cleaner thread scans looking for dirty pages to flush. This is a background operation performed by a page cleaner thread once per second.

A setting smaller than the default is generally suitable for most workloads. A value that is significantly higher than necessary may impact performance. Only consider increasing the value if you have spare I/O capacity under a typical workload. Conversely, if a write-intensive workload saturates your I/O capacity, decrease the value, especially in the case of a large buffer pool.

When tuning `innodb_lru_scan_depth`, start with a low value and configure the setting upward with the goal of rarely seeing zero free pages. Also, consider adjusting `innodb_lru_scan_depth` when changing the number of buffer pool instances, since `innodb_lru_scan_depth` * `innodb_buffer_pool_instances` defines the amount of work performed by the page cleaner thread each second.

The `innodb_flush_neighbors` and `innodb_lru_scan_depth` variables are primarily intended for write-intensive workloads. With heavy DML activity, flushing can fall behind if it is not aggressive enough, or disk writes can saturate I/O capacity if flushing is too aggressive. The ideal settings depend on your workload, data access patterns, and storage configuration (for example, whether data is stored on HDD or SSD devices).

Adaptive Flushing

InnoDB uses an adaptive flushing algorithm to dynamically adjust the 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 flushing activity keeps pace with the current workload. Automatically adjusting the flushing rate helps avoid sudden dips in throughput that can occur when bursts of I/O activity due to buffer pool flushing affects the I/O capacity available for ordinary read and write activity.

Sharp checkpoints, which are typically associated with write-intensive workloads that generate a lot of redo entries, can cause a sudden change in throughput, for example. A sharp checkpoint occurs when **InnoDB** wants to reuse a portion of a log file. Before doing so, all dirty pages with redo entries in that portion of the log file must be flushed. If log files become full, a sharp checkpoint occurs, causing a temporary reduction in throughput. This scenario can occur even if `innodb_max_dirty_pages_pct` threshold is not reached.

The adaptive flushing algorithm helps avoid such scenarios by tracking the number of dirty pages in the buffer pool and the rate at which redo log records are being generated. Based on this information, it decides how many dirty pages to flush from the buffer pool each second, which permits it to manage sudden changes in workload.

The `innodb_adaptive_flushing_lwm` variable defines a low water mark for redo log capacity. When that threshold is crossed, adaptive flushing is enabled, even if the `innodb_adaptive_flushing` variable is disabled.

Internal benchmarking has shown that the algorithm not only maintains throughput over time, but can also improve overall throughput significantly. However, adaptive flushing can affect the I/O pattern of a workload significantly and may not be appropriate in all cases. It gives the most benefit when the redo log is in danger of filling up. If adaptive flushing is not appropriate to the characteristics of your workload, you can disable it. Adaptive flushing controlled by the `innodb_adaptive_flushing` variable, which is enabled by default.

`innodb_flushing_avg_loops` defines the number of iterations that **InnoDB** keeps the previously calculated snapshot of the flushing state, controlling how quickly adaptive flushing responds to foreground workload changes. A high `innodb_flushing_avg_loops` value means that **InnoDB** keeps the previously calculated snapshot longer, so adaptive flushing responds more slowly. When setting a high value it is important to ensure that redo log utilization does not reach 75% (the hardcoded limit at which asynchronous flushing starts), and that the `innodb_max_dirty_pages_pct` threshold keeps the number of dirty pages to a level that is appropriate for the workload.

Systems with consistent workloads, a large log file size (`innodb_log_file_size`), and small spikes that do not reach 75% log space utilization should use a high `innodb_flushing_avg_loops` value to keep flushing as smooth as possible. For systems with extreme load spikes or log files that do not provide a lot of space, a smaller value allows flushing to closely track workload changes, and helps to avoid reaching 75% log space utilization.

Be aware that if flushing falls behind, the rate of buffer pool flushing can exceed the I/O capacity available to **InnoDB**, as defined by `innodb_io_capacity` setting. The `innodb_io_capacity_max`

value defines an upper limit on I/O capacity in such situations, so that a spike in I/O activity does not consume the entire I/O capacity of the server.

The `innodb_io_capacity` setting is applicable to all buffer pool instances. When dirty pages are flushed, I/O capacity is divided equally among buffer pool instances.

Limits Buffer Flushing During Idle Periods

As of MySQL 8.0.18, you can use the `innodb_idle_flush_pct` variable to limit the rate of buffer pool flushing during idle periods, which are periods of time that database pages are not modified. The `innodb_idle_flush_pct` value is a percentage of the `innodb_io_capacity` setting, which defines the number of I/O operations per second available to InnoDB. The default `innodb_idle_flush_pct` value is 100, which is 100 percent of the `innodb_io_capacity` setting. To limit flushing during idle periods, define an `innodb_idle_flush_pct` value less than 100.

Limiting page flushing during idle periods can help extend the life of solid state storage devices. Side effects of limiting page flushing during idle periods may include a longer shutdown time following a lengthy idle period, and a longer recovery period should a server failure occur.

15.8.3.6 Saving and Restoring the Buffer Pool State

To reduce the `warmup` period after restarting the server, InnoDB saves a percentage of the most recently used pages for each buffer pool at server shutdown and restores these pages at server startup. The percentage of recently used pages that is stored is defined by the `innodb_buffer_pool_dump_pct` configuration option.

After restarting a busy server, there is typically a warmup period with steadily increasing throughput, as disk pages that were in the buffer pool are brought back into memory (as the same data is queried, updated, and so on). The ability to restore the buffer pool at startup shortens the warmup period by reloading disk pages that were in the buffer pool before the restart rather than waiting for DML operations to access corresponding rows. Also, I/O requests can be performed in large batches, making the overall I/O faster. Page loading happens in the background, and does not delay database startup.

In addition to saving the buffer pool state at shutdown and restoring it at startup, you can save and restore the buffer pool state at any time, while the server is running. For example, you can save the state of the buffer pool after reaching a stable throughput under a steady workload. You could also restore the previous buffer pool state after running reports or maintenance jobs that bring data pages into the buffer pool that are only required for those operations, or after running some other non-typical workload.

Even though a buffer pool can be many gigabytes in size, the buffer pool data that InnoDB saves to disk is tiny by comparison. Only tablespace IDs and page IDs necessary to locate the appropriate pages are saved to disk. This information is derived from the `INNODB_BUFFER_PAGE_LRU_INFORMATION_SCHEMA` table. By default, tablespace ID and page ID data is saved in a file named `ib_buffer_pool`, which is saved to the InnoDB data directory. The file name and location can be modified using the `innodb_buffer_pool_filename` configuration parameter.

Because data is cached in and aged out of the buffer pool as it is with regular database operations, there is no problem if the disk pages are recently updated, or if a DML operation involves data that has not yet been loaded. The loading mechanism skips requested pages that no longer exist.

The underlying mechanism involves a background thread that is dispatched to perform the dump and load operations.

Disk pages from compressed tables are loaded into the buffer pool in their compressed form. Pages are uncompressed as usual when page contents are accessed during DML operations. Because uncompressing pages is a CPU-intensive process, it is more efficient for concurrency to perform the operation in a connection thread rather than in the single thread that performs the buffer pool restore operation.

Operations related to saving and restoring the buffer pool state are described in the following topics:

- [Configuring the Dump Percentage for Buffer Pool Pages](#)
- [Saving the Buffer Pool State at Shutdown and Restoring it at Startup](#)
- [Saving and Restoring the Buffer Pool State Online](#)
- [Displaying Buffer Pool Dump Progress](#)
- [Displaying Buffer Pool Load Progress](#)
- [Aborting a Buffer Pool Load Operation](#)
- [Monitoring Buffer Pool Load Progress Using Performance Schema](#)

Configuring the Dump Percentage for Buffer Pool Pages

Before dumping pages from the buffer pool, you can configure the percentage of most-recently-used buffer pool pages that you want to dump by setting the `innodb_buffer_pool_dump_pct` option. If you plan to dump buffer pool pages while the server is running, you can configure the option dynamically:

```
SET GLOBAL innodb_buffer_pool_dump_pct=40;
```

If you plan to dump buffer pool pages at server shutdown, set `innodb_buffer_pool_dump_pct` in your configuration file.

```
[mysqld]
innodb_buffer_pool_dump_pct=40
```

The `innodb_buffer_pool_dump_pct` default value is 25 (dump 25% of most-recently-used pages).

Saving the Buffer Pool State at Shutdown and Restoring it at Startup

To save the state of the buffer pool at server shutdown, issue the following statement prior to shutting down the server:

```
SET GLOBAL innodb_buffer_pool_dump_at_shutdown=ON;
```

`innodb_buffer_pool_dump_at_shutdown` is enabled by default.

To restore the buffer pool state at server startup, specify the `--innodb-buffer-pool-load-at-startup` option when starting the server:

```
mysqld --innodb-buffer-pool-load-at-startup=ON;
```

`innodb_buffer_pool_load_at_startup` is enabled by default.

Saving and Restoring the Buffer Pool State Online

To save the state of the buffer pool while MySQL server is running, issue the following statement:

```
SET GLOBAL innodb_buffer_pool_dump_now=ON;
```

To restore the buffer pool state while MySQL is running, issue the following statement:

```
SET GLOBAL innodb_buffer_pool_load_now=ON;
```

Displaying Buffer Pool Dump Progress

To display progress when saving the buffer pool state to disk, issue the following statement:

```
SHOW STATUS LIKE 'Innodb_buffer_pool_dump_status';
```

If the operation has not yet started, “not started” is returned. If the operation is complete, the completion time is printed (e.g. Finished at 110505 12:18:02). If the operation is in progress, status information is provided (e.g. Dumping buffer pool 5/7, page 237/2873).

Displaying Buffer Pool Load Progress

To display progress when loading the buffer pool, issue the following statement:

```
SHOW STATUS LIKE 'Innodb_buffer_pool_load_status';
```

If the operation has not yet started, “not started” is returned. If the operation is complete, the completion time is printed (e.g. Finished at 110505 12:23:24). If the operation is in progress, status information is provided (e.g. Loaded 123/22301 pages).

Aborting a Buffer Pool Load Operation

To abort a buffer pool load operation, issue the following statement:

```
SET GLOBAL innodb_buffer_pool_load_abort=ON;
```

Monitoring Buffer Pool Load Progress Using Performance Schema

You can monitor buffer pool load progress using [Performance Schema](#).

The following example demonstrates how to enable the `stage/innodb/buffer pool load` stage event instrument and related consumer tables to monitor buffer pool load progress.

For information about buffer pool dump and load procedures used in this example, see [Section 15.8.3.6, “Saving and Restoring the Buffer Pool State”](#). For information about Performance Schema stage event instruments and related consumers, see [Section 27.12.5, “Performance Schema Stage Event Tables”](#).

1. Enable the `stage/innodb/buffer pool load` instrument:

```
mysql> UPDATE performance_schema.setup_instruments SET ENABLED = 'YES'
      WHERE NAME LIKE 'stage/innodb/buffer%';
```

2. Enable the stage event consumer tables, which include `events_stages_current`, `events_stages_history`, and `events_stages_history_long`.

```
mysql> UPDATE performance_schema.setup_consumers SET ENABLED = 'YES'
      WHERE NAME LIKE '%stages%';
```

3. Dump the current buffer pool state by enabling `innodb_buffer_pool_dump_now`.

```
mysql> SET GLOBAL innodb_buffer_pool_dump_now=ON;
```

4. Check the buffer pool dump status to ensure that the operation has completed.

```
mysql> SHOW STATUS LIKE 'Innodb_buffer_pool_dump_status'\G
***** 1. row *****
Variable_name: Innodb_buffer_pool_dump_status
Value: Buffer pool(s) dump completed at 150202 16:38:58
```

5. Load the buffer pool by enabling `innodb_buffer_pool_load_now`:

```
mysql> SET GLOBAL innodb_buffer_pool_load_now=ON;
```

6. Check the current status of the buffer pool load operation by querying the Performance Schema `events_stages_current` table. The `WORK_COMPLETED` column shows the number of buffer pool pages loaded. The `WORK_ESTIMATED` column provides an estimate of the remaining work, in pages.

```
mysql> SELECT EVENT_NAME, WORK_COMPLETED, WORK_ESTIMATED
      FROM performance_schema.events_stages_current;
+-----+-----+-----+
```

EVENT_NAME	WORK_COMPLETED	WORK_ESTIMATED
stage/innodb/buffer pool load	5353	7167

The `events_stages_current` table returns an empty set if the buffer pool load operation has completed. In this case, you can check the `events_stages_history` table to view data for the completed event. For example:

EVENT_NAME	WORK_COMPLETED	WORK_ESTIMATED
stage/innodb/buffer pool load	7167	7167



Note

You can also monitor buffer pool load progress using Performance Schema when loading the buffer pool at startup using `innodb_buffer_pool_load_at_startup`. In this case, the `stage/innodb/buffer pool load` instrument and related consumers must be enabled at startup. For more information, see [Section 27.3, “Performance Schema Startup Configuration”](#).

15.8.3.7 Excluding Buffer Pool Pages from Core Files

A core file records the status and memory image of a running process. Because the buffer pool resides in main memory, and the memory image of a running process is dumped to the core file, systems with large buffer pools can produce large core files when the `mysqld` process dies.

Large core files can be problematic for a number of reasons including the time it takes to write them, the amount of disk space they consume, and the challenges associated with transferring large files.

To reduce core file size, you can disable the `innodb_buffer_pool_in_core_file` variable to omit buffer pool pages from core dumps. The `innodb_buffer_pool_in_core_file` variable was introduced in MySQL 8.0.14 and is enabled by default.

Excluding buffer pool pages may also be desirable from a security perspective if you have concerns about dumping database pages to core files that may be shared inside or outside of your organization for debugging purposes.



Note

Access to the data present in buffer pool pages at the time the `mysqld` process died may be beneficial in some debugging scenarios. If in doubt whether to include or exclude buffer pool pages, consult MySQL Support.

Disabling `innodb_buffer_pool_in_core_file` takes effect only if the `core_file` variable is enabled and the operating system supports the `MADV_DONTDUMP` non-POSIX extension to the `madvise()` system call, which is supported in Linux 3.4 and later. The `MADV_DONTDUMP` extension causes pages in a specified range to be excluded from core dumps.

Assuming the operating system supports the `MADV_DONTDUMP` extension, start the server with the `--core-file` and `--innodb-buffer-pool-in-core-file=OFF` options to generate core files without buffer pool pages.

```
$> mysqld --core-file --innodb-buffer-pool-in-core-file=OFF
```

The `core_file` variable is read only and disabled by default. It is enabled by specifying the `--core-file` option at startup. The `innodb_buffer_pool_in_core_file` variable is dynamic. It can be specified at startup or configured at runtime using a `SET` statement.

```
mysql> SET GLOBAL innodb_buffer_pool_in_core_file=OFF;
```

If the `innodb_buffer_pool_in_core_file` variable is disabled but `MADV_DONTDUMP` is not supported by the operating system, or an `madvise()` failure occurs, a warning is written to the MySQL server error log and the `core_file` variable is disabled to prevent writing core files that unintentionally include buffer pool pages. If the read-only `core_file` variable becomes disabled, the server must be restarted to enable it again.

The following table shows configuration and `MADV_DONTDUMP` support scenarios that determine whether core files are generated and whether they include buffer pool pages.

Table 15.4 Core File Configuration Scenarios

<code>core_file</code> variable	<code>innodb_buffer_pool</code> variable	<code>madvise()</code> file MADV_DONTDUMP Support	Outcome
OFF (default)	Not relevant to outcome	Not relevant to outcome	Core file is not generated
ON	ON (default)	Not relevant to outcome	Core file is generated with buffer pool pages
ON	OFF	Yes	Core file is generated without buffer pool pages
ON	OFF	No	Core file is not generated, <code>core_file</code> is disabled, and a warning is written to the server error log

The reduction in core file size achieved by disabling the `innodb_buffer_pool_in_core_file` variable depends on the size of the buffer pool, but it is also affected by the `InnoDB` page size. A smaller page size means more pages are required for the same amount of data, and more pages means more page metadata. The following table provides size reduction examples that you might see for a 1GB buffer pool with different pages sizes.

Table 15.5 Core File Size with Buffer Pool Pages Included and Excluded

<code>innodb_page_size</code> Setting	Buffer Pool Pages Included (<code>innodb_buffer_pool_in_core_file=ON</code>)	Buffer Pool Pages Excluded (<code>innodb_buffer_pool_in_core_file=OFF</code>)
4KB	2.1GB	0.9GB
64KB	1.7GB	0.7GB

15.8.4 Configuring Thread Concurrency for InnoDB

`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.

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.

You can set the configuration option `innodb_adaptive_max_sleep_delay` to the highest value you would allow for `innodb_thread_sleep_delay`, and InnoDB automatically adjusts `innodb_thread_sleep_delay` up or down depending on the current thread-scheduling activity. This dynamic adjustment helps the thread scheduling mechanism to work smoothly during times when the system is lightly loaded and when it is operating near full capacity.

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. The default value of `innodb_thread_concurrency` is 0, so that by default there is no limit on the number of concurrently executing threads.

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 (when `innodb_thread_concurrency` is > 0), 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 a specified number of “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 to perform row operations. If the tickets run out, the thread is evicted, and `innodb_thread_concurrency` is observed again which may place the thread back into the first-in/first-out queue of waiting threads. When the thread is once again entitled to enter InnoDB, tickets are assigned again. The number of tickets assigned is specified by the global option `innodb_concurrency_tickets`, which is 5000 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_adaptive_hash_index`.

For general performance information about MySQL thread handling, see [Section 5.1.12.1, “Connection Interfaces”](#).

15.8.5 Configuring the Number of 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 `innodb_read_io_threads` and `innodb_write_io_threads` configuration parameters. 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 values for these parameters in the MySQL option file (`my.cnf` or `my.ini`); you cannot change values dynamically. The default value for these parameters is 4 and permissible values range from 1–64.

The purpose of these configuration options 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 read-ahead requests. InnoDB tries to balance the load of incoming requests in such way that most 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. If you have a high

end I/O subsystem and you see more than $64 \times \text{innodb_read_io_threads}$ pending read requests in `SHOW ENGINE INNODB STATUS` output, you might improve performance by increasing the value of `innodb_read_io_threads`.

On Linux systems, `InnoDB` uses the asynchronous I/O subsystem by default to perform read-ahead and write requests for data file pages, which changes the way that `InnoDB` background threads service these types of I/O requests. For more information, see [Section 15.8.6, “Using Asynchronous I/O on Linux”](#).

For more information about `InnoDB` I/O performance, see [Section 8.5.8, “Optimizing InnoDB Disk I/O”](#).

15.8.6 Using Asynchronous I/O on Linux

`InnoDB` uses the asynchronous I/O subsystem (native AIO) on Linux to perform read-ahead and write requests for data file pages. This behavior is controlled by the `innodb_use_native_aio` configuration option, which applies to Linux systems only and is enabled by default. On other Unix-like systems, `InnoDB` uses synchronous I/O only. Historically, `InnoDB` only used asynchronous I/O on Windows systems. Using the asynchronous I/O subsystem on Linux requires the `libaio` library.

With synchronous I/O, query threads queue I/O requests, and `InnoDB` background threads retrieve the queued requests one at a time, issuing a synchronous I/O call for each. When an I/O request is completed and the I/O call returns, the `InnoDB` background thread that is handling the request calls an I/O completion routine and returns to process the next request. The number of requests that can be processed in parallel is `n`, where `n` is the number of `InnoDB` background threads. The number of `InnoDB` background threads is controlled by `innodb_read_io_threads` and `innodb_write_io_threads`. See [Section 15.8.5, “Configuring the Number of Background InnoDB I/O Threads”](#).

With native AIO, query threads dispatch I/O requests directly to the operating system, thereby removing the limit imposed by the number of background threads. `InnoDB` background threads wait for I/O events to signal completed requests. When a request is completed, a background thread calls an I/O completion routine and resumes waiting for I/O events.

The advantage of native AIO is scalability for heavily I/O-bound systems that typically show many pending reads/writes in `SHOW ENGINE INNODB STATUS\G` output. The increase in parallel processing when using native AIO means that the type of I/O scheduler or properties of the disk array controller have a greater influence on I/O performance.

A potential disadvantage of native AIO for heavily I/O-bound systems is lack of control over the number of I/O write requests dispatched to the operating system at once. Too many I/O write requests dispatched to the operating system for parallel processing could, in some cases, result in I/O read starvation, depending on the amount of I/O activity and system capabilities.

If a problem with the asynchronous I/O subsystem in the OS prevents `InnoDB` from starting, you can start the server with `innodb_use_native_aio=0`. This option may also be disabled automatically during startup if `InnoDB` detects a potential problem such as a combination of `tmpdir` location, `tmpfs` file system, and Linux kernel that does not support asynchronous I/O on `tmpfs`.

15.8.7 Configuring InnoDB I/O Capacity

The `InnoDB` master thread and other threads perform various tasks in the background, most of which are I/O related, such as flushing dirty pages from the buffer pool and writing changes from the change buffer to the appropriate secondary indexes. `InnoDB` attempts to perform these tasks in a way that does not adversely affect the normal working of the server. It tries to estimate the available I/O bandwidth and tune its activities to take advantage of available capacity.

The `innodb_io_capacity` variable defines the overall I/O capacity available to `InnoDB`. It should be set to approximately the number of I/O operations that the system can perform per second (IOPS). When `innodb_io_capacity` is set, `InnoDB` estimates the I/O bandwidth available for background tasks based on the set value.

You can set `innodb_io_capacity` to a value of 100 or greater. The default value is 200. Typically, values around 100 are appropriate for consumer-level storage devices, such as hard drives up to 7200 RPMs. Faster hard drives, RAID configurations, and solid state drives (SSDs) benefit from higher values.

Ideally, keep the setting as low as practical, but not so low that background activities fall behind. If the value is too high, data is removed from the buffer pool and change buffer too quickly for caching to provide a significant benefit. For busy systems capable of higher I/O rates, you can set a higher value to help the server handle the background maintenance work associated with a high rate of row changes. Generally, you can increase the value as a function of the number of drives used for InnoDB I/O. For example, you can increase the value on systems that use multiple disks or SSDs.

The default setting of 200 is generally sufficient for a lower-end SSD. For a higher-end, bus-attached SSD, consider a higher setting such as 1000, for example. For systems with individual 5400 RPM or 7200 RPM drives, you might lower the value to 100, which represents an estimated proportion of the I/O operations per second (IOPS) available to older-generation disk drives that can perform about 100 IOPS.

Although you can specify a high value such as a million, in practice such large values have little benefit. Generally, a value higher than 20000 is not recommended unless you are certain that lower values are insufficient for your workload.

Consider write workload when tuning `innodb_io_capacity`. Systems with large write workloads are likely to benefit from a higher setting. A lower setting may be sufficient for systems with a small write workload.

The `innodb_io_capacity` setting is not a per buffer pool instance setting. Available I/O capacity is distributed equally among buffer pool instances for flushing activities.

You can set the `innodb_io_capacity` value in the MySQL option file (`my.cnf` or `my.ini`) or modify it at runtime using a `SET GLOBAL` statement, which requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#).

Ignoring I/O Capacity at Checkpoints

The `innodb_flush_sync` variable, which is enabled by default, causes the `innodb_io_capacity` setting to be ignored during bursts of I/O activity that occur at `checkpoints`. To adhere to the I/O rate defined by the `innodb_io_capacity` setting, disable `innodb_flush_sync`.

You can set the `innodb_flush_sync` value in the MySQL option file (`my.cnf` or `my.ini`) or modify it at runtime using a `SET GLOBAL` statement, which requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#).

Configuring an I/O Capacity Maximum

If flushing activity falls behind, InnoDB can flush more aggressively, at a higher rate of I/O operations per second (IOPS) than defined by the `innodb_io_capacity` variable. The `innodb_io_capacity_max` variable defines a maximum number of IOPS performed by InnoDB background tasks in such situations.

If you specify an `innodb_io_capacity` setting at startup but do not specify a value for `innodb_io_capacity_max`, `innodb_io_capacity_max` defaults to twice the value of `innodb_io_capacity` or 2000, whichever value is greater.

When configuring `innodb_io_capacity_max`, twice the `innodb_io_capacity` is often a good starting point. The default value of 2000 is intended for workloads that use an SSD or more than one regular disk drive. A setting of 2000 is likely too high for workloads that do not use SSDs or multiple disk drives, and could allow too much flushing. For a single regular disk drive, a setting between 200 and 400 is recommended. For a high-end, bus-attached SSD, consider a higher setting such as 2500. As with the `innodb_io_capacity` setting, keep the setting as low as practical, but not so low that InnoDB cannot sufficiently extend rate of IOPS beyond the `innodb_io_capacity` setting.

Consider write workload when tuning `innodb_io_capacity_max`. Systems with large write workloads may benefit from a higher setting. A lower setting may be sufficient for systems with a small write workload.

`innodb_io_capacity_max` cannot be set to a value lower than the `innodb_io_capacity` value.

Setting `innodb_io_capacity_max` to `DEFAULT` using a `SET` statement (`SET GLOBAL innodb_io_capacity_max=DEFAULT`) sets `innodb_io_capacity_max` to the maximum value.

The `innodb_io_capacity_max` limit applies to all buffer pool instances. It is not a per buffer pool instance setting.

15.8.8 Configuring Spin Lock Polling

`InnoDB mutexes` and `rw-locks` are typically reserved for short intervals. 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 period of time before it sleeps. If the mutex or rw-lock becomes available during this period, the thread can continue immediately, in the same time slice. However, too-frequent polling of a shared object such as a mutex or rw-lock by multiple threads can cause “cache ping pong”, which results in processors invalidating portions of each other’s cache. `InnoDB` minimizes this issue by forcing a random delay between polls to desynchronize polling activity. The random delay is implemented as a spin-wait loop.

The duration of a spin-wait loop is determined by the number of PAUSE instructions that occur in the loop. That number is generated by randomly selecting an integer ranging from 0 up to but not including the `innodb_spin_wait_delay` value, and multiplying that value by 50. (The multiplier value, 50, is hardcoded before MySQL 8.0.16, and configurable thereafter.) For example, an integer is randomly selected from the following range for an `innodb_spin_wait_delay` setting of 6:

```
{0,1,2,3,4,5}
```

The selected integer is multiplied by 50, resulting in one of six possible PAUSE instruction values:

```
{0,50,100,150,200,250}
```

For that set of values, 250 is the maximum number of PAUSE instructions that can occur in a spin-wait loop. An `innodb_spin_wait_delay` setting of 5 results in a set of five possible values `{0,50,100,150,200}`, where 200 is the maximum number of PAUSE instructions, and so on. In this way, the `innodb_spin_wait_delay` setting controls the maximum delay between spin lock polls.

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.

In the 100MHz Pentium era, an `innodb_spin_wait_delay` unit was calibrated to be equivalent to one microsecond. That time equivalence did not hold, but PAUSE instruction duration remained fairly constant in terms of processor cycles relative to other CPU instructions until the introduction of the Skylake generation of processors, which have a comparatively longer PAUSE instruction. The `innodb_spin_wait_pause_multiplier` variable was introduced in MySQL 8.0.16 to provide a way to account for differences in PAUSE instruction duration.

The `innodb_spin_wait_pause_multiplier` variable controls the size of PAUSE instruction values. For example, assuming an `innodb_spin_wait_delay` setting of 6, decreasing the `innodb_spin_wait_pause_multiplier` value from 50 (the default and previously hardcoded value) to 5 generates a set of smaller PAUSE instruction values:

```
{0,5,10,15,20,25}
```

The ability to increase or decrease PAUSE instruction values permits fine tuning `InnoDB` for different processor architectures. Smaller PAUSE instruction values would be appropriate for processor architectures with a comparatively longer PAUSE instruction, for example.

The `innodb_spin_wait_delay` and `innodb_spin_wait_pause_multiplier` variables are dynamic. They can be specified in a MySQL option file or modified at runtime using a `SET GLOBAL` statement. Modifying the variables at runtime requires privileges sufficient to set global system variables. See [Section 5.1.9.1, “System Variable Privileges”](#).

15.8.9 Purge Configuration

InnoDB does not physically remove a row from the database immediately when you delete it with an SQL statement. A row and its index records are only physically removed when InnoDB discards the undo log record written for the deletion. This removal operation, which only occurs after the row is no longer required for multi-version concurrency control (MVCC) or rollback, is called a purge.

Purge runs on a periodic schedule. It parses and processes undo log pages from the history list, which is a list of undo log pages for committed transactions that is maintained by the InnoDB transaction system. Purge frees the undo log pages from the history list after processing them.

Configuring Purge Threads

Purge operations are performed in the background by one or more purge threads. The number of purge threads is controlled by the `innodb_purge_threads` variable. The default value is 4.

If DML action is concentrated on a single table, purge operations for the table are performed by a single purge thread, which can result in slowed purge operations, increased purge lag, and increased tablespace file size if the DML operations involve large object values. From MySQL 8.0.26, if the `innodb_max_purge_lag` setting is exceeded, purge work is automatically redistributed among available purge threads. Too many active purge threads in this scenario can cause contention with user threads, so manage the `innodb_purge_threads` setting accordingly. The `innodb_max_purge_lag` variable is set to 0 by default, which means that there is no maximum purge lag by default.

If DML action is concentrated on few tables, keep the `innodb_purge_threads` setting low so that the threads do not contend with each other for access to the busy tables. If DML operations are spread across many tables, consider a higher `innodb_purge_threads` setting. The maximum number of purge threads is 32.

The `innodb_purge_threads` setting is the maximum number of purge threads permitted. The purge system automatically adjusts the number of purge threads that are used.

Configuring Purge Batch Size

The `innodb_purge_batch_size` variable defines the number of undo log pages that purge parses and processes in one batch from the history list. The default value is 300. In a multithreaded purge configuration, the coordinator purge thread divides `innodb_purge_batch_size` by `innodb_purge_threads` and assigns that number of pages to each purge thread.

The purge system also frees the undo log pages that are no longer required. It does so every 128 iterations through the undo logs. In addition to defining the number of undo log pages parsed and processed in a batch, the `innodb_purge_batch_size` variable defines the number of undo log pages that purge frees every 128 iterations through the undo logs.

The `innodb_purge_batch_size` variable is intended for advanced performance tuning and experimentation. Most users need not change `innodb_purge_batch_size` from its default value.

Configuring the Maximum Purge Lag

The `innodb_max_purge_lag` variable defines the desired maximum purge lag. When the purge lag exceeds the `innodb_max_purge_lag` threshold, a delay is imposed on `INSERT`, `UPDATE`, and `DELETE` operations to allow time for purge operations to catch up. The default value is 0, which means there is no maximum purge lag and no delay.

The [InnoDB](#) transaction system maintains a list of transactions that have index records delete-marked by [UPDATE](#) or [DELETE](#) operations. The length of the list is the purge lag. Prior to MySQL 8.0.14, the purge lag delay is calculated by the following formula, which results in a minimum delay of 5000 microseconds:

```
(purge_lag/innodb_max_purge_lag - 0.5) * 10000
```

As of MySQL 8.0.14, the purge lag delay is calculated by the following revised formula, which reduces the minimum delay to 5 microseconds. A delay of 5 microseconds is more appropriate for modern systems.

```
(purge_lag/innodb_max_purge_lag - 0.9995) * 10000
```

The delay is calculated at the beginning of a purge batch.

A typical [innodb_max_purge_lag](#) setting for a problematic workload might be 1000000 (1 million), assuming that transactions are small, only 100 bytes in size, and it is permissible to have 100MB of unpurged table rows.

The purge lag is presented as the [History list length](#) value in the [TRANSACTIONS](#) section of [SHOW ENGINE INNODB STATUS](#) output.

```
mysql> SHOW ENGINE INNODB STATUS;
...
-----
TRANSACTIONS
-----
Trx id counter 0 290328385
Purge done for trx's n:o < 0 290315608 undo n:o < 0 17
History list length 20
```

The [History list length](#) is typically a low value, usually less than a few thousand, but a write-heavy workload or long running transactions can cause it to increase, even for transactions that are read only. The reason that a long running transaction can cause the [History list length](#) to increase is that under a consistent read transaction isolation level such as [REPEATABLE READ](#), a transaction must return the same result as when the read view for that transaction was created. Consequently, the [InnoDB](#) multi-version concurrency control (MVCC) system must keep a copy of the data in the undo log until all transactions that depend on that data have completed. The following are examples of long running transactions that could cause the [History list length](#) to increase:

- A [mysqldump](#) operation that uses the [--single-transaction](#) option while there is a significant amount of concurrent DML.
- Running a [SELECT](#) query after disabling [autocommit](#), and forgetting to issue an explicit [COMMIT](#) or [ROLLBACK](#).

To prevent excessive delays in extreme situations where the purge lag becomes huge, you can limit the delay by setting the [innodb_max_purge_lag_delay](#) variable. The [innodb_max_purge_lag_delay](#) variable specifies the maximum delay in microseconds for the delay imposed when the [innodb_max_purge_lag](#) threshold is exceeded. The specified [innodb_max_purge_lag_delay](#) value is an upper limit on the delay period calculated by the [innodb_max_purge_lag](#) formula.

Purge and Undo Tablespace Truncation

The purge system is also responsible for truncating undo tablespaces. You can configure the [innodb_purge_rseg_truncate_frequency](#) variable to control the frequency with which the purge system looks for undo tablespaces to truncate. For more information, see [Truncating Undo Tablespace](#)s.

15.8.10 Configuring Optimizer Statistics for InnoDB

This section describes how to configure persistent and non-persistent optimizer statistics for [InnoDB](#) tables.

Persistent optimizer statistics are persisted across server restarts, allowing for greater [plan stability](#) and more consistent query performance. Persistent optimizer statistics also provide control and flexibility with these additional benefits:

- You can use the `innodb_stats_auto_recalc` configuration option to control whether statistics are updated automatically after substantial changes to a table.
- You can use the `STATS_PERSISTENT`, `STATS_AUTO_RECALC`, and `STATS_SAMPLE_PAGES` clauses with `CREATE TABLE` and `ALTER TABLE` statements to configure optimizer statistics for individual tables.
- You can query optimizer statistics data in the `mysql.innodb_table_stats` and `mysql.innodb_index_stats` tables.
- You can view the `last_update` column of the `mysql.innodb_table_stats` and `mysql.innodb_index_stats` tables to see when statistics were last updated.
- You can manually modify the `mysql.innodb_table_stats` and `mysql.innodb_index_stats` tables to force a specific query optimization plan or to test alternative plans without modifying the database.

The persistent optimizer statistics feature is enabled by default (`innodb_stats_persistent=ON`).

Non-persistent optimizer statistics are cleared on each server restart and after some other operations, and recomputed on the next table access. As a result, different estimates could be produced when recomputing statistics, leading to different choices in execution plans and variations in query performance.

This section also provides information about estimating `ANALYZE TABLE` complexity, which may be useful when attempting to achieve a balance between accurate statistics and `ANALYZE TABLE` execution time.

15.8.10.1 Configuring Persistent Optimizer Statistics Parameters

The persistent optimizer statistics feature improves [plan stability](#) by storing statistics to disk and making them persistent across server restarts so that the [optimizer](#) is more likely to make consistent choices each time for a given query.

Optimizer statistics are persisted to disk when `innodb_stats_persistent=ON` or when individual tables are defined with `STATS_PERSISTENT=1`. `innodb_stats_persistent` is enabled by default.

Formerly, optimizer statistics were cleared when restarting the server and after some other types of operations, and recomputed on the next table access. Consequently, different estimates could be produced when recalculating statistics leading to different choices in query execution plans and variation in query performance.

Persistent statistics are stored in the `mysql.innodb_table_stats` and `mysql.innodb_index_stats` tables. See [InnoDB Persistent Statistics Tables](#).

If you prefer not to persist optimizer statistics to disk, see [Section 15.8.10.2, “Configuring Non-Persistent Optimizer Statistics Parameters”](#)

Configuring Automatic Statistics Calculation for Persistent Optimizer Statistics

The `innodb_stats_auto_recalc` variable, which is enabled by default, controls whether statistics are calculated automatically when a table undergoes changes to more than 10% of its rows. You can also configure automatic statistics recalculation for individual tables by specifying the `STATS_AUTO_RECALC` clause when creating or altering a table.

Because of the asynchronous nature of automatic statistics recalculation, which occurs in the background, statistics may not be recalculated instantly after running a DML operation that affects

more than 10% of a table, even when `innodb_stats_auto_recalc` is enabled. Statistics recalculation can be delayed by few seconds in some cases. If up-to-date statistics are required immediately, run `ANALYZE TABLE` to initiate a synchronous (foreground) recalculation of statistics.

If `innodb_stats_auto_recalc` is disabled, you can ensure the accuracy of optimizer statistics by executing the `ANALYZE TABLE` statement after making substantial changes to indexed columns. You might also consider adding `ANALYZE TABLE` to setup scripts that you run after loading data, and running `ANALYZE TABLE` on a schedule at times of low activity.

When an index is added to an existing table, or when a column is added or dropped, index statistics are calculated and added to the `innodb_index_stats` table regardless of the value of `innodb_stats_auto_recalc`.

Configuring Optimizer Statistics Parameters for Individual Tables

`innodb_stats_persistent`, `innodb_stats_auto_recalc`, and `innodb_stats_persistent_sample_pages` are global variables. To override these system-wide settings and configure optimizer statistics parameters for individual tables, you can define `STATS_PERSISTENT`, `STATS_AUTO_RECALC`, and `STATS_SAMPLE_PAGES` clauses in `CREATE TABLE` or `ALTER TABLE` statements.

- `STATS_PERSISTENT` specifies whether to enable `persistent statistics` for an `InnoDB` table. The value `DEFAULT` causes the persistent statistics setting for the table to be determined by the `innodb_stats_persistent` setting. A value of `1` enables persistent statistics for the table, while a value of `0` disables the feature. After enabling persistent statistics for an individual table, use `ANALYZE TABLE` to calculate statistics after table data is loaded.
- `STATS_AUTO_RECALC` specifies whether to automatically recalculate `persistent statistics`. The value `DEFAULT` causes the persistent statistics setting for the table to be determined by the `innodb_stats_auto_recalc` setting. A value of `1` causes statistics to be recalculated when 10% of table data has changed. A value of `0` prevents automatic recalculation for the table. When using a value of `0`, use `ANALYZE TABLE` to recalculate statistics after making substantial changes to the table.
- `STATS_SAMPLE_PAGES` specifies the number of index pages to sample when cardinality and other statistics are calculated for an indexed column, by an `ANALYZE TABLE` operation, for example.

All three clauses are specified in the following `CREATE TABLE` example:

```
CREATE TABLE `t1` (
  `id` int(8) NOT NULL auto_increment,
  `data` varchar(255),
  `date` datetime,
  PRIMARY KEY (`id`),
  INDEX `DATE_IX` (`date`)
) ENGINE=InnoDB,
  STATS_PERSISTENT=1,
  STATS_AUTO_RECALC=1,
  STATS_SAMPLE_PAGES=25;
```

Configuring the Number of Sampled Pages for InnoDB Optimizer Statistics

The optimizer uses estimated `statistics` about key distributions to choose the indexes for an execution plan, based on the relative `selectivity` of the index. Operations such as `ANALYZE TABLE` cause `InnoDB` to sample random pages from each index on a table to estimate the `cardinality` of the index. This sampling technique is known as a `random dive`.

The `innodb_stats_persistent_sample_pages` controls the number of sampled pages. You can adjust the setting at runtime to manage the quality of statistics estimates used by the optimizer. The default value is 20. Consider modifying the setting when encountering the following issues:

1. *Statistics are not accurate enough and the optimizer chooses suboptimal plans*, as shown in `EXPLAIN` output. You can check the accuracy of statistics by comparing the actual cardinality of an

index (determined by running `SELECT DISTINCT` on the index columns) with the estimates in the `mysql.innodb_index_stats` table.

If it is determined that statistics are not accurate enough, the value of `innodb_stats_persistent_sample_pages` should be increased until the statistics estimates are sufficiently accurate. Increasing `innodb_stats_persistent_sample_pages` too much, however, could cause `ANALYZE TABLE` to run slowly.

2. `ANALYZE TABLE` is too slow. In this case `innodb_stats_persistent_sample_pages` should be decreased until `ANALYZE TABLE` execution time is acceptable. Decreasing the value too much, however, could lead to the first problem of inaccurate statistics and suboptimal query execution plans.

If a balance cannot be achieved between accurate statistics and `ANALYZE TABLE` execution time, consider decreasing the number of indexed columns in the table or limiting the number of partitions to reduce `ANALYZE TABLE` complexity. The number of columns in the table's primary key is also important to consider, as primary key columns are appended to each nonunique index.

For related information, see [Section 15.8.10.3, “Estimating ANALYZE TABLE Complexity for InnoDB Tables”](#).

Including Delete-marked Records in Persistent Statistics Calculations

By default, `InnoDB` reads uncommitted data when calculating statistics. In the case of an uncommitted transaction that deletes rows from a table, delete-marked records are excluded when calculating row estimates and index statistics, which can lead to non-optimal execution plans for other transactions that are operating on the table concurrently using a transaction isolation level other than `READ UNCOMMITTED`. To avoid this scenario, `innodb_stats_include_delete_marked` can be enabled to ensure that delete-marked records are included when calculating persistent optimizer statistics.

When `innodb_stats_include_delete_marked` is enabled, `ANALYZE TABLE` considers delete-marked records when recalculating statistics.

`innodb_stats_include_delete_marked` is a global setting that affects all `InnoDB` tables, and it is only applicable to persistent optimizer statistics.

InnoDB Persistent Statistics Tables

The persistent statistics feature relies on the internally managed tables in the `mysql` database, named `innodb_table_stats` and `innodb_index_stats`. These tables are set up automatically in all install, upgrade, and build-from-source procedures.

Table 15.6 Columns of innodb_table_stats

Column name	Description
<code>database_name</code>	Database name
<code>table_name</code>	Table name, partition name, or subpartition name
<code>last_update</code>	A timestamp indicating the last time that <code>InnoDB</code> updated this row
<code>n_rows</code>	The number of rows in the table
<code>clustered_index_size</code>	The size of the primary index, in pages
<code>sum_of_other_index_sizes</code>	The total size of other (non-primary) indexes, in pages

Table 15.7 Columns of innodb_index_stats

Column name	Description
<code>database_name</code>	Database name

Column name	Description
table_name	Table name, partition name, or subpartition name
index_name	Index name
last_update	A timestamp indicating the last time the row was updated
stat_name	The name of the statistic, whose value is reported in the stat_value column
stat_value	The value of the statistic that is named in stat_name column
sample_size	The number of pages sampled for the estimate provided in the stat_value column
stat_description	Description of the statistic that is named in the stat_name column

The `innodb_table_stats` and `innodb_index_stats` tables include a `last_update` column that shows when index statistics were last updated:

```
mysql> SELECT * FROM innodb_table_stats \G
***** 1. row *****
    database_name: sakila
        table_name: actor
        last_update: 2014-05-28 16:16:44
        n_rows: 200
    clustered_index_size: 1
sum_of_other_index_sizes: 1
...
mysql> SELECT * FROM innodb_index_stats \G
***** 1. row *****
    database_name: sakila
        table_name: actor
        index_name: PRIMARY
        last_update: 2014-05-28 16:16:44
        stat_name: n_diff_pfx01
        stat_value: 200
        sample_size: 1
...

```

The `innodb_table_stats` and `innodb_index_stats` tables can be updated manually, which makes it possible to force a specific query optimization plan or test alternative plans without modifying the database. If you manually update statistics, use the `FLUSH TABLE tbl_name` statement to load the updated statistics.

Persistent statistics are considered local information, because they relate to the server instance. The `innodb_table_stats` and `innodb_index_stats` tables are therefore not replicated when automatic statistics recalculation takes place. If you run `ANALYZE TABLE` to initiate a synchronous recalculation of statistics, the statement is replicated (unless you suppressed logging for it), and recalculation takes place on replicas.

InnoDB Persistent Statistics Tables Example

The `innodb_table_stats` table contains one row for each table. The following example demonstrates the type of data collected.

Table `t1` contains a primary index (columns `a`, `b`) secondary index (columns `c`, `d`), and unique index (columns `e`, `f`):

```
CREATE TABLE t1 (
  a INT, b INT, c INT, d INT, e INT, f INT,
  PRIMARY KEY (a, b), KEY i1 (c, d), UNIQUE KEY i2uniq (e, f)
) ENGINE=INNODB;
```

After inserting five rows of sample data, table `t1` appears as follows:

```
mysql> SELECT * FROM t1;
+---+---+---+---+---+---+
| a | b | c | d | e | f |
+---+---+---+---+---+---+
| 1 | 1 | 10 | 11 | 100 | 101 |
| 1 | 2 | 10 | 11 | 200 | 102 |
| 1 | 3 | 10 | 11 | 100 | 103 |
| 1 | 4 | 10 | 12 | 200 | 104 |
| 1 | 5 | 10 | 12 | 100 | 105 |
+---+---+---+---+---+---+
```

To immediately update statistics, run `ANALYZE TABLE` (if `innodb_stats_auto_recalc` is enabled, statistics are updated automatically within a few seconds assuming that the 10% threshold for changed table rows is reached):

```
mysql> ANALYZE TABLE t1;
+-----+-----+-----+
| Table | Op    | Msg_type | Msg_text |
+-----+-----+-----+
| test.t1 | analyze | status   | OK      |
+-----+-----+-----+
```

Table statistics for table `t1` show the last time InnoDB updated the table statistics (2014-03-14 14:36:34), the number of rows in the table (5), the clustered index size (1 page), and the combined size of the other indexes (2 pages).

```
mysql> SELECT * FROM mysql.innodb_table_stats WHERE table_name like 't1'\G
***** 1. row *****
database_name: test
table_name: t1
last_update: 2014-03-14 14:36:34
n_rows: 5
clustered_index_size: 1
sum_of_other_index_sizes: 2
```

The `innodb_index_stats` table contains multiple rows for each index. Each row in the `innodb_index_stats` table provides data related to a particular index statistic which is named in the `stat_name` column and described in the `stat_description` column. For example:

```
mysql> SELECT index_name, stat_name, stat_value, stat_description
      FROM mysql.innodb_index_stats WHERE table_name like 't1';
+-----+-----+-----+-----+
| index_name | stat_name | stat_value | stat_description |
+-----+-----+-----+-----+
| PRIMARY    | n_diff_pfx01 | 1          | a                |
| PRIMARY    | n_diff_pfx02 | 5          | a,b              |
| PRIMARY    | n_leaf_pages | 1          | Number of leaf pages in the index |
| PRIMARY    | size         | 1          | Number of pages in the index        |
| i1          | n_diff_pfx01 | 1          | c                |
| i1          | n_diff_pfx02 | 2          | c,d              |
| i1          | n_diff_pfx03 | 2          | c,d,a            |
| i1          | n_diff_pfx04 | 5          | c,d,a,b          |
| i1          | n_leaf_pages | 1          | Number of leaf pages in the index |
| i1          | size         | 1          | Number of pages in the index        |
| i2uniq     | n_diff_pfx01 | 2          | e                |
| i2uniq     | n_diff_pfx02 | 5          | e,f              |
| i2uniq     | n_leaf_pages | 1          | Number of leaf pages in the index |
| i2uniq     | size         | 1          | Number of pages in the index        |
+-----+-----+-----+-----+
```

The `stat_name` column shows the following types of statistics:

- `size`: Where `stat_name=size`, the `stat_value` column displays the total number of pages in the index.
- `n_leaf_pages`: Where `stat_name=n_leaf_pages`, the `stat_value` column displays the number of leaf pages in the index.

- `n_diff_pfxNN`: Where `stat_name=n_diff_pfx01`, the `stat_value` column displays the number of distinct values in the first column of the index. Where `stat_name=n_diff_pfx02`, the `stat_value` column displays the number of distinct values in the first two columns of the index, and so on. Where `stat_name=n_diff_pfxNN`, the `stat_description` column shows a comma separated list of the index columns that are counted.

To further illustrate the `n_diff_pfxNN` statistic, which provides cardinality data, consider once again the `t1` table example that was introduced previously. As shown below, the `t1` table is created with a primary index (columns `a`, `b`), a secondary index (columns `c`, `d`), and a unique index (columns `e`, `f`):

```
CREATE TABLE t1 (
  a INT, b INT, c INT, d INT, e INT, f INT,
  PRIMARY KEY (a, b), KEY i1 (c, d), UNIQUE KEY i2uniq (e, f)
) ENGINE=INNODB;
```

After inserting five rows of sample data, table `t1` appears as follows:

```
mysql> SELECT * FROM t1;
+---+---+---+---+---+---+
| a | b | c   | d   | e   | f   |
+---+---+---+---+---+---+
| 1 | 1 | 10  | 11  | 100 | 101 |
| 1 | 2 | 10  | 11  | 200 | 102 |
| 1 | 3 | 10  | 11  | 100 | 103 |
| 1 | 4 | 10  | 12  | 200 | 104 |
| 1 | 5 | 10  | 12  | 100 | 105 |
+---+---+---+---+---+---+
```

When you query the `index_name`, `stat_name`, `stat_value`, and `stat_description`, where `stat_name LIKE 'n_diff%`', the following result set is returned:

```
mysql> SELECT index_name, stat_name, stat_value, stat_description
  FROM mysql.innodb_index_stats
 WHERE table_name like 't1' AND stat_name LIKE 'n_diff%';
+-----+-----+-----+-----+
| index_name | stat_name | stat_value | stat_description |
+-----+-----+-----+-----+
| PRIMARY    | n_diff_pfx01 |      1 | a
| PRIMARY    | n_diff_pfx02 |      5 | a,b
| i1         | n_diff_pfx01 |      1 | c
| i1         | n_diff_pfx02 |      2 | c,d
| i1         | n_diff_pfx03 |      2 | c,d,a
| i1         | n_diff_pfx04 |      5 | c,d,a,b
| i2uniq    | n_diff_pfx01 |      2 | e
| i2uniq    | n_diff_pfx02 |      5 | e,f
+-----+-----+-----+-----+
```

For the `PRIMARY` index, there are two `n_diff%` rows. The number of rows is equal to the number of columns in the index.



Note

For nonunique indexes, InnoDB appends the columns of the primary key.

- Where `index_name=PRIMARY` and `stat_name=n_diff_pfx01`, the `stat_value` is 1, which indicates that there is a single distinct value in the first column of the index (column `a`). The number of distinct values in column `a` is confirmed by viewing the data in column `a` in table `t1`, in which there is a single distinct value (1). The counted column (`a`) is shown in the `stat_description` column of the result set.
- Where `index_name=PRIMARY` and `stat_name=n_diff_pfx02`, the `stat_value` is 5, which indicates that there are five distinct values in the two columns of the index (`a,b`). The number of distinct values in columns `a` and `b` is confirmed by viewing the data in columns `a` and `b` in table `t1`, in which there are five distinct values: (1,1), (1,2), (1,3), (1,4) and (1,5). The counted columns (`a,b`) are shown in the `stat_description` column of the result set.

For the secondary index (`i1`), there are four `n_diff%` rows. Only two columns are defined for the secondary index (`c,d`) but there are four `n_diff%` rows for the secondary index because InnoDB suffixes all nonunique indexes with the primary key. As a result, there are four `n_diff%` rows instead of two to account for the both the secondary index columns (`c,d`) and the primary key columns (`a,b`).

- Where `index_name=i1` and `stat_name=n_diff_pfx01`, the `stat_value` is 1, which indicates that there is a single distinct value in the first column of the index (column `c`). The number of distinct values in column `c` is confirmed by viewing the data in column `c` in table `t1`, in which there is a single distinct value: (10). The counted column (`c`) is shown in the `stat_description` column of the result set.
- Where `index_name=i1` and `stat_name=n_diff_pfx02`, the `stat_value` is 2, which indicates that there are two distinct values in the first two columns of the index (`c,d`). The number of distinct values in columns `c` and `d` is confirmed by viewing the data in columns `c` and `d` in table `t1`, in which there are two distinct values: (10,11) and (10,12). The counted columns (`c,d`) are shown in the `stat_description` column of the result set.
- Where `index_name=i1` and `stat_name=n_diff_pfx03`, the `stat_value` is 2, which indicates that there are two distinct values in the first three columns of the index (`c,d,a`). The number of distinct values in columns `c,d`, and `a` is confirmed by viewing the data in column `c,d`, and `a` in table `t1`, in which there are two distinct values: (10,11,1) and (10,12,1). The counted columns (`c,d,a`) are shown in the `stat_description` column of the result set.
- Where `index_name=i1` and `stat_name=n_diff_pfx04`, the `stat_value` is 5, which indicates that there are five distinct values in the four columns of the index (`c,d,a,b`). The number of distinct values in columns `c,d,a`, and `b` is confirmed by viewing the data in columns `c,d,a`, and `b` in table `t1`, in which there are five distinct values: (10,11,1,1), (10,11,1,2), (10,11,1,3), (10,12,1,4), and (10,12,1,5). The counted columns (`c,d,a,b`) are shown in the `stat_description` column of the result set.

For the unique index (`i2uniq`), there are two `n_diff%` rows.

- Where `index_name=i2uniq` and `stat_name=n_diff_pfx01`, the `stat_value` is 2, which indicates that there are two distinct values in the first column of the index (column `e`). The number of distinct values in column `e` is confirmed by viewing the data in column `e` in table `t1`, in which there are two distinct values: (100) and (200). The counted column (`e`) is shown in the `stat_description` column of the result set.
- Where `index_name=i2uniq` and `stat_name=n_diff_pfx02`, the `stat_value` is 5, which indicates that there are five distinct values in the two columns of the index (`e,f`). The number of distinct values in columns `e` and `f` is confirmed by viewing the data in columns `e` and `f` in table `t1`, in which there are five distinct values: (100,101), (200,102), (100,103), (200,104), and (100,105). The counted columns (`e,f`) are shown in the `stat_description` column of the result set.

Retrieving Index Size Using the `innodb_index_stats` Table

You can retrieve the index size for tables, partitions, or subpartitions can using the `innodb_index_stats` table. In the following example, index sizes are retrieved for table `t1`. For a definition of table `t1` and corresponding index statistics, see [InnoDB Persistent Statistics Tables Example](#).

```
mysql> SELECT SUM(stat_value) pages, index_name,
    SUM(stat_value)*@innodb_page_size size
    FROM mysql.innodb_index_stats WHERE table_name='t1'
    AND stat_name = 'size' GROUP BY index_name;
+-----+-----+-----+
| pages | index_name | size   |
+-----+-----+-----+
|     1 | PRIMARY    | 16384 |
|     1 | i1          | 16384 |
|     1 | i2uniq      | 16384 |
+-----+-----+-----+
```