# Chapter 10 Optimization

## **Table of Contents**

	Optimization Overview	
10.2	Optimizing SQL Statements	
	10.2.1 Optimizing SELECT Statements	1634
	10.2.2 Optimizing Subqueries, Derived Tables, View References, and Common Table	
	Expressions	
	10.2.3 Optimizing INFORMATION_SCHEMA Queries	
	10.2.4 Optimizing Performance Schema Queries	1700
	10.2.5 Optimizing Data Change Statements	1702
	10.2.6 Optimizing Database Privileges	1703
	10.2.7 Other Optimization Tips	1703
10.3	Optimization and Indexes	1704
	10.3.1 How MySQL Uses Indexes	1704
	10.3.2 Primary Key Optimization	1705
	10.3.3 SPATIAL Index Optimization	1705
	10.3.4 Foreign Key Optimization	1706
	10.3.5 Column Indexes	1706
	10.3.6 Multiple-Column Indexes	1707
	10.3.7 Verifying Index Usage	1709
	10.3.8 InnoDB and MyISAM Index Statistics Collection	1709
	10.3.9 Comparison of B-Tree and Hash Indexes	
	10.3.10 Use of Index Extensions	
	10.3.11 Optimizer Use of Generated Column Indexes	
	10.3.12 Invisible Indexes	
	10.3.13 Descending Indexes	1717
	10.3.14 Indexed Lookups from TIMESTAMP Columns	
10.4	Optimizing Database Structure	
	10.4.1 Optimizing Data Size	
	10.4.2 Optimizing MySQL Data Types	
	10.4.3 Optimizing for Many Tables	
	10.4.4 Internal Temporary Table Use in MySQL	
	10.4.5 Limits on Number of Databases and Tables	
	10.4.6 Limits on Table Size	
	10.4.7 Limits on Table Column Count and Row Size	
10.5	Optimizing for InnoDB Tables	
	10.5.1 Optimizing Storage Layout for InnoDB Tables	
	10.5.2 Optimizing InnoDB Transaction Management	
	10.5.3 Optimizing InnoDB Read-Only Transactions	
	10.5.4 Optimizing InnoDB Redo Logging	
	10.5.5 Bulk Data Loading for InnoDB Tables	
	10.5.6 Optimizing InnoDB Queries	
	10.5.7 Optimizing InnoDB DDL Operations	
	10.5.8 Optimizing InnoDB Disk I/O	
	10.5.9 Optimizing InnoDB Configuration Variables	
	10.5.10 Optimizing InnoDB for Systems with Many Tables	
10.6	Optimizing for MyISAM Tables	
	10.6.1 Optimizing MylSAM Queries	
	· · · · · · · · · · · · · · · · · · ·	1745
	10.6.3 Optimizing REPAIR TABLE Statements	
10 7	Optimizing for MEMORY Tables	
	Understanding the Query Execution Plan	
	10.8.1 Optimizing Queries with EXPLAIN	
	10.8.2 EXPLAIN Output Format	

10.8.3 Extended EXPLAIN Output Format	1762
10.8.4 Obtaining Execution Plan Information for a Named Connection	1764
10.8.5 Estimating Query Performance	1765
10.9 Controlling the Query Optimizer	1765
10.9.1 Controlling Query Plan Evaluation	1765
10.9.2 Switchable Optimizations	
10.9.3 Optimizer Hints	
10.9.4 Index Hints	
10.9.5 The Optimizer Cost Model	
10.9.6 Optimizer Statistics	
10.10 Buffering and Caching	
10.10.1 InnoDB Buffer Pool Optimization	
10.10.2 The MylSAM Key Cache	
10.10.3 Caching of Prepared Statements and Stored Programs	
10.10.3 Caching of Frepared Statements and Stored Frograms	
10.11.1 Internal Locking Methods	
10.11.2 Table Locking Issues	
10.11.3 Concurrent Inserts	
10.11.4 Metadata Locking	
10.11.5 External Locking	
10.12 Optimizing the MySQL Server	
10.12.1 Optimizing Disk I/O	
10.12.2 Using Symbolic Links	
10.12.3 Optimizing Memory Use	
10.13 Measuring Performance (Benchmarking)	
10.13.1 Measuring the Speed of Expressions and Functions	1825
10.13.2 Using Your Own Benchmarks	1825
10.13.3 Measuring Performance with performance_schema	1825
10.14 Examining Server Thread (Process) Information	1826
10.14.1 Accessing the Process List	
10.14.2 Thread Command Values	
10.14.3 General Thread States	
10.14.4 Replication Source Thread States	
10.14.5 Replication I/O (Receiver) Thread States	
10.14.6 Replication SQL Thread States	
10.14.7 Replication Connection Thread States	
10.14.8 NDB Cluster Thread States	
10.14.9 Event Scheduler Thread States	
10.14.9 Event Scrieduler Thread States	
· · · · · · · · · · · · · · · · · · ·	
10.15.1 Typical Usage	
10.15.2 System Variables Controlling Tracing	
10.15.3 Traceable Statements	
5 5 5	1843
	1844
	1844
10.15.7 Interaction with thedebug Option	
10.15.8 The optimizer_trace System Variable	
10.15.9 The end_markers_in_json System Variable	
10.15.10 Selecting Optimizer Features to Trace	
10.15.11 Trace General Structure	1845
10.15.12 Example	1845
10.15.13 Displaying Traces in Other Applications	1855
10.15.14 Preventing the Use of Optimizer Trace	
10.15.15 Testing Optimizer Trace	
	1856

This chapter explains how to optimize MySQL performance and provides examples. Optimization involves configuring, tuning, and measuring performance, at several levels. Depending on your job

role (developer, DBA, or a combination of both), you might optimize at the level of individual SQL statements, entire applications, a single database server, or multiple networked database servers. Sometimes you can be proactive and plan in advance for performance, while other times you might troubleshoot a configuration or code issue after a problem occurs. Optimizing CPU and memory usage can also improve scalability, allowing the database to handle more load without slowing down.

## 10.1 Optimization Overview

Database performance depends on several factors at the database level, such as tables, queries, and configuration settings. These software constructs result in CPU and I/O operations at the hardware level, which you must minimize and make as efficient as possible. As you work on database performance, you start by learning the high-level rules and guidelines for the software side, and measuring performance using wall-clock time. As you become an expert, you learn more about what happens internally, and start measuring things such as CPU cycles and I/O operations.

Typical users aim to get the best database performance out of their existing software and hardware configurations. Advanced users look for opportunities to improve the MySQL software itself, or develop their own storage engines and hardware appliances to expand the MySQL ecosystem.

- · Optimizing at the Database Level
- Optimizing at the Hardware Level
- · Balancing Portability and Performance

## Optimizing at the Database Level

The most important factor in making a database application fast is its basic design:

- Are the tables structured properly? In particular, do the columns have the right data types, and
  does each table have the appropriate columns for the type of work? For example, applications that
  perform frequent updates often have many tables with few columns, while applications that analyze
  large amounts of data often have few tables with many columns.
- Are the right indexes in place to make queries efficient?
- Are you using the appropriate storage engine for each table, and taking advantage of the strengths
  and features of each storage engine you use? In particular, the choice of a transactional storage
  engine such as InnoDB or a nontransactional one such as MyISAM can be very important for
  performance and scalability.



#### Note

InnoDB is the default storage engine for new tables. In practice, the advanced InnoDB performance features mean that InnoDB tables often outperform the simpler MyISAM tables, especially for a busy database.

- Does each table use an appropriate row format? This choice also depends on the storage engine
  used for the table. In particular, compressed tables use less disk space and so require less disk I/O
  to read and write the data. Compression is available for all kinds of workloads with InnoDB tables,
  and for read-only MyISAM tables.
- Does the application use an appropriate locking strategy? For example, by allowing shared access
  when possible so that database operations can run concurrently, and requesting exclusive access
  when appropriate so that critical operations get top priority. Again, the choice of storage engine is
  significant. The Innode storage engine handles most locking issues without involvement from you,
  allowing for better concurrency in the database and reducing the amount of experimentation and
  tuning for your code.
- Are all memory areas used for caching sized correctly? That is, large enough to hold frequently
  accessed data, but not so large that they overload physical memory and cause paging. The main
  memory areas to configure are the InnoDB buffer pool and the MyISAM key cache.

### **Optimizing at the Hardware Level**

Any database application eventually hits hardware limits as the database becomes more and more busy. A DBA must evaluate whether it is possible to tune the application or reconfigure the server to avoid these bottlenecks, or whether more hardware resources are required. System bottlenecks typically arise from these sources:

- Disk seeks. It takes time for the disk to find a piece of data. With modern disks, the mean time
  for this is usually lower than 10ms, so we can in theory do about 100 seeks a second. This time
  improves slowly with new disks and is very hard to optimize for a single table. The way to optimize
  seek time is to distribute the data onto more than one disk.
- Disk reading and writing. When the disk is at the correct position, we need to read or write the data.
   With modern disks, one disk delivers at least 10–20MB/s throughput. This is easier to optimize than seeks because you can read in parallel from multiple disks.
- CPU cycles. When the data is in main memory, we must process it to get our result. Having large
  tables compared to the amount of memory is the most common limiting factor. But with small tables,
  speed is usually not the problem.
- Memory bandwidth. When the CPU needs more data than can fit in the CPU cache, main memory bandwidth becomes a bottleneck. This is an uncommon bottleneck for most systems, but one to be aware of.

## **Balancing Portability and Performance**

To use performance-oriented SQL extensions in a portable MySQL program, you can wrap MySQL-specific keywords in a statement within /\*! \*/ comment delimiters. Other SQL servers ignore the commented keywords. For information about writing comments, see Section 11.7, "Comments".

## 10.2 Optimizing SQL Statements

The core logic of a database application is performed through SQL statements, whether issued directly through an interpreter or submitted behind the scenes through an API. The tuning guidelines in this section help to speed up all kinds of MySQL applications. The guidelines cover SQL operations that read and write data, the behind-the-scenes overhead for SQL operations in general, and operations used in specific scenarios such as database monitoring.

## 10.2.1 Optimizing SELECT Statements

Queries, in the form of SELECT statements, perform all the lookup operations in the database. Tuning these statements is a top priority, whether to achieve sub-second response times for dynamic web pages, or to chop hours off the time to generate huge overnight reports.

Besides SELECT statements, the tuning techniques for queries also apply to constructs such as CREATE TABLE...AS SELECT, INSERT INTO...SELECT, and WHERE clauses in DELETE statements. Those statements have additional performance considerations because they combine write operations with the read-oriented query operations.

NDB Cluster supports a join pushdown optimization whereby a qualifying join is sent in its entirety to NDB Cluster data nodes, where it can be distributed among them and executed in parallel. For more information about this optimization, see Conditions for NDB pushdown joins.

The main considerations for optimizing queries are:

• To make a slow SELECT ... WHERE query faster, the first thing to check is whether you can add an index. Set up indexes on columns used in the WHERE clause, to speed up evaluation, filtering, and the final retrieval of results. To avoid wasted disk space, construct a small set of indexes that speed up many related queries used in your application.

Indexes are especially important for queries that reference different tables, using features such as joins and foreign keys. You can use the EXPLAIN statement to determine which indexes are used for a SELECT. See Section 10.3.1, "How MySQL Uses Indexes" and Section 10.8.1, "Optimizing Queries with EXPLAIN".

- Isolate and tune any part of the query, such as a function call, that takes excessive time. Depending
  on how the query is structured, a function could be called once for every row in the result set, or even
  once for every row in the table, greatly magnifying any inefficiency.
- Minimize the number of full table scans in your queries, particularly for big tables.
- Keep table statistics up to date by using the ANALYZE TABLE statement periodically, so the
  optimizer has the information needed to construct an efficient execution plan.
- Learn the tuning techniques, indexing techniques, and configuration parameters that are specific to
  the storage engine for each table. Both InnoDB and MyISAM have sets of guidelines for enabling
  and sustaining high performance in queries. For details, see Section 10.5.6, "Optimizing InnoDB
  Queries" and Section 10.6.1, "Optimizing MyISAM Queries".
- You can optimize single-query transactions for InnoDB tables, using the technique in Section 10.5.3, "Optimizing InnoDB Read-Only Transactions".
- Avoid transforming the query in ways that make it hard to understand, especially if the optimizer does some of the same transformations automatically.
- If a performance issue is not easily solved by one of the basic guidelines, investigate the internal details of the specific query by reading the EXPLAIN plan and adjusting your indexes, WHERE clauses, join clauses, and so on. (When you reach a certain level of expertise, reading the EXPLAIN plan might be your first step for every query.)
- Adjust the size and properties of the memory areas that MySQL uses for caching. With efficient use
  of the InnoDB buffer pool, MyISAM key cache, and the MySQL query cache, repeated queries run
  faster because the results are retrieved from memory the second and subsequent times.
- Even for a query that runs fast using the cache memory areas, you might still optimize further so that
  they require less cache memory, making your application more scalable. Scalability means that your
  application can handle more simultaneous users, larger requests, and so on without experiencing a
  big drop in performance.
- Deal with locking issues, where the speed of your query might be affected by other sessions accessing the tables at the same time.

#### 10.2.1.1 WHERE Clause Optimization

This section discusses optimizations that can be made for processing where clauses. The examples use Select statements, but the same optimizations apply for where clauses in Delete and UPDATE statements.



#### Note

Because work on the MySQL optimizer is ongoing, not all of the optimizations that MySQL performs are documented here.

You might be tempted to rewrite your queries to make arithmetic operations faster, while sacrificing readability. Because MySQL does similar optimizations automatically, you can often avoid this work, and leave the query in a more understandable and maintainable form. Some of the optimizations performed by MySQL follow:

· Removal of unnecessary parentheses:

```
((a AND b) AND c OR (((a AND b) AND (c AND d))))
-> (a AND b AND c) OR (a AND b AND c AND d)
```

· Constant folding:

```
(a<b AND b=c) AND a=5
-> b>5 AND b=c AND a=5
```

· Constant condition removal:

```
(b>=5 AND b=5) OR (b=6 AND 5=5) OR (b=7 AND 5=6)
-> b=5 OR b=6
```

This takes place during preparation rather than during the optimization phase, which helps in simplification of joins. See Section 10.2.1.9, "Outer Join Optimization", for further information and examples.

- · Constant expressions used by indexes are evaluated only once.
- Comparisons of columns of numeric types with constant values are checked and folded or removed for invalid or out-of-rage values:

```
# CREATE TABLE t (c TINYINT UNSIGNED NOT NULL);

SELECT * FROM t WHERE c < 256;

→≫ SELECT * FROM t WHERE 1;
```

See Section 10.2.1.14, "Constant-Folding Optimization", for more information.

- COUNT(\*) on a single table without a WHERE is retrieved directly from the table information for MyISAM and MEMORY tables. This is also done for any NOT NULL expression when used with only one table.
- Early detection of invalid constant expressions. MySQL quickly detects that some SELECT statements are impossible and returns no rows.
- HAVING is merged with WHERE if you do not use GROUP BY or aggregate functions (COUNT(), MIN(), and so on).
- For each table in a join, a simpler WHERE is constructed to get a fast WHERE evaluation for the table and also to skip rows as soon as possible.
- All constant tables are read first before any other tables in the query. A constant table is any of the following:
  - · An empty table or a table with one row.
  - A table that is used with a WHERE clause on a PRIMARY KEY or a UNIQUE index, where all index
    parts are compared to constant expressions and are defined as NOT NULL.

All of the following tables are used as constant tables:

```
SELECT * FROM t WHERE primary_key=1;
SELECT * FROM t1,t2
WHERE t1.primary_key=1 AND t2.primary_key=t1.id;
```

- The best join combination for joining the tables is found by trying all possibilities. If all columns in ORDER BY and GROUP BY clauses come from the same table, that table is preferred first when joining.
- If there is an ORDER BY clause and a different GROUP BY clause, or if the ORDER BY OR GROUP BY contains columns from tables other than the first table in the join queue, a temporary table is created.
- If you use the SQL\_SMALL\_RESULT modifier, MySQL uses an in-memory temporary table.
- Each table index is queried, and the best index is used unless the optimizer believes that it is more
  efficient to use a table scan. At one time, a scan was used based on whether the best index spanned
  more than 30% of the table, but a fixed percentage no longer determines the choice between using

an index or a scan. The optimizer now is more complex and bases its estimate on additional factors such as table size, number of rows, and I/O block size.

- In some cases, MySQL can read rows from the index without even consulting the data file. If all columns used from the index are numeric, only the index tree is used to resolve the query.
- Before each row is output, those that do not match the HAVING clause are skipped.

Some examples of queries that are very fast:

```
SELECT COUNT(*) FROM tbl_name;

SELECT MIN(key_part1), MAX(key_part1) FROM tbl_name;

SELECT MAX(key_part2) FROM tbl_name
WHERE key_part1=constant;

SELECT ... FROM tbl_name
ORDER BY key_part1, key_part2,... LIMIT 10;

SELECT ... FROM tbl_name
ORDER BY key_part1 DESC, key_part2 DESC, ... LIMIT 10;
```

MySQL resolves the following queries using only the index tree, assuming that the indexed columns are numeric:

```
SELECT key_part1, key_part2 FROM tbl_name WHERE key_part1=val;

SELECT COUNT(*) FROM tbl_name
WHERE key_part1=val1 AND key_part2=val2;

SELECT MAX(key_part2) FROM tbl_name GROUP BY key_part1;
```

The following queries use indexing to retrieve the rows in sorted order without a separate sorting pass:

```
SELECT ... FROM tbl_name
ORDER BY key_part1, key_part2,...;

SELECT ... FROM tbl_name
ORDER BY key_part1 DESC, key_part2 DESC, ...;
```

#### 10.2.1.2 Range Optimization

The range access method uses a single index to retrieve a subset of table rows that are contained within one or several index value intervals. It can be used for a single-part or multiple-part index. The following sections describe conditions under which the optimizer uses range access.

- Range Access Method for Single-Part Indexes
- Range Access Method for Multiple-Part Indexes
- Equality Range Optimization of Many-Valued Comparisons
- Skip Scan Range Access Method
- Range Optimization of Row Constructor Expressions
- · Limiting Memory Use for Range Optimization

#### Range Access Method for Single-Part Indexes

For a single-part index, index value intervals can be conveniently represented by corresponding conditions in the WHERE clause, denoted as *range conditions* rather than "intervals."

The definition of a range condition for a single-part index is as follows:

• For both BTREE and HASH indexes, comparison of a key part with a constant value is a range condition when using the =, <=>, IN(), IS NULL, or IS NOT NULL operators.

- Additionally, for BTREE indexes, comparison of a key part with a constant value is a range condition
  when using the >, <, >=, <=, BETWEEN, !=, or <> operators, or LIKE comparisons if the argument to
  LIKE is a constant string that does not start with a wildcard character.
- For all index types, multiple range conditions combined with OR or AND form a range condition.

"Constant value" in the preceding descriptions means one of the following:

- · A constant from the query string
- A column of a const or system table from the same join
- The result of an uncorrelated subquery
- Any expression composed entirely from subexpressions of the preceding types

Here are some examples of queries with range conditions in the WHERE clause:

```
SELECT * FROM t1

WHERE key_col > 1

AND key_col < 10;

SELECT * FROM t1

WHERE key_col = 1

OR key_col IN (15,18,20);

SELECT * FROM t1

WHERE key_col LIKE 'ab%'

OR key_col BETWEEN 'bar' AND 'foo';
```

Some nonconstant values may be converted to constants during the optimizer constant propagation phase.

MySQL tries to extract range conditions from the WHERE clause for each of the possible indexes. During the extraction process, conditions that cannot be used for constructing the range condition are dropped, conditions that produce overlapping ranges are combined, and conditions that produce empty ranges are removed.

Consider the following statement, where key1 is an indexed column and nonkey is not indexed:

```
SELECT * FROM t1 WHERE
  (key1 < 'abc' AND (key1 LIKE 'abcde%' OR key1 LIKE '%b')) OR
  (key1 < 'bar' AND nonkey = 4) OR
  (key1 < 'uux' AND key1 > 'z');
```

The extraction process for key key1 is as follows:

1. Start with original WHERE clause:

```
(key1 < 'abc' AND (key1 LIKE 'abcde%' OR key1 LIKE '%b')) OR
(key1 < 'bar' AND nonkey = 4) OR
(key1 < 'uux' AND key1 > 'z')
```

2. Remove nonkey = 4 and key1 LIKE '%b' because they cannot be used for a range scan. The correct way to remove them is to replace them with TRUE, so that we do not miss any matching rows when doing the range scan. Replacing them with TRUE yields:

```
(key1 < 'abc' AND (key1 LIKE 'abcde%' OR TRUE)) OR
(key1 < 'bar' AND TRUE) OR
(key1 < 'uux' AND key1 > 'z')
```

- 3. Collapse conditions that are always true or false:
  - (key1 LIKE 'abcde%' OR TRUE) is always true
  - (key1 < 'uux' AND key1 > 'z') is always false

Replacing these conditions with constants yields:

```
(key1 < 'abc' AND TRUE) OR (key1 < 'bar' AND TRUE) OR (FALSE)
```

Removing unnecessary TRUE and FALSE constants yields:

```
(key1 < 'abc') OR (key1 < 'bar')
```

4. Combining overlapping intervals into one yields the final condition to be used for the range scan:

```
(key1 < 'bar')
```

In general (and as demonstrated by the preceding example), the condition used for a range scan is less restrictive than the WHERE clause. MySQL performs an additional check to filter out rows that satisfy the range condition but not the full WHERE clause.

The range condition extraction algorithm can handle nested AND/OR constructs of arbitrary depth, and its output does not depend on the order in which conditions appear in WHERE clause.

MySQL does not support merging multiple ranges for the range access method for spatial indexes. To work around this limitation, you can use a UNION with identical SELECT statements, except that you put each spatial predicate in a different SELECT.

#### Range Access Method for Multiple-Part Indexes

Range conditions on a multiple-part index are an extension of range conditions for a single-part index. A range condition on a multiple-part index restricts index rows to lie within one or several key tuple intervals. Key tuple intervals are defined over a set of key tuples, using ordering from the index.

For example, consider a multiple-part index defined as key1(key\_part1, key\_part2, key\_part3), and the following set of key tuples listed in key order:

```
key_part1 key_part2 key_part3
  NULL
                          'abc'
             1
                          'xyz'
  NULT.
             1
  NULL
             2
                          'foo'
             1
                          'abc'
   1
   1
              1
                          'xyz'
                          'abc'
   1
              2
   2
                          'aaa'
```

The condition  $key\_part1 = 1$  defines this interval:

```
(1,-inf,-inf) <= (key_part1,key_part2,key_part3) < (1,+inf,+inf)
```

The interval covers the 4th, 5th, and 6th tuples in the preceding data set and can be used by the range access method.

By contrast, the condition  $key\_part3 = 'abc'$  does not define a single interval and cannot be used by the range access method.

The following descriptions indicate how range conditions work for multiple-part indexes in greater detail.

• For HASH indexes, each interval containing identical values can be used. This means that the interval can be produced only for conditions in the following form:

```
key_part1 cmp const1
AND key_part2 cmp const2
AND ...
AND key_partN cmp constN;
```

Here, const1, const2, ... are constants, cmp is one of the =, <=>, or IS NULL comparison operators, and the conditions cover all index parts. (That is, there are N conditions, one for each part of an N-part index.) For example, the following is a range condition for a three-part HASH index:

```
key_part1 = 1 AND key_part2 IS NULL AND key_part3 = 'foo'
```

For the definition of what is considered to be a constant, see Range Access Method for Single-Part Indexes.

• For a BTREE index, an interval might be usable for conditions combined with AND, where each condition compares a key part with a constant value using =, <=>, IS NULL, >, <, >=, <=, !=, <>, BETWEEN, or LIKE 'pattern' (where 'pattern' does not start with a wildcard). An interval can be used as long as it is possible to determine a single key tuple containing all rows that match the condition (or two intervals if <> or != is used).

The optimizer attempts to use additional key parts to determine the interval as long as the comparison operator is =, <=>, or IS NULL. If the operator is >, <, >=, <=, !=, <>, BETWEEN, or LIKE, the optimizer uses it but considers no more key parts. For the following expression, the optimizer uses = from the first comparison. It also uses >= from the second comparison but considers no further key parts and does not use the third comparison for interval construction:

```
key_part1 = 'foo' AND key_part2 >= 10 AND key_part3 > 10
```

The single interval is:

```
('foo',10,-inf) < (key_part1,key_part2,key_part3) < ('foo',+inf,+inf)
```

It is possible that the created interval contains more rows than the initial condition. For example, the preceding interval includes the value ('foo', 11, 0), which does not satisfy the original condition.

If conditions that cover sets of rows contained within intervals are combined with OR, they form a
condition that covers a set of rows contained within the union of their intervals. If the conditions are
combined with AND, they form a condition that covers a set of rows contained within the intersection
of their intervals. For example, for this condition on a two-part index:

```
(key_part1 = 1 AND key_part2 < 2) OR (key_part1 > 5)
```

The intervals are:

```
(1,-inf) < (key_part1,key_part2) < (1,2) 
(5,-inf) < (key_part1,key_part2)
```

In this example, the interval on the first line uses one key part for the left bound and two key parts for the right bound. The interval on the second line uses only one key part. The key\_len column in the EXPLAIN output indicates the maximum length of the key prefix used.

In some cases,  $key\_len$  may indicate that a key part was used, but that might be not what you would expect. Suppose that  $key\_part1$  and  $key\_part2$  can be NULL. Then the  $key\_len$  column displays two key part lengths for the following condition:

```
key_part1 >= 1 AND key_part2 < 2</pre>
```

But, in fact, the condition is converted to this:

```
key_part1 >= 1 AND key_part2 IS NOT NULL
```

For a description of how optimizations are performed to combine or eliminate intervals for range conditions on a single-part index, see Range Access Method for Single-Part Indexes. Analogous steps are performed for range conditions on multiple-part indexes.

#### **Equality Range Optimization of Many-Valued Comparisons**

Consider these expressions, where col\_name is an indexed column:

```
col_name IN(val1, ..., valN)
col_name = val1 OR ... OR col_name = valN
```

Each expression is true if col\_name is equal to any of several values. These comparisons are equality range comparisons (where the "range" is a single value). The optimizer estimates the cost of reading qualifying rows for equality range comparisons as follows:

- If there is a unique index on col\_name, the row estimate for each range is 1 because at most one row can have the given value.
- Otherwise, any index on col\_name is nonunique and the optimizer can estimate the row count for each range using dives into the index or index statistics.

With index dives, the optimizer makes a dive at each end of a range and uses the number of rows in the range as the estimate. For example, the expression <code>col\_name IN (10, 20, 30)</code> has three equality ranges and the optimizer makes two dives per range to generate a row estimate. Each pair of dives yields an estimate of the number of rows that have the given value.

Index dives provide accurate row estimates, but as the number of comparison values in the expression increases, the optimizer takes longer to generate a row estimate. Use of index statistics is less accurate than index dives but permits faster row estimation for large value lists.

The eq\_range\_index\_dive\_limit system variable enables you to configure the number of values at which the optimizer switches from one row estimation strategy to the other. To permit use of index dives for comparisons of up to N equality ranges, set eq\_range\_index\_dive\_limit to N+1. To disable use of statistics and always use index dives regardless of N, set eq\_range\_index\_dive\_limit to N.

To update table index statistics for best estimates, use ANALYZE TABLE.

Prior to MySQL 8.4, there is no way of skipping the use of index dives to estimate index usefulness, except by using the eq\_range\_index\_dive\_limit system variable. In MySQL 8.4, index dive skipping is possible for queries that satisfy all these conditions:

- The query is for a single table, not a join on multiple tables.
- A single-index FORCE INDEX index hint is present. The idea is that if index use is forced, there is nothing to be gained from the additional overhead of performing dives into the index.
- The index is nonunique and not a FULLTEXT index.
- · No subquery is present.
- No distinct, group by, or order by clause is present.

For EXPLAIN FOR CONNECTION, the output changes as follows if index dives are skipped:

- For traditional output, the rows and filtered values are NULL.
- For JSON output, rows\_examined\_per\_scan and rows\_produced\_per\_join do not appear, skip\_index\_dive\_due\_to\_force is true, and cost calculations are not accurate.

Without FOR CONNECTION, EXPLAIN output does not change when index dives are skipped.

After execution of a query for which index dives are skipped, the corresponding row in the Information Schema OPTIMIZER\_TRACE table contains an index\_dives\_for\_range\_access value of skipped\_due\_to\_force\_index.

#### **Skip Scan Range Access Method**

Consider the following scenario:

CREATE TABLE t1 (f1 INT NOT NULL, f2 INT NOT NULL, PRIMARY KEY(f1, f2));

```
INSERT INTO t1 VALUES
  (1,1), (1,2), (1,3), (1,4), (1,5),
  (2,1), (2,2), (2,3), (2,4), (2,5);
INSERT INTO t1 SELECT f1, f2 + 5 FROM t1;
INSERT INTO t1 SELECT f1, f2 + 10 FROM t1;
INSERT INTO t1 SELECT f1, f2 + 20 FROM t1;
INSERT INTO t1 SELECT f1, f2 + 40 FROM t1;
ANALYZE TABLE t1;
EXPLAIN SELECT f1, f2 FROM t1 WHERE f2 > 40;
```

To execute this query, MySQL can choose an index scan to fetch all rows (the index includes all columns to be selected), then apply the f2 > 40 condition from the WHERE clause to produce the final result set.

A range scan is more efficient than a full index scan, but cannot be used in this case because there is no condition on fl, the first index column. The optimizer can perform multiple range scans, one for each value of fl, using a method called Skip Scan that is similar to Loose Index Scan (see Section 10.2.1.17, "GROUP BY Optimization"):

- 1. Skip between distinct values of the first index part, £1 (the index prefix).
- 2. Perform a subrange scan on each distinct prefix value for the £2 > 40 condition on the remaining index part.

For the data set shown earlier, the algorithm operates like this:

- 1. Get the first distinct value of the first key part (f1 = 1).
- 2. Construct the range based on the first and second key parts (f1 = 1 AND f2 > 40).
- 3. Perform a range scan.
- 4. Get the next distinct value of the first key part (f1 = 2).
- 5. Construct the range based on the first and second key parts (f1 = 2 AND f2 > 40).
- Perform a range scan.

Using this strategy decreases the number of accessed rows because MySQL skips the rows that do not qualify for each constructed range. This Skip Scan access method is applicable under the following conditions:

- Table T has at least one compound index with key parts of the form ([A\_1, ..., A\_k,] B\_1, ..., B\_m, C [, D\_1, ..., D\_n]). Key parts A and D may be empty, but B and C must be nonempty.
- The query references only one table.
- The query does not use GROUP BY or DISTINCT.
- The query references only columns in the index.
- The predicates on A\_1, ..., A\_k must be equality predicates and they must be constants. This includes the IN() operator.
- The query must be a conjunctive query; that is, an AND of OR conditions: (cond1(key\_part1) OR cond2(key\_part1)) AND (cond1(key\_part2) OR ...) AND ...
- There must be a range condition on C.
- Conditions on D columns are permitted. Conditions on D must be in conjunction with the range condition on C.

Use of Skip Scan is indicated in EXPLAIN output as follows:

- Using index for skip scan in the Extra column indicates that the loose index Skip Scan access method is used.
- If the index can be used for Skip Scan, the index should be visible in the possible keys column.

Use of Skip Scan is indicated in optimizer trace output by a "skip scan" element of this form:

```
"skip_scan_range": {
   "type": "skip_scan",
   "index": index_used_for_skip_scan,
   "key_parts_used_for_access": [key_parts_used_for_access],
   "range": [range]
}
```

You may also see a "best\_skip\_scan\_summary" element. If Skip Scan is chosen as the best range access variant, a "chosen\_range\_access\_summary" is written. If Skip Scan is chosen as the overall best access method, a "best\_access\_path" element is present.

Use of Skip Scan is subject to the value of the <code>skip\_scan</code> flag of the <code>optimizer\_switch</code> system variable. See Section 10.9.2, "Switchable Optimizations". By default, this flag is on. To disable it, set <code>skip\_scan</code> to <code>off</code>.

In addition to using the <code>optimizer\_switch</code> system variable to control optimizer use of Skip Scan session-wide, MySQL supports optimizer hints to influence the optimizer on a per-statement basis. See Section 10.9.3, "Optimizer Hints".

#### **Range Optimization of Row Constructor Expressions**

The optimizer is able to apply the range scan access method to queries of this form:

```
SELECT ... FROM t1 WHERE ( col_1, col_2 ) IN (( 'a', 'b' ), ( 'c', 'd' ));
```

Previously, for range scans to be used, it was necessary to write the guery as:

```
SELECT ... FROM t1 WHERE ( col_1 = 'a' AND col_2 = 'b' )
OR ( col_1 = 'c' AND col_2 = 'd' );
```

For the optimizer to use a range scan, queries must satisfy these conditions:

- Only IN() predicates are used, not NOT IN().
- On the left side of the IN() predicate, the row constructor contains only column references.
- On the right side of the IN() predicate, row constructors contain only runtime constants, which are either literals or local column references that are bound to constants during execution.
- On the right side of the IN() predicate, there is more than one row constructor.

For more information about the optimizer and row constructors, see Section 10.2.1.22, "Row Constructor Expression Optimization"

#### **Limiting Memory Use for Range Optimization**

To control the memory available to the range optimizer, use the range\_optimizer\_max\_mem\_size system variable:

- A value of 0 means "no limit."
- With a value greater than 0, the optimizer tracks the memory consumed when considering the
  range access method. If the specified limit is about to be exceeded, the range access method
  is abandoned and other methods, including a full table scan, are considered instead. This
  could be less optimal. If this happens, the following warning occurs (where N is the current
  range\_optimizer\_max\_mem\_size value):

```
Warning 3170 Memory capacity of N bytes for 
'range_optimizer_max_mem_size' exceeded. Range 
optimization was not done for this query.
```

• For UPDATE and DELETE statements, if the optimizer falls back to a full table scan and the sql\_safe\_updates system variable is enabled, an error occurs rather than a warning because, in effect, no key is used to determine which rows to modify. For more information, see Using Safe-Updates Mode (--safe-updates).

For individual queries that exceed the available range optimization memory and for which the optimizer falls back to less optimal plans, increasing the range\_optimizer\_max\_mem\_size value may improve performance.

To estimate the amount of memory needed to process a range expression, use these guidelines:

 For a simple query such as the following, where there is one candidate key for the range access method, each predicate combined with OR uses approximately 230 bytes:

```
SELECT COUNT(*) FROM t
WHERE a=1 OR a=2 OR a=3 OR .. . a=N;
```

 Similarly for a query such as the following, each predicate combined with AND uses approximately 125 bytes:

```
SELECT COUNT(*) FROM t
WHERE a=1 AND b=1 AND c=1 ... N;
```

For a query with IN() predicates:

```
SELECT COUNT(*) FROM t WHERE a IN (1,2, \ldots, M) AND b IN (1,2, \ldots, N);
```

Each literal value in an IN() list counts as a predicate combined with OR. If there are two IN() lists, the number of predicates combined with OR is the product of the number of literal values in each list. Thus, the number of predicates combined with OR in the preceding case is  $M \times N$ .

#### 10.2.1.3 Index Merge Optimization

The *Index Merge* access method retrieves rows with multiple range scans and merges their results into one. This access method merges index scans from a single table only, not scans across multiple tables. The merge can produce unions, intersections, or unions-of-intersections of its underlying scans.

Example queries for which Index Merge may be used:

```
SELECT * FROM tbl_name WHERE key1 = 10 OR key2 = 20;

SELECT * FROM tbl_name

WHERE (key1 = 10 OR key2 = 20) AND non_key = 30;

SELECT * FROM t1, t2

WHERE (t1.key1 IN (1,2) OR t1.key2 LIKE 'value%')

AND t2.key1 = t1.some_col;

SELECT * FROM t1, t2

WHERE t1.key1 = 1

AND (t2.key1 = t1.some_col OR t2.key2 = t1.some_col2);
```



#### Note

The Index Merge optimization algorithm has the following known limitations:

 If your query has a complex WHERE clause with deep AND/OR nesting and MySQL does not choose the optimal plan, try distributing terms using the following identity transformations:

```
(x \text{ AND } y) \text{ OR } z \Rightarrow (x \text{ OR } z) \text{ AND } (y \text{ OR } z)
(x \text{ OR } y) \text{ AND } z \Rightarrow (x \text{ AND } z) \text{ OR } (y \text{ AND } z)
```

• Index Merge is not applicable to full-text indexes.

In EXPLAIN output, the Index Merge method appears as index\_merge in the type column. In this case, the key column contains a list of indexes used, and key\_len contains a list of the longest key parts for those indexes.

The Index Merge access method has several algorithms, which are displayed in the Extra field of EXPLAIN output:

- Using intersect(...)
- Using union(...)
- Using sort\_union(...)

The following sections describe these algorithms in greater detail. The optimizer chooses between different possible Index Merge algorithms and other access methods based on cost estimates of the various available options.

- Index Merge Intersection Access Algorithm
- Index Merge Union Access Algorithm
- Index Merge Sort-Union Access Algorithm
- Influencing Index Merge Optimization

#### **Index Merge Intersection Access Algorithm**

This access algorithm is applicable when a WHERE clause is converted to several range conditions on different keys combined with AND, and each condition is one of the following:

 An N-part expression of this form, where the index has exactly N parts (that is, all index parts are covered):

```
key_part1 = const1 AND key_part2 = const2 ... AND key_partN = constN
```

• Any range condition over the primary key of an InnoDB table.

#### Examples:

```
SELECT * FROM innodb_table
WHERE primary_key < 10 AND key_col1 = 20;

SELECT * FROM tbl_name
WHERE key1_part1 = 1 AND key1_part2 = 2 AND key2 = 2;</pre>
```

The Index Merge intersection algorithm performs simultaneous scans on all used indexes and produces the intersection of row sequences that it receives from the merged index scans.

If all columns used in the query are covered by the used indexes, full table rows are not retrieved (EXPLAIN output contains Using index in Extra field in this case). Here is an example of such a query:

```
SELECT COUNT(*) FROM t1 WHERE key1 = 1 AND key2 = 1;
```

If the used indexes do not cover all columns used in the query, full rows are retrieved only when the range conditions for all used keys are satisfied.

If one of the merged conditions is a condition over the primary key of an InnoDB table, it is not used for row retrieval, but is used to filter out rows retrieved using other conditions.

#### **Index Merge Union Access Algorithm**

The criteria for this algorithm are similar to those for the Index Merge intersection algorithm. The algorithm is applicable when the table's WHERE clause is converted to several range conditions on different keys combined with OR, and each condition is one of the following:

 An N-part expression of this form, where the index has exactly N parts (that is, all index parts are covered):

```
key_part1 = const1 OR key_part2 = const2 ... OR key_partN = constN
```

- Any range condition over a primary key of an InnoDB table.
- A condition for which the Index Merge intersection algorithm is applicable.

#### Examples:

```
SELECT * FROM t1
WHERE key1 = 1 OR key2 = 2 OR key3 = 3;

SELECT * FROM innodb\_table
WHERE (key1 = 1 \text{ AND } key2 = 2)
OR (key3 = 'foo' \text{ AND } key4 = 'bar') AND key5 = 5;
```

#### **Index Merge Sort-Union Access Algorithm**

This access algorithm is applicable when the WHERE clause is converted to several range conditions combined by OR, but the Index Merge union algorithm is not applicable.

#### Examples:

```
SELECT * FROM tbl_name
WHERE key_col1 < 10 OR key_col2 < 20;

SELECT * FROM tbl_name
WHERE (key_col1 > 10 OR key_col2 = 20) AND nonkey_col = 30;
```

The difference between the sort-union algorithm and the union algorithm is that the sort-union algorithm must first fetch row IDs for all rows and sort them before returning any rows.

#### **Influencing Index Merge Optimization**

Use of Index Merge is subject to the value of the index\_merge, index\_merge\_intersection, index\_merge\_union, and index\_merge\_sort\_union flags of the optimizer\_switch system variable. See Section 10.9.2, "Switchable Optimizations". By default, all those flags are on. To enable only certain algorithms, set index\_merge to off, and enable only such of the others as should be permitted.

In addition to using the <code>optimizer\_switch</code> system variable to control optimizer use of the Index Merge algorithms session-wide, MySQL supports optimizer hints to influence the optimizer on a perstatement basis. See Section 10.9.3, "Optimizer Hints".

#### 10.2.1.4 Hash Join Optimization

By default, MySQL employs hash joins whenever possible. It is possible to control whether hash joins are employed using one of the BNL and NO\_BNL optimizer hints, or by setting block\_nested\_loop=on or block\_nested\_loop=off as part of the setting for the optimizer switch server system variable.

MySQL employs a hash join for any query for which each join has an equi-join condition, and in which there are no indexes that can be applied to any join conditions, such as this one:

```
SELECT *
```

```
FROM t1

JOIN t2

ON t1.c1=t2.c1;
```

A hash join can also be used when there are one or more indexes that can be used for single-table predicates.

In the example just shown and the remaining examples in this section, we assume that the three tables t1, t2, and t3 have been created using the following statements:

```
CREATE TABLE t1 (c1 INT, c2 INT);
CREATE TABLE t2 (c1 INT, c2 INT);
CREATE TABLE t3 (c1 INT, c2 INT);
```

You can see that a hash join is being employed by using EXPLAIN, like this:

```
mysql> EXPLAIN
   -> SELECT * FROM t1
       JOIN t2 ON t1.c1=t2.c1\G
id: 1
 select_type: SIMPLE
      table: t1
  partitions: NULL
       type: ALL
possible_keys: NULL
        key: NULL
    key_len: NULL
       ref: NULL
       rows: 1
    filtered: 100.00
      Extra: NULL
             ******* 2. row ****************
        id: 1
 select type: SIMPLE
      table: t2
  partitions: NULL
       type: ALL
possible_keys: NULL
        key: NULL
     key_len: NULL
       ref: NULL
       rows: 1
    filtered: 100.00
      Extra: Using where; Using join buffer (hash join)
```

EXPLAIN ANALYZE also displays information about hash joins used.

The hash join is used for queries involving multiple joins as well, as long as at least one join condition for each pair of tables is an equi-join, like the query shown here:

```
SELECT * FROM t1

JOIN t2 ON (t1.c1 = t2.c1 AND t1.c2 < t2.c2)

JOIN t3 ON (t2.c1 = t3.c1);
```

In cases like the one just shown, which makes use of an inner join, any extra conditions which are not equi-joins are applied as filters after the join is executed. (For outer joins, such as left joins, semijoins, and antijoins, they are printed as part of the join.) This can be seen here in the output of EXPLAIN:

```
mysql> EXPLAIN FORMAT=TREE
    -> SELECT *
    -> FROM t1
    -> JOIN t2
    -> ON (t1.c1 = t2.c1 AND t1.c2 < t2.c2)
    -> JOIN t3
    -> ON (t2.c1 = t3.c1)\G
********************************
EXPLAIN: -> Inner hash join (t3.c1 = t1.c1) (cost=1.05 rows=1)
    -> Table scan on t3 (cost=0.35 rows=1)
```

```
-> Hash
-> Filter: (t1.c2 < t2.c2) (cost=0.70 rows=1)
-> Inner hash join (t2.c1 = t1.c1) (cost=0.70 rows=1)
-> Table scan on t2 (cost=0.35 rows=1)
-> Hash
-> Table scan on t1 (cost=0.35 rows=1)
```

As also can be seen from the output just shown, multiple hash joins can be (and are) used for joins having multiple equi-join conditions.

A hash join is used even if any pair of joined tables does not have at least one equi-join condition, as shown here:

```
mysql> EXPLAIN FORMAT=TREE
    -> SELECT * FROM t1
    -> JOIN t2 ON (t1.c1 = t2.c1)
    -> JOIN t3 ON (t2.c1 < t3.c1)\G

********************************

EXPLAIN: -> Filter: (t1.c1 < t3.c1) (cost=1.05 rows=1)
    -> Inner hash join (no condition) (cost=1.05 rows=1)
    -> Table scan on t3 (cost=0.35 rows=1)
    -> Hash
    -> Inner hash join (t2.c1 = t1.c1) (cost=0.70 rows=1)
    -> Table scan on t2 (cost=0.35 rows=1)
    -> Hash
    -> Table scan on t1 (cost=0.35 rows=1)
```

(Additional examples are provided later in this section.)

A hash join is also applied for a Cartesian product—that is, when no join condition is specified, as shown here:

```
mysql> EXPLAIN FORMAT=TREE
    -> SELECT *
    -> FROM t1
    -> JOIN t2
    -> WHERE t1.c2 > 50\G
*****************************
EXPLAIN: -> Inner hash join (cost=0.70 rows=1)
    -> Table scan on t2 (cost=0.35 rows=1)
    -> Hash
    -> Filter: (t1.c2 > 50) (cost=0.35 rows=1)
    -> Table scan on t1 (cost=0.35 rows=1)
```

It is not necessary for the join to contain at least one equi-join condition in order for a hash join to be used. This means that the types of queries which can be optimized using hash joins include those in the following list (with examples):

Inner non-equi-join:

```
mysql> EXPLAIN FORMAT=TREE SELECT * FROM t1 JOIN t2 ON t1.c1 < t2.c1\G
**********
EXPLAIN: -> Filter: (t1.c1 < t2.c1) (cost=4.70 rows=12)
-> Inner hash join (no condition) (cost=4.70 rows=12)
-> Table scan on t2 (cost=0.08 rows=6)
-> Hash
-> Table scan on t1 (cost=0.85 rows=6)
```

• Semijoin:

• Antijoin:

· Left outer join:

```
mysql> EXPLAIN FORMAT=TREE SELECT * FROM t1 LEFT JOIN t2 ON t1.c1 = t2.c1\G
******************************
EXPLAIN: -> Left hash join (t2.c1 = t1.c1) (cost=0.70 rows=1)
    -> Table scan on t1 (cost=0.35 rows=1)
    -> Hash
    -> Table scan on t2 (cost=0.35 rows=1)
```

Right outer join (observe that MySQL rewrites all right outer joins as left outer joins):

```
mysql> EXPLAIN FORMAT=TREE SELECT * FROM t1 RIGHT JOIN t2 ON t1.c1 = t2.c1\G
***********
EXPLAIN: -> Left hash join (t1.c1 = t2.c1) (cost=0.70 rows=1)
    -> Table scan on t2 (cost=0.35 rows=1)
    -> Hash
    -> Table scan on t1 (cost=0.35 rows=1)
```

By default, MySQL employs hash joins whenever possible. It is possible to control whether hash joins are employed using one of the BNL and NO\_BNL optimizer hints.

Memory usage by hash joins can be controlled using the <code>join\_buffer\_size</code> system variable; a hash join cannot use more memory than this amount. When the memory required for a hash join exceeds the amount available, MySQL handles this by using files on disk. If this happens, you should be aware that the join may not succeed if a hash join cannot fit into memory and it creates more files than set for <code>open\_files\_limit</code>. To avoid such problems, make either of the following changes:

- Increase join buffer size so that the hash join does not spill over to disk.
- Increase open\_files\_limit.

Join buffers for hash joins are allocated incrementally; thus, you can set <code>join\_buffer\_size</code> higher without small queries allocating very large amounts of RAM, but outer joins allocate the entire buffer. Hash joins are used for outer joins (including antijoins and semijoins) as well, so this is no longer an issue.

#### 10.2.1.5 Engine Condition Pushdown Optimization

This optimization improves the efficiency of direct comparisons between a nonindexed column and a constant. In such cases, the condition is "pushed down" to the storage engine for evaluation. This optimization can be used only by the NDB storage engine.

For NDB Cluster, this optimization can eliminate the need to send nonmatching rows over the network between the cluster's data nodes and the MySQL server that issued the query, and can speed up queries where it is used by a factor of 5 to 10 times over cases where condition pushdown could be but is not used.

Suppose that an NDB Cluster table is defined as follows:

```
CREATE TABLE t1 (
```

```
a INT,
b INT,
KEY(a)
) ENGINE=NDB;
```

Engine condition pushdown can be used with queries such as the one shown here, which includes a comparison between a nonindexed column and a constant:

```
SELECT a, b FROM t1 WHERE b = 10;
```

The use of engine condition pushdown can be seen in the output of EXPLAIN:

```
mysql> EXPLAIN SELECT a, b FROM t1 WHERE b = 10\G

*****************************
    id: 1
select_type: SIMPLE
    table: t1
    type: ALL
possible_keys: NULL
    key: NULL
    key: NULL
    key-len: NULL
    ref: NULL
    rows: 10
    Extra: Using where with pushed condition
```

However, engine condition pushdown cannot be used with the following query:

```
SELECT a,b FROM t1 WHERE a = 10;
```

Engine condition pushdown is not applicable here because an index exists on column a. (An index access method would be more efficient and so would be chosen in preference to condition pushdown.)

Engine condition pushdown may also be employed when an indexed column is compared with a constant using a > or < operator:

```
mysql> EXPLAIN SELECT a, b FROM t1 WHERE a < 2\G
**************************
    id: 1
select_type: SIMPLE
    table: t1
    type: range
possible_keys: a
    key: a
    key_len: 5
    ref: NULL
    rows: 2
    Extra: Using where with pushed condition</pre>
```

Other supported comparisons for engine condition pushdown include the following:

• column [NOT] LIKE pattern

pattern must be a string literal containing the pattern to be matched; for syntax, see Section 14.8.1, "String Comparison Functions and Operators".

- column IS [NOT] NULL
- column IN (value\_list)

Each item in the value\_list must be a constant, literal value.

• column BETWEEN constant1 AND constant2

constant1 and constant2 must each be a constant, literal value.

In all of the cases in the preceding list, it is possible for the condition to be converted into the form of one or more direct comparisons between a column and a constant.

Engine condition pushdown is enabled by default. To disable it at server startup, set the optimizer\_switch system variable's engine\_condition\_pushdown flag to off. For example, in a my.cnf file, use these lines:

```
[mysqld]
optimizer_switch=engine_condition_pushdown=off
```

At runtime, disable condition pushdown like this:

```
SET optimizer_switch='engine_condition_pushdown=off';
```

**Limitations.** Engine condition pushdown is subject to the following limitations:

- Engine condition pushdown is supported only by the NDB storage engine.
- In NDB 8.4, columns can be compared with one another as long as they are of exactly the same type, including the same signedness, length, character set, precision, and scale, where these are applicable.
- Columns used in comparisons cannot be of any of the BLOB or TEXT types. This exclusion extends to JSON, BIT, and ENUM columns as well.
- A string value to be compared with a column must use the same collation as the column.
- Joins are not directly supported; conditions involving multiple tables are pushed separately where
  possible. Use extended EXPLAIN output to determine which conditions are actually pushed down.
  See Section 10.8.3, "Extended EXPLAIN Output Format".

Previously, engine condition pushdown was limited to terms referring to column values from the same table to which the condition was being pushed. In NDB 8.4, column values from tables earlier in the query plan can also be referred to from pushed conditions. This reduces the number of rows which must be handled by the SQL node during join processing. Filtering can be also performed in parallel in the LDM threads, rather than in a single <code>mysqld</code> process. This has the potential to improve performance of queries by a significant margin.

NDB can push an outer join using a scan if there are no unpushable conditions on any table used in the same join nest, or on any table in join nests above it on which it depends. This is also true for a semijoin, provided the optimization strategy employed is firstMatch (see Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations).

Join algorithms cannot be combined with referring columns from previous tables in the following two situations:

- 1. When any of the referred previous tables are in a join buffer. In this case, each row retrieved from the scan-filtered table is matched against every row in the buffer. This means that there is no single specific row from which column values can be fetched from when generating the scan filter.
- 2. When the column originates from a child operation in a pushed join. This is because rows referenced from ancestor operations in the join have not yet been retrieved when the scan filter is generated.

Columns from ancestor tables in a join can be pushed down, provided that they meet the requirements listed previously. An example of such a query, using the table t1 created previously, is shown here:

```
mysql> EXPLAIN
    -> SELECT * FROM t1 AS x
    -> LEFT JOIN t1 AS y
    -> ON x.a=0 AND y.b>=3\G
**************************
    id: 1
select_type: SIMPLE
    table: x
partitions: p0,p1
    type: ALL
```

```
possible_keys: NULL
         key: NULL
     key_len: NULL
         ref: NULL
        rows: 4
     filtered: 100.00
       Extra: NULL
                    ***** 2. row *****
          id: 1
  select_type: SIMPLE
       table: v
  partitions: p0,p1
        type: ALL
possible_keys: NULL
         key: NULL
     key_len: NULL
         ref: NULL
        rows: 4
     filtered: 100.00
       Extra: Using where; Using pushed condition ('test'.'y'.'b' >= 3); Using join buffer (hash join)
2 rows in set, 2 warnings (0.00 sec)
```

#### 10.2.1.6 Index Condition Pushdown Optimization

Index Condition Pushdown (ICP) is an optimization for the case where MySQL retrieves rows from a table using an index. Without ICP, the storage engine traverses the index to locate rows in the base table and returns them to the MySQL server which evaluates the WHERE condition for the rows. With ICP enabled, and if parts of the WHERE condition can be evaluated by using only columns from the index, the MySQL server pushes this part of the WHERE condition down to the storage engine. The storage engine then evaluates the pushed index condition by using the index entry and only if this is satisfied is the row read from the table. ICP can reduce the number of times the storage engine must access the base table and the number of times the MySQL server must access the storage engine.

Applicability of the Index Condition Pushdown optimization is subject to these conditions:

- ICP is used for the range, ref, eq\_ref, and ref\_or\_null access methods when there is a need to access full table rows.
- ICP can be used for InnoDB and MyISAM tables, including partitioned InnoDB and MyISAM tables.
- For InnoDB tables, ICP is used only for secondary indexes. The goal of ICP is to reduce the number of full-row reads and thereby reduce I/O operations. For InnoDB clustered indexes, the complete record is already read into the InnoDB buffer. Using ICP in this case does not reduce I/O.
- ICP is not supported with secondary indexes created on virtual generated columns. InnoDB supports secondary indexes on virtual generated columns.
- Conditions that refer to subqueries cannot be pushed down.
- Conditions that refer to stored functions cannot be pushed down. Storage engines cannot invoke stored functions.
- Triggered conditions cannot be pushed down. (For information about triggered conditions, see Section 10.2.2.3, "Optimizing Subqueries with the EXISTS Strategy".)
- Conditions cannot be pushed down to derived tables containing references to system variables.

To understand how this optimization works, first consider how an index scan proceeds when Index Condition Pushdown is not used:

- 1. Get the next row, first by reading the index tuple, and then by using the index tuple to locate and read the full table row.
- 2. Test the part of the WHERE condition that applies to this table. Accept or reject the row based on the test result.

Using Index Condition Pushdown, the scan proceeds like this instead:

- 1. Get the next row's index tuple (but not the full table row).
- 2. Test the part of the WHERE condition that applies to this table and can be checked using only index columns. If the condition is not satisfied, proceed to the index tuple for the next row.
- 3. If the condition is satisfied, use the index tuple to locate and read the full table row.
- 4. Test the remaining part of the WHERE condition that applies to this table. Accept or reject the row based on the test result.

EXPLAIN output shows Using index condition in the Extra column when Index Condition Pushdown is used. It does not show Using index because that does not apply when full table rows must be read.

Suppose that a table contains information about people and their addresses and that the table has an index defined as INDEX (zipcode, lastname, firstname). If we know a person's zipcode value but are not sure about the last name, we can search like this:

```
SELECT * FROM people
WHERE zipcode='95054'
AND lastname LIKE '%etrunia%'
AND address LIKE '%Main Street%';
```

MySQL can use the index to scan through people with <code>zipcode='95054'</code>. The second part (lastname LIKE '%etrunia%') cannot be used to limit the number of rows that must be scanned, so without Index Condition Pushdown, this query must retrieve full table rows for all people who have <code>zipcode='95054'</code>.

With Index Condition Pushdown, MySQL checks the lastname LIKE '%etrunia%' part before reading the full table row. This avoids reading full rows corresponding to index tuples that match the zipcode condition but not the lastname condition.

Index Condition Pushdown is enabled by default. It can be controlled with the <code>optimizer\_switch</code> system variable by setting the <code>index\_condition\_pushdown</code> flag:

```
SET optimizer_switch = 'index_condition_pushdown=off';
SET optimizer_switch = 'index_condition_pushdown=on';
```

See Section 10.9.2, "Switchable Optimizations".

## 10.2.1.7 Nested-Loop Join Algorithms

MySQL executes joins between tables using a nested-loop algorithm or variations on it.

Nested-Loop Join Algorithm

#### **Nested-Loop Join Algorithm**

A simple nested-loop join (NLJ) algorithm reads rows from the first table in a loop one at a time, passing each row to a nested loop that processes the next table in the join. This process is repeated as many times as there remain tables to be joined.

Assume that a join between three tables t1, t2, and t3 is to be executed using the following join types:

```
Table Join Type
t1 range
t2 ref
t3 ALL
```

If a simple NLJ algorithm is used, the join is processed like this:

```
for each row in t1 matching range {
  for each row in t2 matching reference key {
   for each row in t3 {
     if row satisfies join conditions, send to client
   }
  }
}
```

Because the NLJ algorithm passes rows one at a time from outer loops to inner loops, it typically reads tables processed in the inner loops many times.

#### 10.2.1.8 Nested Join Optimization

The syntax for expressing joins permits nested joins. The following discussion refers to the join syntax described in Section 15.2.13.2, "JOIN Clause".

The syntax of <code>table\_factor</code> is extended in comparison with the SQL Standard. The latter accepts only <code>table\_reference</code>, not a list of them inside a pair of parentheses. This is a conservative extension if we consider each comma in a list of <code>table\_reference</code> items as equivalent to an inner join. For example:

```
SELECT * FROM t1 LEFT JOIN (t2, t3, t4)
ON (t2.a=t1.a AND t3.b=t1.b AND t4.c=t1.c)
```

Is equivalent to:

```
SELECT * FROM t1 LEFT JOIN (t2 CROSS JOIN t3 CROSS JOIN t4)
ON (t2.a=t1.a AND t3.b=t1.b AND t4.c=t1.c)
```

In MySQL, CROSS JOIN is syntactically equivalent to INNER JOIN; they can replace each other. In standard SQL, they are not equivalent. INNER JOIN is used with an ON clause; CROSS JOIN is used otherwise.

In general, parentheses can be ignored in join expressions containing only inner join operations. Consider this join expression:

```
t1 LEFT JOIN (t2 LEFT JOIN t3 ON t2.b=t3.b OR t2.b IS NULL)
ON t1.a=t2.a
```

After removing parentheses and grouping operations to the left, that join expression transforms into this expression:

```
(t1 LEFT JOIN t2 ON t1.a=t2.a) LEFT JOIN t3
ON t2.b=t3.b OR t2.b IS NULL
```

Yet, the two expressions are not equivalent. To see this, suppose that the tables t1, t2, and t3 have the following state:

- Table t1 contains rows (1), (2)
- Table t2 contains row (1,101)
- Table t3 contains row (101)

In this case, the first expression returns a result set including the rows (1,1,101,101), (2,NULL,NULL,NULL), whereas the second expression returns the rows (1,1,101,101), (2,NULL,NULL,101):

```
+----+
| 1 | 1 | 101 | 101 |
| 2 | NULL | NULL | NULL |
+----+
| mysql> SELECT *
| FROM (t1 LEFT JOIN t2 ON t1.a=t2.a)
| LEFT JOIN t3
| ON t2.b=t3.b OR t2.b IS NULL;
| a | a | b | b |
| +----+
| 1 | 1 | 101 | 101 |
| 2 | NULL | NULL | 101 |
| +-----+
```

In the following example, an outer join operation is used together with an inner join operation:

```
t1 LEFT JOIN (t2, t3) ON t1.a=t2.a
```

That expression cannot be transformed into the following expression:

```
t1 LEFT JOIN t2 ON t1.a=t2.a, t3
```

For the given table states, the two expressions return different sets of rows:

Therefore, if we omit parentheses in a join expression with outer join operators, we might change the result set for the original expression.

More exactly, we cannot ignore parentheses in the right operand of the left outer join operation and in the left operand of a right join operation. In other words, we cannot ignore parentheses for the inner table expressions of outer join operations. Parentheses for the other operand (operand for the outer table) can be ignored.

The following expression:

```
(t1,t2) LEFT JOIN t3 ON P(t2.b,t3.b)
```

Is equivalent to this expression for any tables t1, t2, t3 and any condition P over attributes t2.b and t3.b:

```
t1, t2 LEFT JOIN t3 ON P(t2.b,t3.b)
```

Whenever the order of execution of join operations in a join expression (<code>joined\_table</code>) is not from left to right, we talk about nested joins. Consider the following queries:

```
SELECT * FROM t1 LEFT JOIN (t2 LEFT JOIN t3 ON t2.b=t3.b) ON t1.a=t2.a
WHERE t1.a > 1

SELECT * FROM t1 LEFT JOIN (t2, t3) ON t1.a=t2.a
WHERE (t2.b=t3.b OR t2.b IS NULL) AND t1.a > 1
```

Those queries are considered to contain these nested joins:

```
t2 LEFT JOIN t3 ON t2.b=t3.b
t2, t3
```

In the first query, the nested join is formed with a left join operation. In the second query, it is formed with an inner join operation.

In the first query, the parentheses can be omitted: The grammatical structure of the join expression dictates the same order of execution for join operations. For the second query, the parentheses cannot be omitted, although the join expression here can be interpreted unambiguously without them. In our extended syntax, the parentheses in (t2, t3) of the second query are required, although theoretically the query could be parsed without them: We still would have unambiguous syntactical structure for the query because LEFT JOIN and ON play the role of the left and right delimiters for the expression (t2,t3).

The preceding examples demonstrate these points:

- For join expressions involving only inner joins (and not outer joins), parentheses can be removed and joins evaluated left to right. In fact, tables can be evaluated in any order.
- The same is not true, in general, for outer joins or for outer joins mixed with inner joins. Removal of parentheses may change the result.

Queries with nested outer joins are executed in the same pipeline manner as queries with inner joins. More exactly, a variation of the nested-loop join algorithm is exploited. Recall the algorithm by which the nested-loop join executes a query (see Section 10.2.1.7, "Nested-Loop Join Algorithms"). Suppose that a join query over 3 tables T1, T2, T3 has this form:

```
SELECT * FROM T1 INNER JOIN T2 ON P1(T1,T2)
INNER JOIN T3 ON P2(T2,T3)
WHERE P(T1,T2,T3)
```

Here, P1(T1,T2) and P2(T3,T3) are some join conditions (on expressions), whereas P(T1,T2,T3) is a condition over columns of tables T1,T2,T3.

The nested-loop join algorithm would execute this query in the following manner:

```
FOR each row t1 in T1 {
   FOR each row t2 in T2 such that P1(t1,t2) {
     FOR each row t3 in T3 such that P2(t2,t3) {
        IF P(t1,t2,t3) {
             t:=t1||t2||t3; OUTPUT t;
        }
     }
   }
}
```

The notation  $t1 \mid t2 \mid t3$  indicates a row constructed by concatenating the columns of rows t1, t2, and t3. In some of the following examples, NULL where a table name appears means a row in which NULL is used for each column of that table. For example,  $t1 \mid t2 \mid NULL$  indicates a row constructed by concatenating the columns of rows t1 and t2, and NULL for each column of t3. Such a row is said to be NULL-complemented.

Now consider a query with nested outer joins:

```
SELECT * FROM T1 LEFT JOIN

(T2 LEFT JOIN T3 ON P2(T2,T3))

ON P1(T1,T2)

WHERE P(T1,T2,T3)
```

For this query, modify the nested-loop pattern to obtain:

```
FOR each row t1 in T1 {
```

```
BOOL f1:=FALSE;
FOR each row t2 in T2 such that P1(t1,t2) {
 BOOL f2:=FALSE;
 FOR each row t3 in T3 such that P2(t2,t3) {
    IF P(t1,t2,t3) {
     t:=t1||t2||t3; OUTPUT t;
    f2=TRUE;
   f1=TRUE;
  IF (!f2) {
    IF P(t1,t2,NULL) {
     t:=t1||t2||NULL; OUTPUT t;
    f1=TRUE;
IF (!f1) {
 IF P(t1,NULL,NULL) {
   t:=t1||NULL||NULL; OUTPUT t;
}
```

In general, for any nested loop for the first inner table in an outer join operation, a flag is introduced that is turned off before the loop and is checked after the loop. The flag is turned on when for the current row from the outer table a match from the table representing the inner operand is found. If at the end of the loop cycle the flag is still off, no match has been found for the current row of the outer table. In this case, the row is complemented by NULL values for the columns of the inner tables. The result row is passed to the final check for the output or into the next nested loop, but only if the row satisfies the join condition of all embedded outer joins.

In the example, the outer join table expressed by the following expression is embedded:

```
(T2 LEFT JOIN T3 ON P2(T2,T3))
```

For the query with inner joins, the optimizer could choose a different order of nested loops, such as this one:

```
FOR each row t3 in T3 {
   FOR each row t2 in T2 such that P2(t2,t3) {
    FOR each row t1 in T1 such that P1(t1,t2) {
        IF P(t1,t2,t3) {
            t:=t1||t2||t3; OUTPUT t;
        }
     }
   }
}
```

For queries with outer joins, the optimizer can choose only such an order where loops for outer tables precede loops for inner tables. Thus, for our query with outer joins, only one nesting order is possible. For the following query, the optimizer evaluates two different nestings. In both nestings, T1 must be processed in the outer loop because it is used in an outer join. T2 and T3 are used in an inner join, so that join must be processed in the inner loop. However, because the join is an inner join, T2 and T3 can be processed in either order.

```
SELECT * T1 LEFT JOIN (T2,T3) ON P1(T1,T2) AND P2(T1,T3)
WHERE P(T1,T2,T3)
```

One nesting evaluates T2, then T3:

```
FOR each row t1 in T1 {
   BOOL f1:=FALSE;
   FOR each row t2 in T2 such that P1(t1,t2) {
      FOR each row t3 in T3 such that P2(t1,t3) {
        IF P(t1,t2,t3) {
            t:=t1||t2||t3; OUTPUT t;
      }
}
```

```
f1:=TRUE
}
}
IF (!f1) {
    IF P(t1,NULL,NULL) {
        t:=t1||NULL||NULL; OUTPUT t;
    }
}
```

The other nesting evaluates T3, then T2:

```
FOR each row t1 in T1 {
    BOOL f1:=FALSE;
    FOR each row t3 in T3 such that P2(t1,t3) {
        FOR each row t2 in T2 such that P1(t1,t2) {
            IF P(t1,t2,t3) {
                t:=t1||t2||t3; OUTPUT t;
            }
            f1:=TRUE
        }
    }
    IF (!f1) {
        IF P(t1,NULL,NULL) {
            t:=t1||NULL||NULL; OUTPUT t;
        }
    }
}
```

When discussing the nested-loop algorithm for inner joins, we omitted some details whose impact on the performance of query execution may be huge. We did not mention so-called "pushed-down" conditions. Suppose that our WHERE condition P(T1,T2,T3) can be represented by a conjunctive formula:

```
P(T1,T2,T2) = C1(T1) AND C2(T2) AND C3(T3).
```

In this case, MySQL actually uses the following nested-loop algorithm for the execution of the query with inner joins:

```
FOR each row t1 in T1 such that C1(t1) {
   FOR each row t2 in T2 such that P1(t1,t2) AND C2(t2) {
     FOR each row t3 in T3 such that P2(t2,t3) AND C3(t3) {
        IF P(t1,t2,t3) {
            t:=t1||t2||t3; OUTPUT t;
        }
    }
}
```

You see that each of the conjuncts C1(T1), C2(T2), C3(T3) are pushed out of the most inner loop to the most outer loop where it can be evaluated. If C1(T1) is a very restrictive condition, this condition pushdown may greatly reduce the number of rows from table T1 passed to the inner loops. As a result, the execution time for the query may improve immensely.

For a query with outer joins, the WHERE condition is to be checked only after it has been found that the current row from the outer table has a match in the inner tables. Thus, the optimization of pushing conditions out of the inner nested loops cannot be applied directly to queries with outer joins. Here we must introduce conditional pushed-down predicates guarded by the flags that are turned on when a match has been encountered.

Recall this example with outer joins:

```
P(T1,T2,T3)=C1(T1) AND C(T2) AND C3(T3)
```

For that example, the nested-loop algorithm using guarded pushed-down conditions looks like this:

```
FOR each row tl in Tl such that Cl(tl) {
BOOL fl:=FALSE;
```

In general, pushed-down predicates can be extracted from join conditions such as P1(T1,T2) and P(T2,T3). In this case, a pushed-down predicate is guarded also by a flag that prevents checking the predicate for the NULL-complemented row generated by the corresponding outer join operation.

Access by key from one inner table to another in the same nested join is prohibited if it is induced by a predicate from the WHERE condition.

#### 10.2.1.9 Outer Join Optimization

Outer joins include LEFT JOIN and RIGHT JOIN.

MySQL implements an A LEFT JOIN B join\_specification as follows:

- Table B is set to depend on table A and all tables on which A depends.
- Table A is set to depend on all tables (except B) that are used in the LEFT JOIN condition.
- The LEFT JOIN condition is used to decide how to retrieve rows from table B. (In other words, any condition in the WHERE clause is not used.)
- All standard join optimizations are performed, with the exception that a table is always read after all tables on which it depends. If there is a circular dependency, an error occurs.
- All standard WHERE optimizations are performed.
- If there is a row in A that matches the WHERE clause, but there is no row in B that matches the ON condition, an extra B row is generated with all columns set to NULL.
- If you use LEFT JOIN to find rows that do not exist in some table and you have the following test: col\_name IS NULL in the WHERE part, where col\_name is a column that is declared as NOT NULL, MySQL stops searching for more rows (for a particular key combination) after it has found one row that matches the LEFT JOIN condition.

The RIGHT JOIN implementation is analogous to that of LEFT JOIN with the table roles reversed. Right joins are converted to equivalent left joins, as described in Section 10.2.1.10, "Outer Join Simplification".

For a LEFT JOIN, if the WHERE condition is always false for the generated NULL row, the LEFT JOIN is changed to an inner join. For example, the WHERE clause would be false in the following query if t2.column1 were NULL:

```
SELECT * FROM t1 LEFT JOIN t2 ON (column1) WHERE t2.column2=5;
```

Therefore, it is safe to convert the query to an inner join:

```
SELECT * FROM t1, t2 WHERE t2.column2=5 AND t1.column1=t2.column1;
```

Trivial WHERE conditions arising from constant literal expressions are removed during preparation, rather than at a later stage in optimization, by which time joins have already been simplified. Earlier removal of trivial conditions allows the optimizer to convert outer joins to inner joins; this can result in improved plans for queries with outer joins containing trivial conditions in the WHERE clause, such as this one:

```
SELECT * FROM t1 LEFT JOIN t2 ON condition_1 WHERE condition_2 OR 0 = 1
```

The optimizer now sees during preparation that 0 = 1 is always false, making OR = 1 redundant, and removes it, leaving this:

```
SELECT * FROM t1 LEFT JOIN t2 ON condition_1 where condition_2
```

Now the optimizer can rewrite the guery as an inner join, like this:

```
SELECT * FROM t1 JOIN t2 WHERE condition_1 AND condition_2
```

Now the optimizer can use table t2 before table t1 if doing so would result in a better query plan. To provide a hint about the table join order, use optimizer hints; see Section 10.9.3, "Optimizer Hints". Alternatively, use STRAIGHT\_JOIN; see Section 15.2.13, "SELECT Statement". However, STRAIGHT\_JOIN may prevent indexes from being used because it disables semijoin transformations; see Optimizing IN and EXISTS Subguery Predicates with Semijoin Transformations.

#### 10.2.1.10 Outer Join Simplification

Table expressions in the FROM clause of a query are simplified in many cases.

At the parser stage, queries with right outer join operations are converted to equivalent queries containing only left join operations. In the general case, the conversion is performed such that this right join:

```
(T1, ...) RIGHT JOIN (T2, ...) ON P(T1, ..., T2, ...)
```

Becomes this equivalent left join:

```
(T2, ...) LEFT JOIN (T1, ...) ON P(T1, ..., T2, ...)
```

All inner join expressions of the form T1 INNER JOIN T2 ON P(T1,T2) are replaced by the list T1,T2, P(T1,T2) being joined as a conjunct to the WHERE condition (or to the join condition of the embedding join, if there is any).

When the optimizer evaluates plans for outer join operations, it takes into consideration only plans where, for each such operation, the outer tables are accessed before the inner tables. The optimizer choices are limited because only such plans enable outer joins to be executed using the nested-loop algorithm.

Consider a query of this form, where  $\mathbb{R}(\mathbb{T}^2)$  greatly narrows the number of matching rows from table  $\mathbb{T}^2$ :

```
SELECT * T1 FROM T1
LEFT JOIN T2 ON P1(T1,T2)
WHERE P(T1,T2) AND R(T2)
```

If the query is executed as written, the optimizer has no choice but to access the less-restricted table T1 before the more-restricted table T2, which may produce a very inefficient execution plan.

Instead, MySQL converts the query to a query with no outer join operation if the WHERE condition is null-rejected. (That is, it converts the outer join to an inner join.) A condition is said to be null-rejected

for an outer join operation if it evaluates to FALSE or UNKNOWN for any NULL-complemented row generated for the operation.

Thus, for this outer join:

```
T1 LEFT JOIN T2 ON T1.A=T2.A
```

Conditions such as these are null-rejected because they cannot be true for any NULL-complemented row (with T2 columns set to NULL):

```
T2.B IS NOT NULL
T2.B > 3
T2.C <= T1.C
T2.B < 2 OR T2.C > 1
```

Conditions such as these are not null-rejected because they might be true for a NULL-complemented row:

```
T2.B IS NULL
T1.B < 3 OR T2.B IS NOT NULL
T1.B < 3 OR T2.B > 3
```

The general rules for checking whether a condition is null-rejected for an outer join operation are simple:

- It is of the form A IS NOT NULL, where A is an attribute of any of the inner tables
- It is a predicate containing a reference to an inner table that evaluates to UNKNOWN when one of its arguments is NULL
- · It is a conjunction containing a null-rejected condition as a conjunct
- It is a disjunction of null-rejected conditions

A condition can be null-rejected for one outer join operation in a query and not null-rejected for another. In this query, the WHERE condition is null-rejected for the second outer join operation but is not null-rejected for the first one:

```
SELECT * FROM T1 LEFT JOIN T2 ON T2.A=T1.A

LEFT JOIN T3 ON T3.B=T1.B

WHERE T3.C > 0
```

If the WHERE condition is null-rejected for an outer join operation in a query, the outer join operation is replaced by an inner join operation.

For example, in the preceding query, the second outer join is null-rejected and can be replaced by an inner join:

```
SELECT * FROM T1 LEFT JOIN T2 ON T2.A=T1.A

INNER JOIN T3 ON T3.B=T1.B

WHERE T3.C > 0
```

For the original query, the optimizer evaluates only plans compatible with the single table-access order  $\mathtt{T1}$ ,  $\mathtt{T2}$ ,  $\mathtt{T3}$ . For the rewritten query, it additionally considers the access order  $\mathtt{T3}$ ,  $\mathtt{T1}$ ,  $\mathtt{T2}$ .

A conversion of one outer join operation may trigger a conversion of another. Thus, the query:

```
SELECT * FROM T1 LEFT JOIN T2 ON T2.A=T1.A

LEFT JOIN T3 ON T3.B=T2.B

WHERE T3.C > 0
```

Is first converted to the query:

```
SELECT * FROM T1 LEFT JOIN T2 ON T2.A=T1.A
INNER JOIN T3 ON T3.B=T2.B
```

```
WHERE T3.C > 0
```

Which is equivalent to the query:

```
SELECT * FROM (T1 LEFT JOIN T2 ON T2.A=T1.A), T3
WHERE T3.C > 0 AND T3.B=T2.B
```

The remaining outer join operation can also be replaced by an inner join because the condition T3.B=T2.B is null-rejected. This results in a query with no outer joins at all:

```
SELECT * FROM (T1 INNER JOIN T2 ON T2.A=T1.A), T3
WHERE T3.C > 0 AND T3.B=T2.B
```

Sometimes the optimizer succeeds in replacing an embedded outer join operation, but cannot convert the embedding outer join. The following query:

```
SELECT * FROM T1 LEFT JOIN

(T2 LEFT JOIN T3 ON T3.B=T2.B)

ON T2.A=T1.A

WHERE T3.C > 0
```

Is converted to:

```
SELECT * FROM T1 LEFT JOIN

(T2 INNER JOIN T3 ON T3.B=T2.B)

ON T2.A=T1.A

WHERE T3.C > 0
```

That can be rewritten only to the form still containing the embedding outer join operation:

```
SELECT * FROM T1 LEFT JOIN
(T2,T3)
ON (T2.A=T1.A AND T3.B=T2.B)
WHERE T3.C > 0
```

Any attempt to convert an embedded outer join operation in a query must take into account the join condition for the embedding outer join together with the WHERE condition. In this query, the WHERE condition is not null-rejected for the embedded outer join, but the join condition of the embedding outer join T2.A=T1.A AND T3.C=T1.C is null-rejected:

```
SELECT * FROM T1 LEFT JOIN

(T2 LEFT JOIN T3 ON T3.B=T2.B)

ON T2.A=T1.A AND T3.C=T1.C

WHERE T3.D > 0 OR T1.D > 0
```

Consequently, the query can be converted to:

```
SELECT * FROM T1 LEFT JOIN

(T2, T3)

ON T2.A=T1.A AND T3.C=T1.C AND T3.B=T2.B

WHERE T3.D > 0 OR T1.D > 0
```

#### 10.2.1.11 Multi-Range Read Optimization

Reading rows using a range scan on a secondary index can result in many random disk accesses to the base table when the table is large and not stored in the storage engine's cache. With the Disk-Sweep Multi-Range Read (MRR) optimization, MySQL tries to reduce the number of random disk access for range scans by first scanning the index only and collecting the keys for the relevant rows. Then the keys are sorted and finally the rows are retrieved from the base table using the order of the primary key. The motivation for Disk-sweep MRR is to reduce the number of random disk accesses and instead achieve a more sequential scan of the base table data.

The Multi-Range Read optimization provides these benefits:

 MRR enables data rows to be accessed sequentially rather than in random order, based on index tuples. The server obtains a set of index tuples that satisfy the query conditions, sorts them according to data row ID order, and uses the sorted tuples to retrieve data rows in order. This makes data access more efficient and less expensive.

MRR enables batch processing of requests for key access for operations that require access to data
rows through index tuples, such as range index scans and equi-joins that use an index for the join
attribute. MRR iterates over a sequence of index ranges to obtain qualifying index tuples. As these
results accumulate, they are used to access the corresponding data rows. It is not necessary to
acquire all index tuples before starting to read data rows.

The MRR optimization is not supported with secondary indexes created on virtual generated columns. InnobB supports secondary indexes on virtual generated columns.

The following scenarios illustrate when MRR optimization can be advantageous:

Scenario A: MRR can be used for InnoDB and MyISAM tables for index range scans and equi-join operations.

- 1. A portion of the index tuples are accumulated in a buffer.
- 2. The tuples in the buffer are sorted by their data row ID.
- 3. Data rows are accessed according to the sorted index tuple sequence.

Scenario B: MRR can be used for NDB tables for multiple-range index scans or when performing an equi-join by an attribute.

- 1. A portion of ranges, possibly single-key ranges, is accumulated in a buffer on the central node where the query is submitted.
- 2. The ranges are sent to the execution nodes that access data rows.
- 3. The accessed rows are packed into packages and sent back to the central node.
- 4. The received packages with data rows are placed in a buffer.
- 5. Data rows are read from the buffer.

When MRR is used, the Extra column in EXPLAIN output shows Using MRR.

InnoDB and MyISAM do not use MRR if full table rows need not be accessed to produce the query result. This is the case if results can be produced entirely on the basis on information in the index tuples (through a covering index); MRR provides no benefit.

Two optimizer\_switch system variable flags provide an interface to the use of MRR optimization. The mrr flag controls whether MRR is enabled. If mrr is enabled (on), the mrr\_cost\_based flag controls whether the optimizer attempts to make a cost-based choice between using and not using MRR (on) or uses MRR whenever possible (off). By default, mrr is on and mrr\_cost\_based is on. See Section 10.9.2, "Switchable Optimizations".

For MRR, a storage engine uses the value of the read\_rnd\_buffer\_size system variable
as a guideline for how much memory it can allocate for its buffer. The engine uses up to
read\_rnd\_buffer\_size bytes and determines the number of ranges to process in a single pass.

#### 10.2.1.12 Block Nested-Loop and Batched Key Access Joins

In MySQL, a Batched Key Access (BKA) Join algorithm is available that uses both index access to the joined table and a join buffer. The BKA algorithm supports inner join, outer join, and semijoin operations, including nested outer joins. Benefits of BKA include improved join performance due to more efficient table scanning. Also, the Block Nested-Loop (BNL) Join algorithm previously used only for inner joins is extended and can be employed for outer join and semijoin operations, including nested outer joins.

The following sections discuss the join buffer management that underlies the extension of the original BNL algorithm, the extended BNL algorithm, and the BKA algorithm. For information about semijoin strategies, see Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations

- Join Buffer Management for Block Nested-Loop and Batched Key Access Algorithms
- Block Nested-Loop Algorithm for Outer Joins and Semijoins
- · Batched Key Access Joins
- · Optimizer Hints for Block Nested-Loop and Batched Key Access Algorithms

#### Join Buffer Management for Block Nested-Loop and Batched Key Access Algorithms

MySQL can employ join buffers to execute not only inner joins without index access to the inner table, but also outer joins and semijoins that appear after subquery flattening. Moreover, a join buffer can be effectively used when there is an index access to the inner table.

The join buffer management code slightly more efficiently utilizes join buffer space when storing the values of the interesting row columns: No additional bytes are allocated in buffers for a row column if its value is NULL, and the minimum number of bytes is allocated for any value of the VARCHAR type.

The code supports two types of buffers, regular and incremental. Suppose that join buffer B1 is employed to join tables t1 and t2 and the result of this operation is joined with table t3 using join buffer B2:

- A regular join buffer contains columns from each join operand. If B2 is a regular join buffer, each
  row r put into B2 is composed of the columns of a row r1 from B1 and the interesting columns of a
  matching row r2 from table t3.
- An incremental join buffer contains only columns from rows of the table produced by the second join operand. That is, it is incremental to a row from the first operand buffer. If B2 is an incremental join buffer, it contains the interesting columns of the row *r2* together with a link to the row *r1* from B1.

Incremental join buffers are always incremental relative to a join buffer from an earlier join operation, so the buffer from the first join operation is always a regular buffer. In the example just given, the buffer B1 used to join tables t1 and t2 must be a regular buffer.

Each row of the incremental buffer used for a join operation contains only the interesting columns of a row from the table to be joined. These columns are augmented with a reference to the interesting columns of the matched row from the table produced by the first join operand. Several rows in the incremental buffer can refer to the same row r whose columns are stored in the previous join buffers insofar as all these rows match row r.

Incremental buffers enable less frequent copying of columns from buffers used for previous join operations. This provides a savings in buffer space because in the general case a row produced by the first join operand can be matched by several rows produced by the second join operand. It is unnecessary to make several copies of a row from the first operand. Incremental buffers also provide a savings in processing time due to the reduction in copying time.

The block\_nested\_loop flag of the optimizer\_switch system variable controls hash joins.

The batched\_key\_access flag controls how the optimizer uses the Batched Key Access join algorithms.

By default, block\_nested\_loop is on and batched\_key\_access is off. See Section 10.9.2, "Switchable Optimizations". Optimizer hints may also be applied; see Optimizer Hints for Block Nested-Loop and Batched Key Access Algorithms.

For information about semijoin strategies, see Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations

#### **Block Nested-Loop Algorithm for Outer Joins and Semijoins**

The original implementation of the MySQL BNL algorithm was extended to support outer join and semijoin operations (and was later superseded by the hash join algorithm; see Section 10.2.1.4, "Hash Join Optimization").

When these operations are executed with a join buffer, each row put into the buffer is supplied with a match flag.

If an outer join operation is executed using a join buffer, each row of the table produced by the second operand is checked for a match against each row in the join buffer. When a match is found, a new extended row is formed (the original row plus columns from the second operand) and sent for further extensions by the remaining join operations. In addition, the match flag of the matched row in the buffer is enabled. After all rows of the table to be joined have been examined, the join buffer is scanned. Each row from the buffer that does not have its match flag enabled is extended by NULL complements (NULL values for each column in the second operand) and sent for further extensions by the remaining join operations.

The block\_nested\_loop flag of the optimizer\_switch system variable controls hash joins.

See Section 10.9.2, "Switchable Optimizations", for more information. Optimizer hints may also be applied; see Optimizer Hints for Block Nested-Loop and Batched Key Access Algorithms.

In EXPLAIN output, use of BNL for a table is signified when the Extra value contains Using join buffer (Block Nested Loop) and the type value is ALL, index, or range.

For information about semijoin strategies, see Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations

#### **Batched Key Access Joins**

MySQL implements a method of joining tables called the Batched Key Access (BKA) join algorithm. BKA can be applied when there is an index access to the table produced by the second join operand. Like the BNL join algorithm, the BKA join algorithm employs a join buffer to accumulate the interesting columns of the rows produced by the first operand of the join operation. Then the BKA algorithm builds keys to access the table to be joined for all rows in the buffer and submits these keys in a batch to the database engine for index lookups. The keys are submitted to the engine through the Multi-Range Read (MRR) interface (see Section 10.2.1.11, "Multi-Range Read Optimization"). After submission of the keys, the MRR engine functions perform lookups in the index in an optimal way, fetching the rows of the joined table found by these keys, and starts feeding the BKA join algorithm with matching rows. Each matching row is coupled with a reference to a row in the join buffer.

When BKA is used, the value of <code>join\_buffer\_size</code> defines how large the batch of keys is in each request to the storage engine. The larger the buffer, the more sequential access is made to the right hand table of a join operation, which can significantly improve performance.

For BKA to be used, the batched\_key\_access flag of the optimizer\_switch system variable must be set to on. BKA uses MRR, so the mrr flag must also be on. Currently, the cost estimation for MRR is too pessimistic. Hence, it is also necessary for mrr\_cost\_based to be off for BKA to be used. The following setting enables BKA:

mysql> SET optimizer\_switch='mrr=on,mrr\_cost\_based=off,batched\_key\_access=on';

There are two scenarios by which MRR functions execute:

• The first scenario is used for conventional disk-based storage engines such as InnoDB and MyISAM. For these engines, usually the keys for all rows from the join buffer are submitted to the MRR interface at once. Engine-specific MRR functions perform index lookups for the submitted keys, get row IDs (or primary keys) from them, and then fetch rows for all these selected row IDs one by one by request from BKA algorithm. Every row is returned with an association reference that enables access to the matched row in the join buffer. The rows are fetched by the MRR functions in an optimal way: They are fetched in the row ID (primary key) order. This improves performance because reads are in disk order rather than random order.

• The second scenario is used for remote storage engines such as NDB. A package of keys for a portion of rows from the join buffer, together with their associations, is sent by a MySQL Server (SQL node) to MySQL Cluster data nodes. In return, the SQL node receives a package (or several packages) of matching rows coupled with corresponding associations. The BKA join algorithm takes these rows and builds new joined rows. Then a new set of keys is sent to the data nodes and the rows from the returned packages are used to build new joined rows. The process continues until the last keys from the join buffer are sent to the data nodes, and the SQL node has received and joined all rows matching these keys. This improves performance because fewer key-bearing packages sent by the SQL node to the data nodes means fewer round trips between it and the data nodes to perform the join operation.

With the first scenario, a portion of the join buffer is reserved to store row IDs (primary keys) selected by index lookups and passed as a parameter to the MRR functions.

There is no special buffer to store keys built for rows from the join buffer. Instead, a function that builds the key for the next row in the buffer is passed as a parameter to the MRR functions.

In EXPLAIN output, use of BKA for a table is signified when the Extra value contains Using join buffer (Batched Key Access) and the type value is ref or eq\_ref.

#### Optimizer Hints for Block Nested-Loop and Batched Key Access Algorithms

In addition to using the <code>optimizer\_switch</code> system variable to control optimizer use of the BNL and BKA algorithms session-wide, MySQL supports optimizer hints to influence the optimizer on a perstatement basis. See Section 10.9.3, "Optimizer Hints".

To use a BNL or BKA hint to enable join buffering for any inner table of an outer join, join buffering must be enabled for all inner tables of the outer join.

#### 10.2.1.13 Condition Filtering

In join processing, prefix rows are those rows passed from one table in a join to the next. In general, the optimizer attempts to put tables with low prefix counts early in the join order to keep the number of row combinations from increasing rapidly. To the extent that the optimizer can use information about conditions on rows selected from one table and passed to the next, the more accurately it can compute row estimates and choose the best execution plan.

Without condition filtering, the prefix row count for a table is based on the estimated number of rows selected by the WHERE clause according to whichever access method the optimizer chooses. Condition filtering enables the optimizer to use other relevant conditions in the WHERE clause not taken into account by the access method, and thus improve its prefix row count estimates. For example, even though there might be an index-based access method that can be used to select rows from the current table in a join, there might also be additional conditions for the table in the WHERE clause that can filter (further restrict) the estimate for qualifying rows passed to the next table.

A condition contributes to the filtering estimate only if:

- It refers to the current table.
- It depends on a constant value or values from earlier tables in the join sequence.
- It was not already taken into account by the access method.

In EXPLAIN output, the rows column indicates the row estimate for the chosen access method, and the filtered column reflects the effect of condition filtering. filtered values are expressed as percentages. The maximum value is 100, which means no filtering of rows occurred. Values decreasing from 100 indicate increasing amounts of filtering.

The prefix row count (the number of rows estimated to be passed from the current table in a join to the next) is the product of the rows and filtered values. That is, the prefix row count is the estimated row count, reduced by the estimated filtering effect. For example, if rows is 1000 and filtered is

20%, condition filtering reduces the estimated row count of 1000 to a prefix row count of  $1000 \times 20\% = 1000 \times .2 = 200$ .

Consider the following query:

```
SELECT *

FROM employee JOIN department ON employee.dept_no = department.dept_no
WHERE employee.first_name = 'John'
AND employee.hire_date BETWEEN '2018-01-01' AND '2018-06-01';
```

Suppose that the data set has these characteristics:

- The employee table has 1024 rows.
- The department table has 12 rows.
- Both tables have an index on dept\_no.
- The employee table has an index on first\_name.
- 8 rows satisfy this condition on employee.first\_name:

```
employee.first_name = 'John'
```

• 150 rows satisfy this condition on employee.hire date:

```
employee.hire_date BETWEEN '2018-01-01' AND '2018-06-01'
```

1 row satisfies both conditions:

```
employee.first_name = 'John'
AND employee.hire_date BETWEEN '2018-01-01' AND '2018-06-01'
```

Without condition filtering, EXPLAIN produces output like this:

id   table	type	+   possible_keys	+   key	+   ref +	+   rows	++   filtered
1		name,h_date,dept   PRIMARY	name   PRIMARY	const dept_no	8   1	100.00     100.00

For employee, the access method on the name index picks up the 8 rows that match a name of 'John'. No filtering is done (filtered is 100%), so all rows are prefix rows for the next table: The prefix row count is rows  $\times$  filtered =  $8 \times 100\% = 8$ .

With condition filtering, the optimizer additionally takes into account conditions from the WHERE clause not taken into account by the access method. In this case, the optimizer uses heuristics to estimate a filtering effect of 16.31% for the BETWEEN condition on employee.hire\_date. As a result, EXPLAIN produces output like this:

++		possible_keys	+   key	+   ref	+   rows	++   filtered
	ref	name,h_date,dept	•	const dept_no	8   1	16.31     100.00

Now the prefix row count is  $rows \times filtered = 8 \times 16.31\% = 1.3$ , which more closely reflects actual data set.

Normally, the optimizer does not calculate the condition filtering effect (prefix row count reduction) for the last joined table because there is no next table to pass rows to. An exception occurs for EXPLAIN: To provide more information, the filtering effect is calculated for all joined tables, including the last one.

To control whether the optimizer considers additional filtering conditions, use the condition\_fanout\_filter flag of the optimizer\_switch system variable (see Section 10.9.2,

"Switchable Optimizations"). This flag is enabled by default but can be disabled to suppress condition filtering (for example, if a particular query is found to yield better performance without it).

If the optimizer overestimates the effect of condition filtering, performance may be worse than if condition filtering is not used. In such cases, these techniques may help:

- If a column is not indexed, index it so that the optimizer has some information about the distribution of column values and can improve its row estimates.
- Similarly, if no column histogram information is available, generate a histogram (see Section 10.9.6, "Optimizer Statistics").
- Change the join order. Ways to accomplish this include join-order optimizer hints (see Section 10.9.3, "Optimizer Hints"), STRAIGHT\_JOIN immediately following the SELECT, and the STRAIGHT\_JOIN join operator.
- Disable condition filtering for the session:

```
SET optimizer_switch = 'condition_fanout_filter=off';
```

Or, for a given query, using an optimizer hint:

```
SELECT /*+ SET_VAR(optimizer_switch = 'condition_fanout_filter=off') */ ...
```

## **10.2.1.14 Constant-Folding Optimization**

Comparisons between constants and column values in which the constant value is out of range or of the wrong type with respect to the column type are now handled once during query optimization rather row-by-row than during execution. The comparisons that can be treated in this manner are >, >=, <, <=, <>/!=, =, and <=>.

Consider the table created by the following statement:

```
CREATE TABLE t (c TINYINT UNSIGNED NOT NULL);
```

The where condition in the query Select \* From t where c < 256 contains the integral constant 256 which is out of range for a TINYINT UNSIGNED column. Previously, this was handled by treating both operands as the larger type, but now, since any allowed value for c is less than the constant, the where expression can instead be folded as where 1, so that the query is rewritten as SELECT \* FROM t where 1.

This makes it possible for the optimizer to remove the WHERE expression altogether. If the column c were nullable (that is, defined only as TINYINT UNSIGNED) the query would be rewritten like this:

```
SELECT * FROM t WHERE ti IS NOT NULL
```

Folding is performed for constants compared to supported MySQL column types as follows:

- **Integer column type.** Integer types are compared with constants of the following types as described here:
  - Integer value. If the constant is out of range for the column type, the comparison is folded to 1
    or IS NOT NULL, as already shown.

If the constant is a range boundary, the comparison is folded to =. For example (using the same table as already defined):

```
mysql> EXPLAIN SELECT * FROM t WHERE c >= 255;

***********
        id: 1
        select_type: SIMPLE
            table: t
        partitions: NULL
            type: ALL
possible_keys: NULL
```

• Floating- or fixed-point value. If the constant is one of the decimal types (such as DECIMAL, REAL, DOUBLE, or FLOAT) and has a nonzero decimal portion, it cannot be equal; fold accordingly. For other comparisons, round up or down to an integer value according to the sign, then perform a range check and handle as already described for integer-integer comparisons.

A REAL value that is too small to be represented as DECIMAL is rounded to .01 or -.01 depending on the sign, then handled as a DECIMAL.

- **String types.** Try to interpret the string value as an integer type, then handle the comparison as between integer values. If this fails, attempt to handle the value as a REAL.
- **DECIMAL or REAL column.** Decimal types are compared with constants of the following types as described here:
  - **Integer value.** Perform a range check against the column value's integer part. If no folding results, convert the constant to DECIMAL with the same number of decimal places as the column value, then check it as a DECIMAL (see next).
  - **DECIMAL or REAL value.** Check for overflow (that is, whether the constant has more digits in its integer part than allowed for the column's decimal type). If so, fold.

If the constant has more significant fractional digits than column's type, truncate the constant. If the comparison operator is = or <>, fold. If the operator is >= or <=, adjust the operator due to truncation. For example, if column's type is DECIMAL(3,1), SELECT \* FROM t WHERE f >= 10.13 becomes SELECT \* FROM t WHERE f > 10.1.

If the constant has fewer decimal digits than the column's type, convert it to a constant with same number of digits. For underflow of a REAL value (that is, too few fractional digits to represent it), convert the constant to decimal 0.

- **String value.** If the value can be interpreted as an integer type, handle it as such. Otherwise, try to handle it as REAL.
- **FLOAT or DOUBLE column.** FLOAT (m,n) or DOUBLE (m,n) values compared with constants are handled as follows:

If the value overflows the range of the column, fold.

If the value has more than n decimals, truncate, compensating during folding. For = and <> comparisons, fold to TRUE, FALSE, or IS <code>[NOT]</code> NULL as described previously; for other operators, adjust the operator.

If the value has more than m integer digits, fold.

**Limitations.** This optimization cannot be used in the following cases:

- 1. With comparisons using BETWEEN or IN.
- 2. With BIT columns or columns using date or time types.

3. During the preparation phase for a prepared statement, although it can be applied during the optimization phase when the prepared statement is actually executed. This due to the fact that, during statement preparation, the value of the constant is not yet known.

## 10.2.1.15 IS NULL Optimization

MySQL can perform the same optimization on <code>col\_name</code> IS <code>NULL</code> that it can use for <code>col\_name = constant\_value</code>. For example, MySQL can use indexes and ranges to search for <code>NULL</code> with <code>IS NULL</code>.

### Examples:

```
SELECT * FROM tbl_name WHERE key_col IS NULL;

SELECT * FROM tbl_name WHERE key_col <=> NULL;

SELECT * FROM tbl_name
WHERE key_col=const1 OR key_col=const2 OR key_col IS NULL;
```

If a WHERE clause includes a col\_name IS NULL condition for a column that is declared as NOT NULL, that expression is optimized away. This optimization does not occur in cases when the column might produce NULL anyway (for example, if it comes from a table on the right side of a LEFT JOIN).

MySQL can also optimize the combination  $col\_name = expr$  OR  $col\_name$  IS NULL, a form that is common in resolved subqueries. EXPLAIN shows ref\_or\_null when this optimization is used.

This optimization can handle one IS NULL for any key part.

Some examples of queries that are optimized, assuming that there is an index on columns a and b of table t2:

```
SELECT * FROM t1 WHERE t1.a=expr OR t1.a IS NULL;

SELECT * FROM t1, t2 WHERE t1.a=t2.a OR t2.a IS NULL;

SELECT * FROM t1, t2
WHERE (t1.a=t2.a OR t2.a IS NULL) AND t2.b=t1.b;

SELECT * FROM t1, t2
WHERE t1.a=t2.a AND (t2.b=t1.b OR t2.b IS NULL);

SELECT * FROM t1, t2
WHERE (t1.a=t2.a AND t2.a IS NULL AND ...)
OR (t1.a=t2.a AND t2.a IS NULL AND ...);
```

ref\_or\_null works by first doing a read on the reference key, and then a separate search for rows with a NULL key value.

The optimization can handle only one IS NULL level. In the following query, MySQL uses key lookups only on the expression (t1.a=t2.a AND t2.a IS NULL) and is not able to use the key part on b:

```
SELECT * FROM t1, t2
WHERE (t1.a=t2.a AND t2.a IS NULL)
OR (t1.b=t2.b AND t2.b IS NULL);
```

## 10.2.1.16 ORDER BY Optimization

This section describes when MySQL can use an index to satisfy an ORDER BY clause, the filesort operation used when an index cannot be used, and execution plan information available from the optimizer about ORDER BY.

An ORDER BY with and without LIMIT may return rows in different orders, as discussed in Section 10.2.1.19, "LIMIT Query Optimization".

Use of Indexes to Satisfy ORDER BY

- · Use of filesort to Satisfy ORDER BY
- · Influencing ORDER BY Optimization
- ORDER BY Execution Plan Information Available

### Use of Indexes to Satisfy ORDER BY

In some cases, MySQL may use an index to satisfy an ORDER BY clause and avoid the extra sorting involved in performing a filesort operation.

The index may also be used even if the ORDER BY does not match the index exactly, as long as all unused portions of the index and all extra ORDER BY columns are constants in the WHERE clause. If the index does not contain all columns accessed by the query, the index is used only if index access is cheaper than other access methods.

Assuming that there is an index on (<code>key\_part1</code>, <code>key\_part2</code>), the following queries may use the index to resolve the <code>ORDER BY</code> part. Whether the optimizer actually does so depends on whether reading the index is more efficient than a table scan if columns not in the index must also be read.

• In this query, the index on (key\_part1, key\_part2) enables the optimizer to avoid sorting:

```
SELECT * FROM t1
ORDER BY key_part1, key_part2;
```

However, the query uses SELECT \*, which may select more columns than  $key\_part1$  and  $key\_part2$ . In that case, scanning an entire index and looking up table rows to find columns not in the index may be more expensive than scanning the table and sorting the results. If so, the optimizer probably does not use the index. If SELECT \* selects only the index columns, the index is used and sorting avoided.

If t1 is an InnoDB table, the table primary key is implicitly part of the index, and the index can be used to resolve the ORDER BY for this query:

```
SELECT pk, key_part1, key_part2 FROM t1
ORDER BY key_part1, key_part2;
```

• In this query,  $key\_part1$  is constant, so all rows accessed through the index are in  $key\_part2$  order, and an index on  $(key\_part1, key\_part2)$  avoids sorting if the WHERE clause is selective enough to make an index range scan cheaper than a table scan:

```
SELECT * FROM t1
WHERE key_part1 = constant
ORDER BY key_part2;
```

• In the next two queries, whether the index is used is similar to the same queries without DESC shown previously:

```
SELECT * FROM t1
  ORDER BY key_part1 DESC, key_part2 DESC;

SELECT * FROM t1
  WHERE key_part1 = constant
  ORDER BY key_part2 DESC;
```

• Two columns in an ORDER BY can sort in the same direction (both ASC, or both DESC) or in opposite directions (one ASC, one DESC). A condition for index use is that the index must have the same homogeneity, but need not have the same actual direction.

If a query mixes ASC and DESC, the optimizer can use an index on the columns if the index also uses corresponding mixed ascending and descending columns:

```
SELECT * FROM t1
ORDER BY key_part1 DESC, key_part2 ASC;
```

The optimizer can use an index on (key\_part1, key\_part2) if key\_part1 is descending and key\_part2 is ascending. It can also use an index on those columns (with a backward scan) if key\_part1 is ascending and key\_part2 is descending. See Section 10.3.13, "Descending Indexes".

• In the next two queries, <a href="key\_part1">key\_part1</a> is compared to a constant. The index is used if the <a href="https://www.next.numer.com/where.c

```
SELECT * FROM t1

WHERE key_part1 > constant

ORDER BY key_part1 ASC;

SELECT * FROM t1

WHERE key_part1 < constant

ORDER BY key_part1 DESC;
```

• In the next query, the ORDER BY does not name  $key\_part1$ , but all rows selected have a constant  $key\_part1$  value, so the index can still be used:

```
SELECT * FROM t1
WHERE key_part1 = constant1 AND key_part2 > constant2
ORDER BY key_part2;
```

In some cases, MySQL *cannot* use indexes to resolve the ORDER BY, although it may still use indexes to find the rows that match the WHERE clause. Examples:

The query uses ORDER BY on different indexes:

```
SELECT * FROM t1 ORDER BY key1, key2;
```

• The query uses ORDER BY on nonconsecutive parts of an index:

```
SELECT * FROM t1 WHERE key2=constant ORDER BY key1_part1, key1_part3;
```

The index used to fetch the rows differs from the one used in the ORDER BY:

```
SELECT * FROM t1 WHERE key2=constant ORDER BY key1;
```

 The query uses ORDER BY with an expression that includes terms other than the index column name:

```
SELECT * FROM t1 ORDER BY ABS(key);
SELECT * FROM t1 ORDER BY -key;
```

- The query joins many tables, and the columns in the ORDER BY are not all from the first nonconstant table that is used to retrieve rows. (This is the first table in the EXPLAIN output that does not have a const join type.)
- The query has different ORDER BY and GROUP BY expressions.
- There is an index on only a prefix of a column named in the ORDER BY clause. In this case, the index
  cannot be used to fully resolve the sort order. For example, if only the first 10 bytes of a CHAR(20)
  column are indexed, the index cannot distinguish values past the 10th byte and a filesort is
  needed.
- The index does not store rows in order. For example, this is true for a HASH index in a MEMORY table.

Availability of an index for sorting may be affected by the use of column aliases. Suppose that the column t1.a is indexed. In this statement, the name of the column in the select list is a. It refers to t1.a, as does the reference to a in the ORDER BY, so the index on t1.a can be used:

```
SELECT a FROM t1 ORDER BY a;
```

In this statement, the name of the column in the select list is also a, but it is the alias name. It refers to ABS(a), as does the reference to a in the ORDER BY, so the index on t1.a cannot be used:

```
SELECT ABS(a) AS a FROM t1 ORDER BY a;
```

In the following statement, the ORDER BY refers to a name that is not the name of a column in the select list. But there is a column in t1 named a, so the ORDER BY refers to t1.a and the index on t1.a can be used. (The resulting sort order may be completely different from the order for ABS(a), of course.)

```
SELECT ABS(a) AS b FROM t1 ORDER BY a;
```

Previously (MySQL 8.3 and lower), GROUP BY sorted implicitly under certain conditions. In MySQL 8.4, that no longer occurs, so specifying ORDER BY NULL at the end to suppress implicit sorting (as was done previously) is no longer necessary. However, query results may differ from previous MySQL versions. To produce a given sort order, provide an ORDER BY clause.

## Use of filesort to Satisfy ORDER BY

If an index cannot be used to satisfy an ORDER BY clause, MySQL performs a filesort operation that reads table rows and sorts them. A filesort constitutes an extra sorting phase in query execution.

To obtain memory for filesort operations, the optimizer allocates memory buffers incrementally as needed, up to the size indicated by the sort\_buffer\_size system variable. This enables users to set sort\_buffer\_size to larger values to speed up larger sorts, without concern for excessive memory use for small sorts. (This benefit may not occur for multiple concurrent sorts on Windows, which has a weak multithreaded malloc.)

A filesort operation uses temporary disk files as necessary if the result set is too large to fit in memory. Some types of queries are particularly suited to completely in-memory filesort operations. For example, the optimizer can use filesort to efficiently handle in memory, without temporary files, the ORDER BY operation for queries (and subqueries) of the following form:

```
SELECT ... FROM single_table ... ORDER BY non_index_column [DESC] LIMIT [M,]N;
```

Such queries are common in web applications that display only a few rows from a larger result set. Examples:

```
SELECT col1, ... FROM t1 ... ORDER BY name LIMIT 10;
SELECT col1, ... FROM t1 ... ORDER BY RAND() LIMIT 15;
```

### Influencing ORDER BY Optimization

To increase ORDER BY speed, check whether you can get MySQL to use indexes rather than an extra sorting phase. If this is not possible, try the following strategies:

• Increase the sort\_buffer\_size variable value. Ideally, the value should be large enough for the entire result set to fit in the sort buffer (to avoid writes to disk and merge passes).

Take into account that the size of column values stored in the sort buffer is affected by the max\_sort\_length system variable value. For example, if tuples store values of long string columns and you increase the value of max\_sort\_length, the size of sort buffer tuples increases as well and may require you to increase sort\_buffer\_size.

To monitor the number of merge passes (to merge temporary files), check the Sort\_merge\_passes status variable.

- Increase the read\_rnd\_buffer\_size variable value so that more rows are read at a time.
- Change the tmpdir system variable to point to a dedicated file system with large amounts of free space. The variable value can list several paths that are used in round-robin fashion; you can use this feature to spread the load across several directories. Separate the paths by colon characters (:) on Unix and semicolon characters (;) on Windows. The paths should name directories in file systems located on different *physical* disks, not different partitions on the same disk.

### **ORDER BY Execution Plan Information Available**

With EXPLAIN (see Section 10.8.1, "Optimizing Queries with EXPLAIN"), you can check whether MySQL can use indexes to resolve an ORDER BY clause:

- If the Extra column of EXPLAIN output does not contain Using filesort, the index is used and a filesort is not performed.
- If the Extra column of EXPLAIN output contains Using filesort, the index is not used and a filesort is performed.

In addition, if a filesort is performed, optimizer trace output includes a filesort\_summary block. For example:

```
"filesort_summary": {
   "rows": 100,
   "examined_rows": 100,
   "number_of_tmp_files": 0,
   "peak_memory_used": 25192,
   "sort_mode": "<sort_key, packed_additional_fields>"
}
```

peak\_memory\_used indicates the maximum memory used at any one time during the sort. This is a value up to but not necessarily as large as the value of the sort\_buffer\_size system variable. The optimizer allocates sort-buffer memory incrementally, beginning with a small amount and adding more as necessary, up to sort\_buffer\_size bytes.)

The sort mode value provides information about the contents of tuples in the sort buffer:

- <sort\_key, rowid>: This indicates that sort buffer tuples are pairs that contain the sort key value
  and row ID of the original table row. Tuples are sorted by sort key value and the row ID is used to
  read the row from the table.
- <sort\_key, additional\_fields>: This indicates that sort buffer tuples contain the sort key
  value and columns referenced by the query. Tuples are sorted by sort key value and column values
  are read directly from the tuple.
- <sort\_key, packed\_additional\_fields>: Like the previous variant, but the additional columns are packed tightly together instead of using a fixed-length encoding.

EXPLAIN does not distinguish whether the optimizer does or does not perform a filesort in memory. Use of an in-memory filesort can be seen in optimizer trace output. Look for filesort\_priority\_queue\_optimization. For information about the optimizer trace, see Section 10.15, "Tracing the Optimizer".

### 10.2.1.17 GROUP BY Optimization

The most general way to satisfy a GROUP BY clause is to scan the whole table and create a new temporary table where all rows from each group are consecutive, and then use this temporary table to discover groups and apply aggregate functions (if any). In some cases, MySQL is able to do much better than that and avoid creation of temporary tables by using index access.

The most important preconditions for using indexes for GROUP BY are that all GROUP BY columns reference attributes from the same index, and that the index stores its keys in order (as is true, for example, for a BTREE index, but not for a HASH index). Whether use of temporary tables can be replaced by index access also depends on which parts of an index are used in a query, the conditions specified for these parts, and the selected aggregate functions.

There are two ways to execute a GROUP BY query through index access, as detailed in the following sections. The first method applies the grouping operation together with all range predicates (if any). The second method first performs a range scan, and then groups the resulting tuples.

- · Loose Index Scan
- Tight Index Scan

Loose Index Scan can also be used in the absence of GROUP BY under some conditions. See Skip Scan Range Access Method.

### **Loose Index Scan**

The most efficient way to process GROUP BY is when an index is used to directly retrieve the grouping columns. With this access method, MySQL uses the property of some index types that the keys are ordered (for example, BTREE). This property enables use of lookup groups in an index without having to consider all keys in the index that satisfy all WHERE conditions. This access method considers only a fraction of the keys in an index, so it is called a *Loose Index Scan*. When there is no WHERE clause, a Loose Index Scan reads as many keys as the number of groups, which may be a much smaller number than that of all keys. If the WHERE clause contains range predicates (see the discussion of the range join type in Section 10.8.1, "Optimizing Queries with EXPLAIN"), a Loose Index Scan looks up the first key of each group that satisfies the range conditions, and again reads the smallest possible number of keys. This is possible under the following conditions:

- The query is over a single table.
- The GROUP BY names only columns that form a leftmost prefix of the index and no other columns. (If, instead of GROUP BY, the query has a DISTINCT clause, all distinct attributes refer to columns that form a leftmost prefix of the index.) For example, if a table t1 has an index on (c1,c2,c3), Loose Index Scan is applicable if the query has GROUP BY c1, c2. It is not applicable if the query has GROUP BY c2, c3 (the columns are not a leftmost prefix) or GROUP BY c1, c2, c4 (c4 is not in the index).
- The only aggregate functions used in the select list (if any) are MIN() and MAX(), and all of them refer to the same column. The column must be in the index and must immediately follow the columns in the GROUP BY.
- Any other parts of the index than those from the GROUP BY referenced in the query must be constants (that is, they must be referenced in equalities with constants), except for the argument of MIN() or MAX() functions.
- For columns in the index, full column values must be indexed, not just a prefix. For example, with c1 VARCHAR(20), INDEX (c1(10)), the index uses only a prefix of c1 values and cannot be used for Loose Index Scan.

If Loose Index Scan is applicable to a query, the EXPLAIN output shows Using index for group-by in the Extra column.

Assume that there is an index idx(c1,c2,c3) on table t1(c1,c2,c3,c4). The Loose Index Scan access method can be used for the following queries:

```
SELECT c1, c2 FROM t1 GROUP BY c1, c2;

SELECT DISTINCT c1, c2 FROM t1;

SELECT c1, MIN(c2) FROM t1 GROUP BY c1;

SELECT c1, c2 FROM t1 WHERE c1 < const GROUP BY c1, c2;

SELECT MAX(c3), MIN(c3), c1, c2 FROM t1 WHERE c2 > const GROUP BY c1, c2;

SELECT c2 FROM t1 WHERE c1 < const GROUP BY c1, c2;

SELECT c1, c2 FROM t1 WHERE c3 = const GROUP BY c1, c2;
```

The following queries cannot be executed with this quick select method, for the reasons given:

• There are aggregate functions other than MIN() or MAX():

```
SELECT c1, SUM(c2) FROM t1 GROUP BY c1;
```

The columns in the GROUP BY clause do not form a leftmost prefix of the index:

```
SELECT c1, c2 FROM t1 GROUP BY c2, c3;
```

• The query refers to a part of a key that comes after the GROUP BY part, and for which there is no equality with a constant:

```
SELECT c1, c3 FROM t1 GROUP BY c1, c2;
```

Were the query to include WHERE c3 = const, Loose Index Scan could be used.

The Loose Index Scan access method can be applied to other forms of aggregate function references in the select list, in addition to the MIN() and MAX() references already supported:

- AVG(DISTINCT), SUM(DISTINCT), and COUNT(DISTINCT) are supported. AVG(DISTINCT) and SUM(DISTINCT) take a single argument. COUNT(DISTINCT) can have more than one column argument.
- There must be no GROUP BY or DISTINCT clause in the query.
- · The Loose Index Scan limitations described previously still apply.

Assume that there is an index idx(c1,c2,c3) on table t1(c1,c2,c3,c4). The Loose Index Scan access method can be used for the following queries:

```
SELECT COUNT(DISTINCT c1), SUM(DISTINCT c1) FROM t1;
SELECT COUNT(DISTINCT c1, c2), COUNT(DISTINCT c2, c1) FROM t1;
```

## **Tight Index Scan**

A Tight Index Scan may be either a full index scan or a range index scan, depending on the query conditions.

When the conditions for a Loose Index Scan are not met, it still may be possible to avoid creation of temporary tables for GROUP BY queries. If there are range conditions in the WHERE clause, this method reads only the keys that satisfy these conditions. Otherwise, it performs an index scan. Because this method reads all keys in each range defined by the WHERE clause, or scans the whole index if there are no range conditions, it is called a *Tight Index Scan*. With a Tight Index Scan, the grouping operation is performed only after all keys that satisfy the range conditions have been found.

For this method to work, it is sufficient that there be a constant equality condition for all columns in a query referring to parts of the key coming before or in between parts of the <code>GROUP BY</code> key. The constants from the equality conditions fill in any "gaps" in the search keys so that it is possible to form complete prefixes of the index. These index prefixes then can be used for index lookups. If the <code>GROUP BY</code> result requires sorting, and it is possible to form search keys that are prefixes of the index, MySQL also avoids extra sorting operations because searching with prefixes in an ordered index already retrieves all the keys in order.

Assume that there is an index idx(c1,c2,c3) on table t1(c1,c2,c3,c4). The following queries do not work with the Loose Index Scan access method described previously, but still work with the Tight Index Scan access method.

There is a gap in the GROUP BY, but it is covered by the condition c2 = 'a':

```
SELECT c1, c2, c3 FROM t1 WHERE c2 = 'a' GROUP BY c1, c3;
```

The GROUP BY does not begin with the first part of the key, but there is a condition that provides a
constant for that part:

```
SELECT c1, c2, c3 FROM t1 WHERE c1 = 'a' GROUP BY c2, c3;
```

## 10.2.1.18 DISTINCT Optimization

DISTINCT combined with ORDER BY needs a temporary table in many cases.

Because DISTINCT may use GROUP BY, learn how MySQL works with columns in ORDER BY or HAVING clauses that are not part of the selected columns. See Section 14.19.3, "MySQL Handling of GROUP BY".

In most cases, a <code>DISTINCT</code> clause can be considered as a special case of <code>GROUP BY</code>. For example, the following two queries are equivalent:

```
SELECT DISTINCT c1, c2, c3 FROM t1
WHERE c1 > const;

SELECT c1, c2, c3 FROM t1
WHERE c1 > const GROUP BY c1, c2, c3;
```

Due to this equivalence, the optimizations applicable to GROUP BY queries can be also applied to queries with a DISTINCT clause. Thus, for more details on the optimization possibilities for DISTINCT queries, see Section 10.2.1.17, "GROUP BY Optimization".

When combining LIMIT row\_count with DISTINCT, MySQL stops as soon as it finds row\_count unique rows.

If you do not use columns from all tables named in a query, MySQL stops scanning any unused tables as soon as it finds the first match. In the following case, assuming that t1 is used before t2 (which you can check with EXPLAIN), MySQL stops reading from t2 (for any particular row in t1) when it finds the first row in t2:

```
SELECT DISTINCT t1.a FROM t1, t2 where t1.a=t2.a;
```

### 10.2.1.19 LIMIT Query Optimization

If you need only a specified number of rows from a result set, use a LIMIT clause in the query, rather than fetching the whole result set and throwing away the extra data.

MySQL sometimes optimizes a query that has a LIMIT row count clause and no HAVING clause:

- If you select only a few rows with LIMIT, MySQL uses indexes in some cases when normally it would prefer to do a full table scan.
- If you combine LIMIT row\_count with ORDER BY, MySQL stops sorting as soon as it has found the first row\_count rows of the sorted result, rather than sorting the entire result. If ordering is done by using an index, this is very fast. If a filesort must be done, all rows that match the query without the LIMIT clause are selected, and most or all of them are sorted, before the first row\_count are found. After the initial rows have been found, MySQL does not sort any remainder of the result set.

One manifestation of this behavior is that an ORDER BY query with and without LIMIT may return rows in different order, as described later in this section.

- If you combine LIMIT row\_count with DISTINCT, MySQL stops as soon as it finds row\_count unique rows.
- In some cases, a GROUP BY can be resolved by reading the index in order (or doing a sort on the index), then calculating summaries until the index value changes. In this case, LIMIT row\_count does not calculate any unnecessary GROUP BY values.
- As soon as MySQL has sent the required number of rows to the client, it aborts the query unless you are using SQL\_CALC\_FOUND\_ROWS. In that case, the number of rows can be retrieved with SELECT FOUND\_ROWS(). See Section 14.15, "Information Functions".
- LIMIT 0 quickly returns an empty set. This can be useful for checking the validity of a query. It can also be employed to obtain the types of the result columns within applications that use a MySQL API that makes result set metadata available. With the mysql client program, you can use the --column-type-info option to display result column types.

- If the server uses temporary tables to resolve a query, it uses the LIMIT row\_count clause to calculate how much space is required.
- If an index is not used for ORDER BY but a LIMIT clause is also present, the optimizer may be able to avoid using a merge file and sort the rows in memory using an in-memory filesort operation.

If multiple rows have identical values in the ORDER BY columns, the server is free to return those rows in any order, and may do so differently depending on the overall execution plan. In other words, the sort order of those rows is nondeterministic with respect to the nonordered columns.

One factor that affects the execution plan is LIMIT, so an ORDER BY query with and without LIMIT may return rows in different orders. Consider this query, which is sorted by the category column but nondeterministic with respect to the id and rating columns:

```
mysql> SELECT * FROM ratings ORDER BY category;
+----+
| id | category | rating |
+----+
| 1 | 1 | 4.5 |
| 5 | 1 | 3.2 |
| 3 | 2 | 3.7 |
| 4 | 2 | 3.5 |
| 6 | 2 | 3.5 |
| 2 | 3 | 5.0 |
| 7 | 3 | 2.7 |
+----+
```

Including LIMIT may affect order of rows within each category value. For example, this is a valid query result:

In each case, the rows are sorted by the ORDER BY column, which is all that is required by the SQL standard.

If it is important to ensure the same row order with and without LIMIT, include additional columns in the ORDER BY clause to make the order deterministic. For example, if id values are unique, you can make rows for a given category value appear in id order by sorting like this:

```
mysql> SELECT * FROM ratings ORDER BY category, id;
| id | category | rating |
 1 |
         1 | 4.5
  5
          1 | 3.2
           2 |
  3 |
                 3.7
                3.5
           2
  4 |
  6 |
          2 |
                3.5
  2 | 7 |
           3 |
                 5.0
           3 |
                 2.7
mysql> SELECT * FROM ratings ORDER BY category, id LIMIT 5;
| id | category | rating |
  1 | 1 | 4.5
  5 |
           1 |
                3.2
  3
          2 | 3.7
```

For a query with an ORDER BY OF GROUP BY and a LIMIT clause, the optimizer tries to choose an ordered index by default when it appears doing so would speed up query execution. In cases where using some other optimization might be faster, it is possible to turn off this optimization by setting the optimizer\_switch system variable's prefer\_ordering\_index flag to off.

*Example*: First we create and populate a table t as shown here:

```
# Create and populate a table t:

mysql> CREATE TABLE t (
    ->     id1 BIGINT NOT NULL,
    ->     id2 BIGINT NOT NULL,
    ->     c1 VARCHAR(50) NOT NULL,
    ->     c2 VARCHAR(50) NOT NULL,
    ->     pRIMARY KEY (id1),
    ->     INDEX i (id2, c1)
    -> );

# [Insert some rows into table t - not shown]
```

Verify that the prefer\_ordering\_index flag is enabled:

```
mysql> SELECT @@optimizer_switch LIKE '%prefer_ordering_index=on%';
+-------+
| @@optimizer_switch LIKE '%prefer_ordering_index=on%' |
+------+
| 1 |
```

Since the following query has a LIMIT clause, we expect it to use an ordered index, if possible. In this case, as we can see from the EXPLAIN output, it uses the table's primary key.

```
mysql> EXPLAIN SELECT c2 FROM t
  -> WHERE id2 > 3
   ->
        ORDER BY id1 ASC LIMIT 2\G
************************ 1. row *****************
         id: 1
 select_type: SIMPLE
       table: t
  partitions: NULL
       type: index
possible_keys: i
        key: PRIMARY
     key_len: 8
        ref: NULL
        rows: 2
    filtered: 70.00
       Extra: Using where
```

Now we disable the prefer\_ordering\_index flag, and re-run the same query; this time it uses the index i (which includes the id2 column used in the WHERE clause), and a filesort:

```
mysql> SET optimizer_switch = "prefer_ordering_index=off";

mysql> EXPLAIN SELECT c2 FROM t
    -> WHERE id2 > 3
    -> ORDER BY id1 ASC LIMIT 2\G

*****************************
    id: 1
    select_type: SIMPLE
        table: t
    partitions: NULL
        type: range
possible_keys: i
        key: i
        key_len: 8
```

```
ref: NULL
rows: 14
filtered: 100.00
Extra: Using index condition; Using filesort
```

See also Section 10.9.2, "Switchable Optimizations".

### 10.2.1.20 Function Call Optimization

MySQL functions are tagged internally as deterministic or nondeterministic. A function is nondeterministic if, given fixed values for its arguments, it can return different results for different invocations. Examples of nondeterministic functions: RAND(), UUID().

If a function is tagged nondeterministic, a reference to it in a WHERE clause is evaluated for every row (when selecting from one table) or combination of rows (when selecting from a multiple-table join).

MySQL also determines when to evaluate functions based on types of arguments, whether the arguments are table columns or constant values. A deterministic function that takes a table column as argument must be evaluated whenever that column changes value.

Nondeterministic functions may affect query performance. For example, some optimizations may not be available, or more locking might be required. The following discussion uses RAND() but applies to other nondeterministic functions as well.

Suppose that a table t has this definition:

```
CREATE TABLE t (id INT NOT NULL PRIMARY KEY, col_a VARCHAR(100));
```

Consider these two queries:

```
SELECT * FROM t WHERE id = POW(1,2);
SELECT * FROM t WHERE id = FLOOR(1 + RAND() * 49);
```

Both queries appear to use a primary key lookup because of the equality comparison against the primary key, but that is true only for the first of them:

- The first query always produces a maximum of one row because POW() with constant arguments is a constant value and is used for index lookup.
- The second query contains an expression that uses the nondeterministic function RAND(), which is not constant in the query but in fact has a new value for every row of table t. Consequently, the query reads every row of the table, evaluates the predicate for each row, and outputs all rows for which the primary key matches the random value. This might be zero, one, or multiple rows, depending on the id column values and the values in the RAND() sequence.

The effects of nondeterminism are not limited to SELECT statements. This UPDATE statement uses a nondeterministic function to select rows to be modified:

```
UPDATE t SET col_a = some_expr WHERE id = FLOOR(1 + RAND() * 49);
```

Presumably the intent is to update at most a single row for which the primary key matches the expression. However, it might update zero, one, or multiple rows, depending on the id column values and the values in the RAND() sequence.

The behavior just described has implications for performance and replication:

- Because a nondeterministic function does not produce a constant value, the optimizer cannot use strategies that might otherwise be applicable, such as index lookups. The result may be a table scan.
- InnoDB might escalate to a range-key lock rather than taking a single row lock for one matching row.
- Updates that do not execute deterministically are unsafe for replication.

The difficulties stem from the fact that the RAND() function is evaluated once for every row of the table. To avoid multiple function evaluations, use one of these techniques:

Move the expression containing the nondeterministic function to a separate statement, saving the
value in a variable. In the original statement, replace the expression with a reference to the variable,
which the optimizer can treat as a constant value:

```
SET @keyval = FLOOR(1 + RAND() * 49);
UPDATE t SET col_a = some_expr WHERE id = @keyval;
```

 Assign the random value to a variable in a derived table. This technique causes the variable to be assigned a value, once, prior to its use in the comparison in the WHERE clause:

```
UPDATE /*+ NO_MERGE(dt) */ t, (SELECT FLOOR(1 + RAND() * 49) AS r) AS dt
SET col_a = some_expr WHERE id = dt.r;
```

As mentioned previously, a nondeterministic expression in the WHERE clause might prevent optimizations and result in a table scan. However, it may be possible to partially optimize the WHERE clause if other expressions are deterministic. For example:

```
SELECT * FROM t WHERE partial_key=5 AND some_column=RAND();
```

If the optimizer can use partial\_key to reduce the set of rows selected, RAND() is executed fewer times, which diminishes the effect of nondeterminism on optimization.

## 10.2.1.21 Window Function Optimization

Window functions affect the strategies the optimizer considers:

- Derived table merging for a subquery is disabled if the subquery has window functions. The subquery is always materialized.
- Semijoins are not applicable to window function optimization because semijoins apply to subqueries in WHERE and JOIN ... ON, which cannot contain window functions.
- The optimizer processes multiple windows that have the same ordering requirements in sequence, so sorting can be skipped for windows following the first one.
- The optimizer makes no attempt to merge windows that could be evaluated in a single step (for example, when multiple OVER clauses contain identical window definitions). The workaround is to define the window in a WINDOW clause and refer to the window name in the OVER clauses.

An aggregate function not used as a window function is aggregated in the outermost possible query. For example, in this query, MySQL sees that COUNT(t1.b) is something that cannot exist in the outer query because of its placement in the WHERE clause:

```
SELECT * FROM t1 WHERE t1.a = (SELECT COUNT(t1.b) FROM t2);
```

Consequently, MySQL aggregates inside the subquery, treating t1.b as a constant and returning the count of rows of t2.

Replacing WHERE with HAVING results in an error:

```
mysql> SELECT * FROM t1 HAVING t1.a = (SELECT COUNT(t1.b) FROM t2);
ERROR 1140 (42000): In aggregated query without GROUP BY, expression #1
of SELECT list contains nonaggregated column 'test.t1.a'; this is
incompatible with sql_mode=only_full_group_by
```

The error occurs because COUNT(t1.b) can exist in the HAVING, and so makes the outer query aggregated.

Window functions (including aggregate functions used as window functions) do not have the preceding complexity. They always aggregate in the subquery where they are written, never in the outer query.

Window function evaluation may be affected by the value of the windowing\_use\_high\_precision system variable, which determines whether to compute window operations without loss of precision. By default, windowing\_use\_high\_precision is enabled.

For some moving frame aggregates, the inverse aggregate function can be applied to remove values from the aggregate. This can improve performance but possibly with a loss of precision. For example, adding a very small floating-point value to a very large value causes the very small value to be "hidden" by the large value. When inverting the large value later, the effect of the small value is lost.

Loss of precision due to inverse aggregation is a factor only for operations on floating-point (approximate-value) data types. For other types, inverse aggregation is safe; this includes <code>DECIMAL</code>, which permits a fractional part but is an exact-value type.

For faster execution, MySQL always uses inverse aggregation when it is safe:

- For floating-point values, inverse aggregation is not always safe and might result in loss of precision. The default is to avoid inverse aggregation, which is slower but preserves precision. If it is permissible to sacrifice safety for speed, windowing\_use\_high\_precision can be disabled to permit inverse aggregation.
- For nonfloating-point data types, inverse aggregation is always safe and is used regardless of the windowing\_use\_high\_precision value.
- windowing\_use\_high\_precision has no effect on MIN() and MAX(), which do not use inverse
  aggregation in any case.

For evaluation of the variance functions STDDEV\_POP(), STDDEV\_SAMP(), VAR\_POP(), VAR\_SAMP(), and their synonyms, evaluation can occur in optimized mode or default mode. Optimized mode may produce slightly different results in the last significant digits. If such differences are permissible, windowing\_use\_high\_precision can be disabled to permit optimized mode.

For EXPLAIN, windowing execution plan information is too extensive to display in traditional output format. To see windowing information, use EXPLAIN FORMAT=JSON and look for the windowing element.

### 10.2.1.22 Row Constructor Expression Optimization

Row constructors permit simultaneous comparisons of multiple values. For example, these two statements are semantically equivalent:

```
SELECT * FROM t1 WHERE (column1,column2) = (1,1);
SELECT * FROM t1 WHERE column1 = 1 AND column2 = 1;
```

In addition, the optimizer handles both expressions the same way.

The optimizer is less likely to use available indexes if the row constructor columns do not cover the prefix of an index. Consider the following table, which has a primary key on (c1, c2, c3):

```
CREATE TABLE t1 (
c1 INT, c2 INT, c3 INT, c4 CHAR(100),
PRIMARY KEY(c1,c2,c3)
);
```

In this query, the WHERE clause uses all columns in the index. However, the row constructor itself does not cover an index prefix, with the result that the optimizer uses only c1 (key\_len=4, the size of c1):

```
possible_keys: PRIMARY
key: PRIMARY
key_len: 4
ref: const
rows: 3
filtered: 100.00
Extra: Using where
```

In such cases, rewriting the row constructor expression using an equivalent nonconstructor expression may result in more complete index use. For the given query, the row constructor and equivalent nonconstructor expressions are:

```
(c2,c3) > (1,1)

c2 > 1 \text{ OR } ((c2 = 1) \text{ AND } (c3 > 1))
```

Rewriting the query to use the nonconstructor expression results in the optimizer using all three columns in the index (key len=12):

Thus, for better results, avoid mixing row constructors with AND/OR expressions. Use one or the other.

Under certain conditions, the optimizer can apply the range access method to IN() expressions that have row constructor arguments. See Range Optimization of Row Constructor Expressions.

### 10.2.1.23 Avoiding Full Table Scans

The output from EXPLAIN shows ALL in the type column when MySQL uses a full table scan to resolve a query. This usually happens under the following conditions:

- The table is so small that it is faster to perform a table scan than to bother with a key lookup. This is common for tables with fewer than 10 rows and a short row length.
- There are no usable restrictions in the ON or WHERE clause for indexed columns.
- You are comparing indexed columns with constant values and MySQL has calculated (based on the index tree) that the constants cover too large a part of the table and that a table scan would be faster. See Section 10.2.1.1, "WHERE Clause Optimization".
- You are using a key with low cardinality (many rows match the key value) through another column.
   In this case, MySQL assumes that by using the key probably requires many key lookups and that a table scan would be faster.

For small tables, a table scan often is appropriate and the performance impact is negligible. For large tables, try the following techniques to avoid having the optimizer incorrectly choose a table scan:

- Use ANALYZE TABLE *tb1\_name* to update the key distributions for the scanned table. See Section 15.7.3.1, "ANALYZE TABLE Statement".
- Use FORCE INDEX for the scanned table to tell MySQL that table scans are very expensive compared to using the given index:

```
SELECT * FROM t1, t2 FORCE INDEX (index_for_column)
WHERE t1.col_name=t2.col_name;
```

See Section 10.9.4, "Index Hints".

• Start mysqld with the --max-seeks-for-key=1000 option or use SET max\_seeks\_for\_key=1000 to tell the optimizer to assume that no key scan causes more than 1,000 key seeks. See Section 7.1.8, "Server System Variables".

# 10.2.2 Optimizing Subqueries, Derived Tables, View References, and Common Table Expressions

The MySQL query optimizer has different strategies available to evaluate subqueries:

- For a subquery used with an IN, = ANY, or EXISTS predicate, the optimizer has these choices:
  - Semijoin
  - Materialization
  - EXISTS strategy
- For a subquery used with a NOT IN, <> ALL or NOT EXISTS predicate, the optimizer has these
  choices:
  - Materialization
  - EXISTS strategy

For a derived table, the optimizer has these choices (which also apply to view references and common table expressions):

- · Merge the derived table into the outer query block
- Materialize the derived table to an internal temporary table

The following discussion provides more information about the preceding optimization strategies.



### Note

A limitation on UPDATE and DELETE statements that use a subquery to modify a single table is that the optimizer does not use semijoin or materialization subquery optimizations. As a workaround, try rewriting them as multiple-table UPDATE and DELETE statements that use a join rather than a subquery.

# 10.2.2.1 Optimizing IN and EXISTS Subquery Predicates with Semijoin and Antijoin Transformations

A semijoin is a preparation-time transformation that enables multiple execution strategies such as table pullout, duplicate weedout, first match, loose scan, and materialization. The optimizer uses semijoin strategies to improve subquery execution, as described in this section.

For an inner join between two tables, the join returns a row from one table as many times as there are matches in the other table. But for some questions, the only information that matters is whether there is a match, not the number of matches. Suppose that there are tables named class and roster that list classes in a course curriculum and class rosters (students enrolled in each class), respectively. To list the classes that actually have students enrolled, you could use this join:

```
SELECT class.class_num, class.class_name
FROM class
INNER JOIN roster
```

```
WHERE class.class_num = roster.class_num;
```

However, the result lists each class once for each enrolled student. For the question being asked, this is unnecessary duplication of information.

Assuming that class\_num is a primary key in the class table, duplicate suppression is possible by using SELECT DISTINCT, but it is inefficient to generate all matching rows first only to eliminate duplicates later.

The same duplicate-free result can be obtained by using a subquery:

```
SELECT class_num, class_name
FROM class
WHERE class_num IN
(SELECT class_num FROM roster);
```

Here, the optimizer can recognize that the IN clause requires the subquery to return only one instance of each class number from the roster table. In this case, the query can use a *semijoin*; that is, an operation that returns only one instance of each row in class that is matched by rows in roster.

The following statement, which contains an EXISTS subquery predicate, is equivalent to the previous statement containing an IN subquery predicate:

```
SELECT class_num, class_name
FROM class
WHERE EXISTS
(SELECT * FROM roster WHERE class.class_num = roster.class_num);
```

Any statement with an EXISTS subquery predicate is subject to the same semijoin transforms as a statement with an equivalent IN subquery predicate.

The following subqueries are transformed into antijoins:

```
NOT IN (SELECT ... FROM ...)
NOT EXISTS (SELECT ... FROM ...)
IN (SELECT ... FROM ...) IS NOT TRUE
EXISTS (SELECT ... FROM ...) IS NOT TRUE.
IN (SELECT ... FROM ...) IS FALSE
EXISTS (SELECT ... FROM ...) IS FALSE.
```

In short, any negation of a subquery of the form IN (SELECT ... FROM ...) or EXISTS (SELECT ... FROM ...) is transformed into an antijoin.

An antijoin is an operation that returns only rows for which there is no match. Consider the query shown here:

```
SELECT class_num, class_name
FROM class
WHERE class_num NOT IN
(SELECT class_num FROM roster);
```

This query is rewritten internally as the antijoin SELECT class\_num, class\_name FROM class ANTIJOIN roster ON class\_num, which returns one instance of each row in class that is *not* matched by any rows in roster. This means that, for each row in class, as soon as a match is found in roster, the row in class can be discarded.

Antijoin transformations cannot in most cases be applied if the expressions being compared are nullable. An exception to this rule is that (... NOT IN (SELECT ...)) IS NOT FALSE and its equivalent (... IN (SELECT ...)) IS NOT TRUE can be transformed into antijoins.

Outer join and inner join syntax is permitted in the outer query specification, and table references may be base tables, derived tables, view references, or common table expressions.

In MySQL, a subquery must satisfy these criteria to be handled as a semijoin (or an antijoin, if NOT modifies the subquery):

• It must be part of an IN, = ANY, or EXISTS predicate that appears at the top level of the WHERE or ON clause, possibly as a term in an AND expression. For example:

```
SELECT ...

FROM otl, ...

WHERE (oel, ...) IN

(SELECT iel, ... FROM itl, ... WHERE ...);
```

Here,  $ot_i$  and  $it_i$  represent tables in the outer and inner parts of the query, and  $oe_i$  and  $ie_i$  represent expressions that refer to columns in the outer and inner tables.

The subquery can also be the argument to an expression modified by NOT, IS [NOT] TRUE, or IS [NOT] FALSE.

- It must be a single SELECT without UNION constructs.
- It must not contain a HAVING clause.
- It must not contain any aggregate functions (whether it is explicitly or implicitly grouped).
- It must not have a LIMIT clause.
- The statement must not use the STRAIGHT\_JOIN join type in the outer query.
- The STRAIGHT JOIN modifier must not be present.
- The number of outer and inner tables together must be less than the maximum number of tables permitted in a join.
- The subquery may be correlated or uncorrelated. Decorrelation looks at trivially correlated predicates in the WHERE clause of a subquery used as the argument to EXISTS, and makes it possible to optimize it as if it was used within IN (SELECT b FROM ...). The term trivially correlated means that the predicate is an equality predicate, that it is the sole predicate in the WHERE clause (or is combined with AND), and that one operand is from a table referenced in the subquery and the other operand is from the outer query block.
- The DISTINCT keyword is permitted but ignored. Semijoin strategies automatically handle duplicate removal.
- A GROUP BY clause is permitted but ignored, unless the subquery also contains one or more aggregate functions.
- An ORDER BY clause is permitted but ignored, since ordering is irrelevant to the evaluation of semijoin strategies.

If a subquery meets the preceding criteria, MySQL converts it to a semijoin (or to an antijoin if applicable) and makes a cost-based choice from these strategies:

- Convert the subquery to a join, or use table pullout and run the query as an inner join between subquery tables and outer tables. Table pullout pulls a table out from the subquery to the outer query.
- Duplicate Weedout: Run the semijoin as if it was a join and remove duplicate records using a temporary table.
- FirstMatch: When scanning the inner tables for row combinations and there are multiple instances of a given value group, choose one rather than returning them all. This "shortcuts" scanning and eliminates production of unnecessary rows.

- LooseScan: Scan a subquery table using an index that enables a single value to be chosen from each subquery's value group.
- Materialize the subquery into an indexed temporary table that is used to perform a join, where the
  index is used to remove duplicates. The index might also be used later for lookups when joining
  the temporary table with the outer tables; if not, the table is scanned. For more information about
  materialization, see Section 10.2.2.2, "Optimizing Subqueries with Materialization".

Each of these strategies can be enabled or disabled using the following <code>optimizer\_switch</code> system variable flags:

- The semijoin flag controls whether semijoins and antijoins are used.
- If semijoin is enabled, the firstmatch, loosescan, duplicateweedout, and materialization flags enable finer control over the permitted semijoin strategies.
- If the duplicateweedout semijoin strategy is disabled, it is not used unless all other applicable strategies are also disabled.
- If duplicateweedout is disabled, on occasion the optimizer may generate a query plan that is far from optimal. This occurs due to heuristic pruning during greedy search, which can be avoided by setting optimizer\_prune\_level=0.

These flags are enabled by default. See Section 10.9.2, "Switchable Optimizations".

The optimizer minimizes differences in handling of views and derived tables. This affects queries that use the STRAIGHT\_JOIN modifier and a view with an IN subquery that can be converted to a semijoin. The following query illustrates this because the change in processing causes a change in transformation, and thus a different execution strategy:

```
CREATE VIEW v AS
SELECT *
FROM t1
WHERE a IN (SELECT b
FROM t2);

SELECT STRAIGHT_JOIN *
FROM t3 JOIN v ON t3.x = v.a;
```

The optimizer first looks at the view and converts the IN subquery to a semijoin, then checks whether it is possible to merge the view into the outer query. Because the STRAIGHT\_JOIN modifier in the outer query prevents semijoin, the optimizer refuses the merge, causing derived table evaluation using a materialized table.

EXPLAIN output indicates the use of semijoin strategies as follows:

- For extended EXPLAIN output, the text displayed by a following SHOW WARNINGS shows the
  rewritten query, which displays the semijoin structure. (See Section 10.8.3, "Extended EXPLAIN
  Output Format".) From this you can get an idea about which tables were pulled out of the semijoin. If
  a subquery was converted to a semijoin, you should see that the subquery predicate is gone and its
  tables and WHERE clause were merged into the outer query join list and WHERE clause.
- Temporary table use for Duplicate Weedout is indicated by Start temporary and End temporary in the Extra column. Tables that were not pulled out and are in the range of EXPLAIN output rows covered by Start temporary and End temporary have their rowid in the temporary table.
- FirstMatch(tbl name) in the Extra column indicates join shortcutting.
- LooseScan(m..n) in the Extra column indicates use of the LooseScan strategy. m and n are key part numbers.
- Temporary table use for materialization is indicated by rows with a select\_type value of MATERIALIZED and rows with a table value of <subqueryN>.

A semijoin transformation can also be applied to a single-table UPDATE or DELETE statement that uses a [NOT] IN or [NOT] EXISTS subquery predicate, provided that the statement does not use ORDER BY or LIMIT, and that semijoin transformations are allowed by an optimizer hint or by the optimizer\_switch setting.

## 10.2.2.2 Optimizing Subqueries with Materialization

The optimizer uses materialization to enable more efficient subquery processing. Materialization speeds up query execution by generating a subquery result as a temporary table, normally in memory. The first time MySQL needs the subquery result, it materializes that result into a temporary table. Any subsequent time the result is needed, MySQL refers again to the temporary table. The optimizer may index the table with a hash index to make lookups fast and inexpensive. The index contains unique values to eliminate duplicates and make the table smaller.

Subquery materialization uses an in-memory temporary table when possible, falling back to on-disk storage if the table becomes too large. See Section 10.4.4, "Internal Temporary Table Use in MySQL".

If materialization is not used, the optimizer sometimes rewrites a noncorrelated subquery as a correlated subquery. For example, the following IN subquery is noncorrelated (where\_condition involves only columns from t2 and not t1):

```
SELECT * FROM t1
WHERE t1.a IN (SELECT t2.b FROM t2 WHERE where_condition);
```

The optimizer might rewrite this as an EXISTS correlated subquery:

```
SELECT * FROM t1
WHERE EXISTS (SELECT t2.b FROM t2 WHERE where_condition AND t1.a=t2.b);
```

Subquery materialization using a temporary table avoids such rewrites and makes it possible to execute the subquery only once rather than once per row of the outer query.

For subquery materialization to be used in MySQL, the <code>optimizer\_switch</code> system variable materialization flag must be enabled. (See Section 10.9.2, "Switchable Optimizations".) With the materialization flag enabled, materialization applies to subquery predicates that appear anywhere (in the select list, <code>WHERE</code>, <code>ON</code>, <code>GROUP</code> <code>BY</code>, <code>HAVING</code>, or <code>ORDER</code> <code>BY</code>), for predicates that fall into any of these use cases:

The predicate has this form, when no outer expression oe\_i or inner expression ie\_i is nullable. N
is 1 or larger.

```
(oe_1, oe_2, ..., oe_N) [NOT] IN (SELECT ie_1, i_2, ..., ie_N ...)
```

• The predicate has this form, when there is a single outer expression *oe* and inner expression *ie*. The expressions can be nullable.

```
oe [NOT] IN (SELECT ie ...)
```

• The predicate is IN or NOT IN and a result of UNKNOWN (NULL) has the same meaning as a result of FALSE.

The following examples illustrate how the requirement for equivalence of UNKNOWN and FALSE predicate evaluation affects whether subquery materialization can be used. Assume that where\_condition involves columns only from t2 and not t1 so that the subquery is noncorrelated.

This query is subject to materialization:

```
SELECT * FROM t1
WHERE t1.a IN (SELECT t2.b FROM t2 WHERE where_condition);
```

Here, it does not matter whether the IN predicate returns UNKNOWN or FALSE. Either way, the row from t1 is not included in the query result.

An example where subquery materialization is not used is the following query, where t2.b is a nullable column:

```
SELECT * FROM t1
WHERE (t1.a,t1.b) NOT IN (SELECT t2.a,t2.b FROM t2
WHERE where_condition);
```

The following restrictions apply to the use of subquery materialization:

- The types of the inner and outer expressions must match. For example, the optimizer might be able to use materialization if both expressions are integer or both are decimal, but cannot if one expression is integer and the other is decimal.
- The inner expression cannot be a BLOB.

Use of EXPLAIN with a query provides some indication of whether the optimizer uses subquery materialization:

- Compared to query execution that does not use materialization, select\_type may change from DEPENDENT SUBQUERY to SUBQUERY. This indicates that, for a subquery that would be executed once per outer row, materialization enables the subquery to be executed just once.
- For extended EXPLAIN output, the text displayed by a following SHOW WARNINGS includes materialize and materialized-subquery.

MySQL can also apply subquery materialization to a single-table UPDATE or DELETE statement that uses a [NOT] IN OR [NOT] EXISTS subquery predicate, provided that the statement does not use ORDER BY OR LIMIT, and that subquery materialization is allowed by an optimizer hint or by the optimizer\_switch setting.

## 10.2.2.3 Optimizing Subqueries with the EXISTS Strategy

Certain optimizations are applicable to comparisons that use the IN (or =ANY) operator to test subquery results. This section discusses these optimizations, particularly with regard to the challenges that NULL values present. The last part of the discussion suggests how you can help the optimizer.

Consider the following subquery comparison:

```
outer_expr IN (SELECT inner_expr FROM ... WHERE subquery_where)
```

MySQL evaluates queries "from outside to inside." That is, it first obtains the value of the outer expression *outer\_expr*, and then runs the subquery and captures the rows that it produces.

A very useful optimization is to "inform" the subquery that the only rows of interest are those where the inner expression <code>inner\_expr</code> is equal to <code>outer\_expr</code>. This is done by pushing down an appropriate equality into the subquery's <code>WHERE</code> clause to make it more restrictive. The converted comparison looks like this:

```
EXISTS (SELECT 1 FROM ... WHERE subquery_where AND outer_expr=inner_expr)
```

After the conversion, MySQL can use the pushed-down equality to limit the number of rows it must examine to evaluate the subquery.

More generally, a comparison of N values to a subquery that returns N-value rows is subject to the same conversion. If  $oe_i$  and  $ie_i$  represent corresponding outer and inner expression values, this subquery comparison:

```
(oe_1, ..., oe_N) IN (SELECT ie_1, ..., ie_N FROM ... WHERE subquery_where)
```

### Becomes:

```
EXISTS (SELECT 1 FROM ... WHERE subquery_where

AND oe_1 = ie_1
```

```
AND ...
AND oe_N = ie_N
```

For simplicity, the following discussion assumes a single pair of outer and inner expression values.

The "pushdown" strategy just described works if either of these conditions is true:

- outer\_expr and inner\_expr cannot be NULL.
- You need not distinguish NULL from FALSE subquery results. If the subquery is a part of an OR
  or AND expression in the WHERE clause, MySQL assumes that you do not care. Another instance
  where the optimizer notices that NULL and FALSE subquery results need not be distinguished is this
  construct:

```
... WHERE outer_expr IN (subquery)
```

In this case, the WHERE clause rejects the row whether IN (subquery) returns NULL or FALSE.

Suppose that  $outer\_expr$  is known to be a non-NULL value but the subquery does not produce a row such that  $outer\_expr$  =  $inner\_expr$ . Then  $outer\_expr$  IN (SELECT ...) evaluates as follows:

- NULL, if the SELECT produces any row where inner\_expr is NULL
- FALSE, if the SELECT produces only non-NULL values or produces nothing

In this situation, the approach of looking for rows with  $outer\_expr = inner\_expr$  is no longer valid. It is necessary to look for such rows, but if none are found, also look for rows where  $inner\_expr$  is NULL. Roughly speaking, the subquery can be converted to something like this:

```
EXISTS (SELECT 1 FROM ... WHERE subquery_where AND (outer_expr=inner_expr OR inner_expr IS NULL))
```

The need to evaluate the extra IS NULL condition is why MySQL has the ref\_or\_null access method:

```
mysql> EXPLAIN
     SELECT outer_expr IN (SELECT t2.maybe_null_key
                       FROM t2, t3 WHERE ...)
     FROM t1;
id: 1
 select_type: PRIMARY
      table: t1
*********************** 2. row ******************
        id: 2
 select_type: DEPENDENT SUBQUERY
      table: t2
       type: ref_or_null
possible_keys: maybe_null_key
       key: maybe_null_key
    key_len: 5
       ref: func
       rows: 2
      Extra: Using where; Using index
```

The unique\_subquery and index\_subquery subquery-specific access methods also have "or NULL" variants.

The additional OR ... IS NULL condition makes query execution slightly more complicated (and some optimizations within the subquery become inapplicable), but generally this is tolerable.

The situation is much worse when *outer\_expr* can be NULL. According to the SQL interpretation of NULL as "unknown value," NULL IN (SELECT *inner expr* ...) should evaluate to:

NULL, if the SELECT produces any rows

• FALSE, if the SELECT produces no rows

For proper evaluation, it is necessary to be able to check whether the SELECT has produced any rows at all, so  $outer\_expr = inner\_expr$  cannot be pushed down into the subquery. This is a problem because many real world subqueries become very slow unless the equality can be pushed down.

Essentially, there must be different ways to execute the subquery depending on the value of <code>outer\_expr</code>.

The optimizer chooses SQL compliance over speed, so it accounts for the possibility that <code>outer\_expr</code> might be <code>NULL</code>:

• If *outer\_expr* is NULL, to evaluate the following expression, it is necessary to execute the SELECT to determine whether it produces any rows:

```
NULL IN (SELECT inner_expr FROM ... WHERE subquery_where)
```

It is necessary to execute the original SELECT here, without any pushed-down equalities of the kind mentioned previously.

• On the other hand, when outer\_expr is not NULL, it is absolutely essential that this comparison:

```
outer_expr IN (SELECT inner_expr FROM ... WHERE subquery_where)
```

Be converted to this expression that uses a pushed-down condition:

```
EXISTS (SELECT 1 FROM ... WHERE subquery_where AND outer_expr=inner_expr)
```

Without this conversion, subqueries are slow.

To solve the dilemma of whether or not to push down conditions into the subquery, the conditions are wrapped within "trigger" functions. Thus, an expression of the following form:

```
outer_expr IN (SELECT inner_expr FROM ... WHERE subquery_where)
```

Is converted into:

```
EXISTS (SELECT 1 FROM ... WHERE subquery_where

AND trigcond(outer_expr=inner_expr))
```

More generally, if the subquery comparison is based on several pairs of outer and inner expressions, the conversion takes this comparison:

```
(oe_1, ..., oe_N) IN (SELECT ie_1, ..., ie_N FROM ... WHERE subquery_where)
```

And converts it to this expression:

Each trigcond(X) is a special function that evaluates to the following values:

- X when the "linked" outer expression oe\_i is not NULL
- TRUE when the "linked" outer expression oe\_i is NULL



#### Note

Trigger functions are *not* triggers of the kind that you create with CREATE TRIGGER.

Equalities that are wrapped within trigcond() functions are not first class predicates for the query optimizer. Most optimizations cannot deal with predicates that may be turned on and off at query

execution time, so they assume any trigcond(X) to be an unknown function and ignore it. Triggered equalities can be used by those optimizations:

- Reference optimizations: trigcond(X=Y [OR Y IS NULL]) can be used to construct ref, eq\_ref, or ref\_or\_null table accesses.
- Index lookup-based subquery execution engines: trigcond(X=Y) can be used to construct unique\_subquery or index\_subquery accesses.
- Table-condition generator: If the subquery is a join of several tables, the triggered condition is checked as soon as possible.

When the optimizer uses a triggered condition to create some kind of index lookup-based access (as for the first two items of the preceding list), it must have a fallback strategy for the case when the condition is turned off. This fallback strategy is always the same: Do a full table scan. In EXPLAIN output, the fallback shows up as Full scan on NULL key in the Extra column:

```
mysql> EXPLAIN SELECT t1.col1,
      t1.col1 IN (SELECT t2.key1 FROM t2 WHERE t2.col2=t1.col2) FROM t1\G
******* 1. row *****
         id: 1
 select_type: PRIMARY
      table: t1
************************ 2. row *****************
        id: 2
 select_type: DEPENDENT SUBQUERY
       table: t2
        type: index_subquery
possible_keys: key1
        key: key1
     key_len: 5
        ref: func
        rows: 2
       Extra: Using where; Full scan on NULL key
```

If you run EXPLAIN followed by SHOW WARNINGS, you can see the triggered condition:

The use of triggered conditions has some performance implications. A NULL IN (SELECT ...) expression now may cause a full table scan (which is slow) when it previously did not. This is the price paid for correct results (the goal of the trigger-condition strategy is to improve compliance, not speed).

For multiple-table subqueries, execution of  $\mathtt{NULL}$  IN (SELECT ...) is particularly slow because the join optimizer does not optimize for the case where the outer expression is  $\mathtt{NULL}$ . It assumes that subquery evaluations with  $\mathtt{NULL}$  on the left side are very rare, even if there are statistics that indicate otherwise. On the other hand, if the outer expression might be  $\mathtt{NULL}$  but never actually is, there is no performance penalty.

To help the query optimizer better execute your queries, use these suggestions:

- Declare a column as NOT NULL if it really is. This also helps other aspects of the optimizer by simplifying condition testing for the column.
- If you need not distinguish a NULL from FALSE subquery result, you can easily avoid the slow execution path. Replace a comparison that looks like this:

```
outer_expr [NOT] IN (SELECT inner_expr FROM ...)
```

with this expression:

```
(outer_expr IS NOT NULL) AND (outer_expr [NOT] IN (SELECT inner_expr FROM ...))
```

Then NULL IN (SELECT ...) is never evaluated because MySQL stops evaluating AND parts as soon as the expression result is clear.

Another possible rewrite:

```
[NOT] EXISTS (SELECT inner_expr FROM ...

WHERE inner_expr=outer_expr)
```

The subquery\_materialization\_cost\_based flag of the optimizer\_switch system variable enables control over the choice between subquery materialization and IN-to-EXISTS subquery transformation. See Section 10.9.2, "Switchable Optimizations".

# 10.2.2.4 Optimizing Derived Tables, View References, and Common Table Expressions with Merging or Materialization

The optimizer can handle derived table references using two strategies (which also apply to view references and common table expressions):

- Merge the derived table into the outer query block
- Materialize the derived table to an internal temporary table

### Example 1:

```
SELECT * FROM (SELECT * FROM t1) AS derived_t1;
```

With merging of the derived table derived\_t1, that query is executed similar to:

```
SELECT * FROM t1;
```

### Example 2:

```
SELECT *
FROM t1 JOIN (SELECT t2.f1 FROM t2) AS derived_t2 ON t1.f2=derived_t2.f1
WHERE t1.f1 > 0;
```

With merging of the derived table derived\_t2, that query is executed similar to:

```
SELECT t1.*, t2.f1
FROM t1 JOIN t2 ON t1.f2=t2.f1
WHERE t1.f1 > 0;
```

With materialization, derived\_t1 and derived\_t2 are each treated as a separate table within their respective queries.

The optimizer handles derived tables, view references, and common table expressions the same way: It avoids unnecessary materialization whenever possible, which enables pushing down conditions from the outer query to derived tables and produces more efficient execution plans. (For an example, see Section 10.2.2.2, "Optimizing Subqueries with Materialization".)

If merging would result in an outer query block that references more than 61 base tables, the optimizer chooses materialization instead.

The optimizer propagates an ORDER BY clause in a derived table or view reference to the outer query block if these conditions are all true:

• The outer query is not grouped or aggregated.

- The outer query does not specify DISTINCT, HAVING, or ORDER BY.
- The outer query has this derived table or view reference as the only source in the FROM clause.

Otherwise, the optimizer ignores the ORDER BY clause.

The following means are available to influence whether the optimizer attempts to merge derived tables, view references, and common table expressions into the outer query block:

- The MERGE and NO\_MERGE optimizer hints can be used. They apply assuming that no other rule prevents merging. See Section 10.9.3, "Optimizer Hints".
- Similarly, you can use the derived\_merge flag of the optimizer\_switch system variable. See Section 10.9.2, "Switchable Optimizations". By default, the flag is enabled to permit merging. Disabling the flag prevents merging and avoids ER\_UPDATE\_TABLE\_USED errors.

The derived\_merge flag also applies to views that contain no ALGORITHM clause. Thus, if an ER\_UPDATE\_TABLE\_USED error occurs for a view reference that uses an expression equivalent to the subquery, adding ALGORITHM=TEMPTABLE to the view definition prevents merging and takes precedence over the derived\_merge value.

- It is possible to disable merging by using in the subquery any constructs that prevent merging, although these are not as explicit in their effect on materialization. Constructs that prevent merging are the same for derived tables, common table expressions, and view references:
  - Aggregate functions or window functions (SUM(), MIN(), MAX(), COUNT(), and so forth)
  - DISTINCT
  - GROUP BY
  - HAVING
  - LIMIT
  - UNION OF UNION ALL
  - · Subqueries in the select list
  - Assignments to user variables
  - References only to literal values (in this case, there is no underlying table)

If the optimizer chooses the materialization strategy rather than merging for a derived table, it handles the query as follows:

- The optimizer postpones derived table materialization until its contents are needed during query
  execution. This improves performance because delaying materialization may result in not having to
  do it at all. Consider a query that joins the result of a derived table to another table: If the optimizer
  processes that other table first and finds that it returns no rows, the join need not be carried out
  further and the optimizer can completely skip materializing the derived table.
- During query execution, the optimizer may add an index to a derived table to speed up row retrieval from it.

Consider the following EXPLAIN statement, for a SELECT query that contains a derived table:

```
EXPLAIN SELECT * FROM (SELECT * FROM t1) AS derived_t1;
```

The optimizer avoids materializing the derived table by delaying it until the result is needed during SELECT execution. In this case, the query is not executed (because it occurs in an EXPLAIN statement), so the result is never needed.

Even for queries that are executed, delay of derived table materialization may enable the optimizer to avoid materialization entirely. When this happens, query execution is quicker by the time needed to perform materialization. Consider the following query, which joins the result of a derived table to another table:

```
SELECT *
  FROM t1 JOIN (SELECT t2.f1 FROM t2) AS derived_t2
          ON t1.f2=derived_t2.f1
WHERE t1.f1 > 0;
```

If the optimization processes t1 first and the WHERE clause produces an empty result, the join must necessarily be empty and the derived table need not be materialized.

For cases when a derived table requires materialization, the optimizer may add an index to the materialized table to speed up access to it. If such an index enables ref access to the table, it can greatly reduce amount of data read during query execution. Consider the following query:

```
SELECT *
FROM t1 JOIN (SELECT DISTINCT f1 FROM t2) AS derived_t2
ON t1.f1=derived_t2.f1;
```

The optimizer constructs an index over column f1 from derived\_t2 if doing so would enable use of ref access for the lowest cost execution plan. After adding the index, the optimizer can treat the materialized derived table the same as a regular table with an index, and it benefits similarly from the generated index. The overhead of index creation is negligible compared to the cost of query execution without the index. If ref access would result in higher cost than some other access method, the optimizer creates no index and loses nothing.

For optimizer trace output, a merged derived table or view reference is not shown as a node. Only its underlying tables appear in the top query's plan.

What is true for materialization of derived tables is also true for common table expressions (CTEs). In addition, the following considerations pertain specifically to CTEs.

If a CTE is materialized by a query, it is materialized once for the query, even if the query references it several times.

A recursive CTE is always materialized.

If a CTE is materialized, the optimizer automatically adds relevant indexes if it estimates that indexing can speed up access by the top-level statement to the CTE. This is similar to automatic indexing of derived tables, except that if the CTE is referenced multiple times, the optimizer may create multiple indexes, to speed up access by each reference in the most appropriate way.

The MERGE and NO\_MERGE optimizer hints can be applied to CTEs. Each CTE reference in the top-level statement can have its own hint, permitting CTE references to be selectively merged or materialized. The following statement uses hints to indicate that ctel should be merged and ctel should be materialized:

```
WITH
  ctel AS (SELECT a, b FROM table1),
  cte2 AS (SELECT c, d FROM table2)
SELECT /*+ MERGE(cte1) NO_MERGE(cte2) */ cte1.b, cte2.d
FROM cte1 JOIN cte2
WHERE cte1.a = cte2.c;
```

The ALGORITHM clause for CREATE VIEW does not affect materialization for any WITH clause preceding the SELECT statement in the view definition. Consider this statement:

```
CREATE ALGORITHM={TEMPTABLE|MERGE} VIEW v1 AS WITH ... SELECT ...
```

The ALGORITHM value affects materialization only of the SELECT, not the WITH clause.

As mentioned previously, a CTE, if materialized, is materialized once, even if referenced multiple times. To indicate one-time materialization, optimizer trace output contains an occurrence of creating\_tmp\_table plus one or more occurrences of reusing\_tmp\_table.

CTEs are similar to derived tables, for which the materialized\_from\_subquery node follows the reference. This is true for a CTE that is referenced multiple times, so there is no duplication of materialized\_from\_subquery nodes (which would give the impression that the subquery is executed multiple times, and produce unnecessarily verbose output). Only one reference to the CTE has a complete materialized\_from\_subquery node with the description of its subquery plan. Other references have a reduced materialized\_from\_subquery node. The same idea applies to EXPLAIN output in TRADITIONAL format: Subqueries for other references are not shown.

## 10.2.2.5 Derived Condition Pushdown Optimization

MySQL supports derived condition pushdown for eligible subqueries. For a query such as SELECT \* FROM (SELECT i, j FROM t1) AS dt WHERE i > constant, it is possible in many cases to push the outer WHERE condition down to the derived table, in this case resulting in SELECT \* FROM (SELECT i, j FROM t1 WHERE i > constant) AS dt. When a derived table cannot be merged into the outer query (for example, if the derived table uses aggregation), pushing the outer WHERE condition down to the derived table should decrease the number of rows that need to be processed and thus speed up execution of the query.

Outer WHERE conditions can be pushed down to derived materialized tables under the following circumstances:

When the derived table uses no aggregate or window functions, the outer WHERE condition can be
pushed down to it directly. This includes WHERE conditions having multiple predicates joined with
AND, OR, or both.

For example, the query SELECT \* FROM (SELECT f1, f2 FROM t1) AS dt WHERE f1 < 3 AND f2 > 11 is rewritten as SELECT f1, f2 FROM (SELECT f1, f2 FROM t1 WHERE f1 < 3 AND f2 > 11) AS dt.

When the derived table has a GROUP BY and uses no window functions, an outer WHERE condition
referencing one or more columns which are not part of the GROUP BY can be pushed down to the
derived table as a HAVING condition.

For example, SELECT \* FROM (SELECT i, j, SUM(k) AS sum FROM t1 GROUP BY i, j) AS dt WHERE sum > 100 is rewritten following derived condition pushdown as SELECT \* FROM (SELECT i, j, SUM(k) AS sum FROM t1 GROUP BY i, j HAVING sum > 100) AS dt.

 When the derived table uses a GROUP BY and the columns in the outer WHERE condition are GROUP BY columns, the WHERE conditions referencing those columns can be pushed down directly to the derived table.

```
For example, the query SELECT * FROM (SELECT i,j, SUM(k) AS sum FROM t1 GROUP BY i,j) AS dt WHERE i > 10 is rewritten as SELECT * FROM (SELECT i,j, SUM(k) AS sum FROM t1 WHERE i > 10 GROUP BY i,j) AS dt.
```

In the event that the outer WHERE condition has predicates referencing columns which are part of the GROUP BY as well as predicates referencing columns which are not, predicates of the former sort are pushed down as WHERE conditions, while those of the latter type are pushed down as HAVING conditions. For example, in the query SELECT \* FROM (SELECT i, j, SUM(k) AS sum FROM t1 GROUP BY i,j) AS dt WHERE i > 10 AND sum > 100, the predicate i > 10 in the outer WHERE clause references a GROUP BY column, whereas the predicate sum > 100 does not reference any GROUP BY column. Thus the derived table pushdown optimization causes the query to be rewritten in a manner similar to what is shown here:

```
SELECT * FROM (
SELECT i, j, SUM(k) AS sum FROM t1
WHERE i > 10
```

```
GROUP BY i, j
HAVING sum > 100
) AS dt;
```

To enable derived condition pushdown, the <code>optimizer\_switch</code> system variable's <code>derived\_condition\_pushdown</code> flag (added in this release) must be set to on, which is the default setting. If this optimization is disabled by <code>optimizer\_switch</code>, you can enable it for a specific query using the <code>DERIVED\_CONDITION\_PUSHDOWN</code> optimizer hint. To disable the optimization for a given query, use the <code>NO\_DERIVED\_CONDITION\_PUSHDOWN</code> optimizer hint.

The following restrictions and limitations apply to the derived table condition pushdown optimization:

- The derived table condition pushdown optimization can be employed with UNION queries, with the following exceptions:
  - Condition pushdown cannot be used with a UNION query if any materialized derived table that
    is part of the UNION is a recursive common table expression (see Recursive Common Table
    Expressions).
  - Conditions containing nondeterministic expressions cannot be pushed down to a derived table.
- The derived table cannot use a LIMIT clause.
- · Conditions containing subqueries cannot be pushed down.
- The optimization cannot be used if the derived table is an inner table of an outer join.
- If a materialized derived table is a common table expression, conditions are not pushed down to it if it is referenced multiple times.
- Conditions using parameters can be pushed down if the condition is of the form <a href="derived\_column">derived\_column</a>
   ?. If a derived column in an outer <a href="https://www.where.condition">where condition is an expression having a ? in the underlying derived table, this condition cannot be pushed down.
- For a query in which the condition is on the tables of a view created using ALGORITHM=TEMPTABLE instead of on the view itself, the multiple equality is not recognized at resolution, and thus the condition cannot be not pushed down. This because, when optimizing a query, condition pushdown takes place during resolution phase while multiple equality propagation occurs during optimization.

This is not an issue in such cases for a view using ALGORITHM=MERGE, where the equality can be propagated and the condition pushed down.

• A condition cannot be pushed down if the derived table's SELECT list contain any assignments to user variables.

## 10.2.3 Optimizing INFORMATION\_SCHEMA Queries

Applications that monitor databases may make frequent use of INFORMATION\_SCHEMA tables. To write queries for these tables most efficiently, use the following general guidelines:

- Try to query only INFORMATION\_SCHEMA tables that are views on data dictionary tables.
- Try to query only for static metadata. Selecting columns or using retrieval conditions for dynamic metadata along with static metadata adds overhead to process the dynamic metadata.



### **Note**

Comparison behavior for database and table names in INFORMATION\_SCHEMA queries might differ from what you expect. For details, see Section 12.8.7, "Using Collation in INFORMATION SCHEMA Searches".

These INFORMATION\_SCHEMA tables are implemented as views on data dictionary tables, so queries on them retrieve information from the data dictionary:

```
CHARACTER_SETS
CHECK_CONSTRAINTS
COLLATIONS
COLLATION_CHARACTER_SET_APPLICABILITY
COLUMNS
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
PARAMETERS
PARTITIONS
REFERENTIAL_CONSTRAINTS
RESOURCE GROUPS
ROUTINES
SCHEMATA
STATISTICS
TABLE_CONSTRAINTS
TRIGGERS
VIEWS
VIEW_ROUTINE_USAGE
VIEW_TABLE_USAGE
```

Some types of values, even for a non-view INFORMATION\_SCHEMA table, are retrieved by lookups from the data dictionary. This includes values such as database and table names, table types, and storage engines.

Some INFORMATION\_SCHEMA tables contain columns that provide table statistics:

```
STATISTICS.CARDINALITY

TABLES.AUTO_INCREMENT

TABLES.AVG_ROW_LENGTH

TABLES.CHECKSUM

TABLES.CHECK_TIME

TABLES.CREATE_TIME

TABLES.DATA_FREE

TABLES.DATA_LENGTH

TABLES.INDEX_LENGTH

TABLES.MAX_DATA_LENGTH

TABLES.TABLE_ROWS

TABLES.UPDATE_TIME
```

Those columns represent dynamic table metadata; that is, information that changes as table contents change.

By default, MySQL retrieves cached values for those columns from the <code>mysql.index\_stats</code> and <code>mysql.innodb\_table\_stats</code> dictionary tables when the columns are queried, which is more efficient than retrieving statistics directly from the storage engine. If cached statistics are not available or have expired, MySQL retrieves the latest statistics from the storage engine and caches them in the <code>mysql.index\_stats</code> and <code>mysql.innodb\_table\_stats</code> dictionary tables. Subsequent queries retrieve the cached statistics until the cached statistics expire. A server restart or the first opening of the <code>mysql.index\_stats</code> and <code>mysql.innodb\_table\_stats</code> tables do not update cached statistics automatically.

The information\_schema\_stats\_expiry session variable defines the period of time before cached statistics expire. The default is 86400 seconds (24 hours), but the time period can be extended to as much as one year.

To update cached values at any time for a given table, use ANALYZE TABLE.

Querying statistics columns does not store or update statistics in the <code>mysql.index\_stats</code> and <code>mysql.innodb\_table\_stats</code> dictionary tables under these circumstances:

- · When cached statistics have not expired.
- When information\_schema\_stats\_expiry is set to 0.
- When the server is in read\_only, super\_read\_only, transaction\_read\_only, or innodb\_read\_only mode.
- When the query also fetches Performance Schema data.

information\_schema\_stats\_expiry is a session variable, and each client session can define its own expiration value. Statistics that are retrieved from the storage engine and cached by one session are available to other sessions.



### **Note**

If the innodb\_read\_only system variable is enabled, ANALYZE TABLE may fail because it cannot update statistics tables in the data dictionary, which use InnoDB. For ANALYZE TABLE operations that update the key distribution, failure may occur even if the operation updates the table itself (for example, if it is a MyISAM table). To obtain the updated distribution statistics, set information schema stats expiry=0.

For INFORMATION\_SCHEMA tables implemented as views on data dictionary tables, indexes on the underlying data dictionary tables permit the optimizer to construct efficient query execution plans. To see the choices made by the optimizer, use EXPLAIN. To also see the query used by the server to execute an INFORMATION\_SCHEMA query, use SHOW WARNINGS immediately following EXPLAIN.

Consider this statement, which identifies collations for the utf8mb4 character set:

How does the server process that statement? To find out, use EXPLAIN:

```
mysql> EXPLAIN SELECT COLLATION_NAME
      FROM INFORMATION_SCHEMA.COLLATION_CHARACTER_SET_APPLICABILITY
      WHERE CHARACTER_SET_NAME = 'utf8mb4'\G
************************ 1. row *****************
         id: 1
 select_type: SIMPLE
       table: cs
  partitions: NULL
        type: const
possible_keys: PRIMARY, name
         kev: name
     key_len: 194
        ref: const
        rows: 1
    filtered: 100.00
      Extra: Using index
  ******************** 2. row *****************
          id: 1
```

```
select_type: SIMPLE
    table: col
partitions: NULL
    type: ref
possible_keys: character_set_id
    key: character_set_id
    key_len: 8
        ref: const
        rows: 68
    filtered: 100.00
        Extra: NULL
2 rows in set, 1 warning (0.01 sec)
```

To see the query used to satisfy that statement, use SHOW WARNINGS:

As indicated by SHOW WARNINGS, the server handles the query on COLLATION\_CHARACTER\_SET\_APPLICABILITY as a query on the character\_sets and collations data dictionary tables in the mysql system database.

## 10.2.4 Optimizing Performance Schema Queries

Applications that monitor databases may make frequent use of Performance Schema tables. To write queries for these tables most efficiently, take advantage of their indexes. For example, include a WHERE clause that restricts retrieved rows based on comparison to specific values in an indexed column.

Most Performance Schema tables have indexes. Tables that do not are those that normally contain few rows or are unlikely to be queried frequently. Performance Schema indexes give the optimizer access to execution plans other than full table scans. These indexes also improve performance for related objects, such as sys schema views that use those tables.

To see whether a given Performance Schema table has indexes and what they are, use SHOW INDEX or SHOW CREATE TABLE:

```
mysgl> SHOW INDEX FROM performance schema.accounts\G
************************ 1. row ******************
      Table: accounts
  Non_unique: 0
   Key_name: ACCOUNT
Seq_in_index: 1
 Column_name: USER
   Collation: NULL
 Cardinality: NULL
    Sub_part: NULL
     Packed: NULL
       Null: YES
  Index_type: HASH
    Comment:
Index comment:
     Visible: YES
************************* 2. row ******************
      Table: accounts
  Non_unique: 0
   Key_name: ACCOUNT
 Seq_in_index: 2
 Column name: HOST
   Collation: NULL
 Cardinality: NULL
    Sub_part: NULL
```

```
Packed: NULL
        Null: YES
  Index_type: HASH
     Comment:
Index_comment:
     Visible: YES
mysql> SHOW CREATE TABLE performance_schema.rwlock_instances\G
           ************** 1. row **************
      Table: rwlock_instances
Create Table: CREATE TABLE `rwlock_instances` (
  `NAME` varchar(128) NOT NULL,
  `OBJECT_INSTANCE_BEGIN` bigint(20) unsigned NOT NULL,
  `WRITE_LOCKED_BY_THREAD_ID` bigint(20) unsigned DEFAULT NULL,
  `READ_LOCKED_BY_COUNT` int(10) unsigned NOT NULL,
 PRIMARY KEY (`OBJECT_INSTANCE_BEGIN`),
 KEY `NAME` (`NAME`),
 KEY `WRITE_LOCKED_BY_THREAD_ID` (`WRITE_LOCKED_BY_THREAD_ID`)
) ENGINE=PERFORMANCE_SCHEMA DEFAULT CHARSET=utf8mb4_COLLATE=utf8mb4_0900_ai_ci
```

To see the execution plan for a Performance Schema query and whether it uses any indexes, use EXPLAIN:

The EXPLAIN output indicates that the optimizer uses the accounts table ACCOUNT index that comprises the USER and HOST columns.

Performance Schema indexes are virtual: They are a construct of the Performance Schema storage engine and use no memory or disk storage. The Performance Schema reports index information to the optimizer so that it can construct efficient execution plans. The Performance Schema in turn uses optimizer information about what to look for (for example, a particular key value), so that it can perform efficient lookups without building actual index structures. This implementation provides two important benefits:

- It entirely avoids the maintenance cost normally incurred for tables that undergo frequent updates.
- It reduces at an early stage of query execution the amount of data retrieved. For conditions on the indexed columns, the Performance Schema efficiently returns only table rows that satisfy the query conditions. Without an index, the Performance Schema would return all rows in the table, requiring that the optimizer later evaluate the conditions against each row to produce the final result.

Performance Schema indexes are predefined and cannot be dropped, added, or altered.

Performance Schema indexes are similar to hash indexes. For example:

- They are used only for equality comparisons that use the = or <=> operators.
- They are unordered. If a query result must have specific row ordering characteristics, include an ORDER BY clause.

For additional information about hash indexes, see Section 10.3.9, "Comparison of B-Tree and Hash Indexes".

## 10.2.5 Optimizing Data Change Statements

This section explains how to speed up data change statements: INSERT, UPDATE, and DELETE. Traditional OLTP applications and modern web applications typically do many small data change operations, where concurrency is vital. Data analysis and reporting applications typically run data change operations that affect many rows at once, where the main considerations is the I/O to write large amounts of data and keep indexes up-to-date. For inserting and updating large volumes of data (known in the industry as ETL, for "extract-transform-load"), sometimes you use other SQL statements or external commands, that mimic the effects of INSERT, UPDATE, and DELETE statements.

## 10.2.5.1 Optimizing INSERT Statements

To optimize insert speed, combine many small operations into a single large operation. Ideally, you make a single connection, send the data for many new rows at once, and delay all index updates and consistency checking until the very end.

The time required for inserting a row is determined by the following factors, where the numbers indicate approximate proportions:

- Connecting: (3)
- Sending query to server: (2)
- Parsing query: (2)
- Inserting row: (1 x size of row)
- Inserting indexes: (1 x number of indexes)
- Closing: (1)

This does not take into consideration the initial overhead to open tables, which is done once for each concurrently running query.

The size of the table slows down the insertion of indexes by  $\log N$ , assuming B-tree indexes.

You can use the following methods to speed up inserts:

- If you are inserting many rows from the same client at the same time, use INSERT statements with
  multiple VALUES lists to insert several rows at a time. This is considerably faster (many times faster
  in some cases) than using separate single-row INSERT statements. If you are adding data to a
  nonempty table, you can tune the bulk\_insert\_buffer\_size variable to make data insertion
  even faster. See Section 7.1.8, "Server System Variables".
- When loading a table from a text file, use LOAD DATA. This is usually 20 times faster than using INSERT statements. See Section 15.2.9, "LOAD DATA Statement".
- Take advantage of the fact that columns have default values. Insert values explicitly only when the
  value to be inserted differs from the default. This reduces the parsing that MySQL must do and
  improves the insert speed.
- See Section 10.5.5, "Bulk Data Loading for InnoDB Tables" for tips specific to InnoDB tables.
- See Section 10.6.2, "Bulk Data Loading for MyISAM Tables" for tips specific to MyISAM tables.

## 10.2.5.2 Optimizing UPDATE Statements

An update statement is optimized like a SELECT query with the additional overhead of a write. The speed of the write depends on the amount of data being updated and the number of indexes that are updated. Indexes that are not changed do not get updated.

Another way to get fast updates is to delay updates and then do many updates in a row later. Performing multiple updates together is much quicker than doing one at a time if you lock the table.

For a MyISAM table that uses dynamic row format, updating a row to a longer total length may split the row. If you do this often, it is very important to use OPTIMIZE TABLE occasionally. See Section 15.7.3.4, "OPTIMIZE TABLE Statement".

### 10.2.5.3 Optimizing DELETE Statements

The time required to delete individual rows in a MyISAM table is exactly proportional to the number of indexes. To delete rows more quickly, you can increase the size of the key cache by increasing the key\_buffer\_size system variable. See Section 7.1.1, "Configuring the Server".

To delete all rows from a MyISAM table, TRUNCATE TABLE *tbl\_name* is faster than DELETE FROM *tbl\_name*. Truncate operations are not transaction-safe; an error occurs when attempting one in the course of an active transaction or active table lock. See Section 15.1.37, "TRUNCATE TABLE Statement".

# 10.2.6 Optimizing Database Privileges

The more complex your privilege setup, the more overhead applies to all SQL statements. Simplifying the privileges established by GRANT statements enables MySQL to reduce permission-checking overhead when clients execute statements. For example, if you do not grant any table-level or column-level privileges, the server need not ever check the contents of the tables\_priv and columns\_priv tables. Similarly, if you place no resource limits on any accounts, the server does not have to perform resource counting. If you have a very high statement-processing load, consider using a simplified grant structure to reduce permission-checking overhead.

## 10.2.7 Other Optimization Tips

This section lists a number of miscellaneous tips for improving query processing speed:

- If your application makes several database requests to perform related updates, combining the
  statements into a stored routine can help performance. Similarly, if your application computes a
  single result based on several column values or large volumes of data, combining the computation
  into a loadable function can help performance. The resulting fast database operations are then
  available to be reused by other queries, applications, and even code written in different programming
  languages. See Section 27.2, "Using Stored Routines" and Adding Functions to MySQL for more
  information.
- To fix any compression issues that occur with ARCHIVE tables, use OPTIMIZE TABLE. See Section 18.5, "The ARCHIVE Storage Engine".
- If possible, classify reports as "live" or as "statistical", where data needed for statistical reports is created only from summary tables that are generated periodically from the live data.
- If you have data that does not conform well to a rows-and-columns table structure, you can pack and store data into a BLOB column. In this case, you must provide code in your application to pack and unpack information, but this might save I/O operations to read and write the sets of related values.
- With Web servers, store images and other binary assets as files, with the path name stored in the database rather than the file itself. Most Web servers are better at caching files than database contents, so using files is generally faster. (Although you must handle backups and storage issues yourself in this case.)
- If you need really high speed, look at the low-level MySQL interfaces. For example, by accessing the MySQL InnoDB or MyISAM storage engine directly, you could get a substantial speed increase compared to using the SQL interface.

Similarly, for databases using the NDBCLUSTER storage engine, you may wish to investigate possible use of the NDB API (see MySQL NDB Cluster API Developer Guide).

 Replication can provide a performance benefit for some operations. You can distribute client retrievals among replicas to split up the load. To avoid slowing down the source while making backups, you can make backups using a replica. See Chapter 19, Replication.

# 10.3 Optimization and Indexes

The best way to improve the performance of SELECT operations is to create indexes on one or more of the columns that are tested in the query. The index entries act like pointers to the table rows, allowing the query to quickly determine which rows match a condition in the WHERE clause, and retrieve the other column values for those rows. All MySQL data types can be indexed.

Although it can be tempting to create an indexes for every possible column used in a query, unnecessary indexes waste space and waste time for MySQL to determine which indexes to use. Indexes also add to the cost of inserts, updates, and deletes because each index must be updated. You must find the right balance to achieve fast queries using the optimal set of indexes.

# 10.3.1 How MySQL Uses Indexes

Indexes are used to find rows with specific column values quickly. Without an index, MySQL must begin with the first row and then read through the entire table to find the relevant rows. The larger the table, the more this costs. If the table has an index for the columns in question, MySQL can quickly determine the position to seek to in the middle of the data file without having to look at all the data. This is much faster than reading every row sequentially.

Most MySQL indexes (PRIMARY KEY, UNIQUE, INDEX, and FULLTEXT) are stored in B-trees. Exceptions: Indexes on spatial data types use R-trees; MEMORY tables also support hash indexes; Innobb uses inverted lists for FULLTEXT indexes.

In general, indexes are used as described in the following discussion. Characteristics specific to hash indexes (as used in MEMORY tables) are described in Section 10.3.9, "Comparison of B-Tree and Hash Indexes".

MySQL uses indexes for these operations:

- To find the rows matching a WHERE clause quickly.
- To eliminate rows from consideration. If there is a choice between multiple indexes, MySQL normally uses the index that finds the smallest number of rows (the most selective index).
- If the table has a multiple-column index, any leftmost prefix of the index can be used by the optimizer to look up rows. For example, if you have a three-column index on (col1, col2, col3), you have indexed search capabilities on (col1), (col1, col2), and (col1, col2, col3). For more information, see Section 10.3.6, "Multiple-Column Indexes".
- To retrieve rows from other tables when performing joins. MySQL can use indexes on columns more efficiently if they are declared as the same type and size. In this context, VARCHAR and CHAR are considered the same if they are declared as the same size. For example, VARCHAR(10) and CHAR(10) are the same size, but VARCHAR(10) and CHAR(15) are not.

For comparisons between nonbinary string columns, both columns should use the same character set. For example, comparing a utf8mb4 column with a latin1 column precludes use of an index.

Comparison of dissimilar columns (comparing a string column to a temporal or numeric column, for example) may prevent use of indexes if values cannot be compared directly without conversion. For a given value such as 1 in the numeric column, it might compare equal to any number of values in the string column such as '1', ' 1', '00001', or '01.e1'. This rules out use of any indexes for the string column.

• To find the MIN() or MAX() value for a specific indexed column key\_col. This is optimized by a preprocessor that checks whether you are using WHERE key\_part\_N = constant on all key

parts that occur before  $key\_col$  in the index. In this case, MySQL does a single key lookup for each MIN() or MAX() expression and replaces it with a constant. If all expressions are replaced with constants, the query returns at once. For example:

```
SELECT MIN(key_part2), MAX(key_part2)
FROM tbl_name WHERE key_part1=10;
```

- To sort or group a table if the sorting or grouping is done on a leftmost prefix of a usable index (for example, ORDER BY key\_part1, key\_part2). If all key parts are followed by DESC, the key is read in reverse order. (Or, if the index is a descending index, the key is read in forward order.) See Section 10.2.1.16, "ORDER BY Optimization", Section 10.2.1.17, "GROUP BY Optimization", and Section 10.3.13, "Descending Indexes".
- In some cases, a query can be optimized to retrieve values without consulting the data rows. (An
  index that provides all the necessary results for a query is called a covering index.) If a query uses
  from a table only columns that are included in some index, the selected values can be retrieved from
  the index tree for greater speed:

```
SELECT key_part3 FROM tbl_name
WHERE key_part1=1
```

Indexes are less important for queries on small tables, or big tables where report queries process most or all of the rows. When a query needs to access most of the rows, reading sequentially is faster than working through an index. Sequential reads minimize disk seeks, even if not all the rows are needed for the query. See Section 10.2.1.23, "Avoiding Full Table Scans" for details.

# 10.3.2 Primary Key Optimization

The primary key for a table represents the column or set of columns that you use in your most vital queries. It has an associated index, for fast query performance. Query performance benefits from the NOT NULL optimization, because it cannot include any NULL values. With the InnoDB storage engine, the table data is physically organized to do ultra-fast lookups and sorts based on the primary key column or columns.

If your table is big and important, but does not have an obvious column or set of columns to use as a primary key, you might create a separate column with auto-increment values to use as the primary key. These unique IDs can serve as pointers to corresponding rows in other tables when you join tables using foreign keys.

# 10.3.3 SPATIAL Index Optimization

MySQL permits creation of SPATIAL indexes on NOT NULL geometry-valued columns (see Section 13.4.10, "Creating Spatial Indexes"). The optimizer checks the SRID attribute for indexed columns to determine which spatial reference system (SRS) to use for comparisons, and uses calculations appropriate to the SRS. (Prior to MySQL 8.4, the optimizer performs comparisons of SPATIAL index values using Cartesian calculations; the results of such operations are undefined if the column contains values with non-Cartesian SRIDs.)

For comparisons to work properly, each column in a SPATIAL index must be SRID-restricted. That is, the column definition must include an explicit SRID attribute, and all column values must have the same SRID.

The optimizer considers SPATIAL indexes only for SRID-restricted columns:

- Indexes on columns restricted to a Cartesian SRID enable Cartesian bounding box computations.
- Indexes on columns restricted to a geographic SRID enable geographic bounding box computations.

The optimizer ignores SPATIAL indexes on columns that have no SRID attribute (and thus are not SRID-restricted). MySQL still maintains such indexes, as follows:

- They are updated for table modifications (INSERT, UPDATE, DELETE, and so forth). Updates occur
  as though the index was Cartesian, even though the column might contain a mix of Cartesian and
  geographical values.
- They exist only for backward compatibility (for example, the ability to perform a dump in MySQL 8.2 and restore in MySQL 8.3). Because SPATIAL indexes on columns that are not SRID-restricted are of no use to the optimizer, each such column should be modified:
  - Verify that all values within the column have the same SRID. To determine the SRIDs contained in a geometry column col\_name, use the following query:

```
SELECT DISTINCT ST_SRID(col_name) FROM tbl_name;
```

If the query returns more than one row, the column contains a mix of SRIDs. In that case, modify its contents so all values have the same SRID.

- Redefine the column to have an explicit SRID attribute.
- Recreate the SPATIAL index.

# 10.3.4 Foreign Key Optimization

If a table has many columns, and you query many different combinations of columns, it might be efficient to split the less-frequently used data into separate tables with a few columns each, and relate them back to the main table by duplicating the numeric ID column from the main table. That way, each small table can have a primary key for fast lookups of its data, and you can query just the set of columns that you need using a join operation. Depending on how the data is distributed, the queries might perform less I/O and take up less cache memory because the relevant columns are packed together on disk. (To maximize performance, queries try to read as few data blocks as possible from disk; tables with only a few columns can fit more rows in each data block.)

## 10.3.5 Column Indexes

The most common type of index involves a single column, storing copies of the values from that column in a data structure, allowing fast lookups for the rows with the corresponding column values. The B-tree data structure lets the index quickly find a specific value, a set of values, or a range of values, corresponding to operators such as =, >,  $\leq$ , BETWEEN, IN, and so on, in a WHERE clause.

The maximum number of indexes per table and the maximum index length is defined per storage engine. See Chapter 17, *The InnoDB Storage Engine*, and Chapter 18, *Alternative Storage Engines*. All storage engines support at least 16 indexes per table and a total index length of at least 256 bytes. Most storage engines have higher limits.

For additional information about column indexes, see Section 15.1.15, "CREATE INDEX Statement".

- Index Prefixes
- FULLTEXT Indexes
- · Spatial Indexes
- Indexes in the MEMORY Storage Engine

### **Index Prefixes**

With  $col_name(N)$  syntax in an index specification for a string column, you can create an index that uses only the first N characters of the column. Indexing only a prefix of column values in this way can make the index file much smaller. When you index a BLOB or TEXT column, you *must* specify a prefix length for the index. For example:

```
CREATE TABLE test (blob_col BLOB, INDEX(blob_col(10)));
```

Prefixes can be up to 767 bytes long for InnoDB tables that use the REDUNDANT or COMPACT row format. The prefix length limit is 3072 bytes for InnoDB tables that use the DYNAMIC or COMPRESSED row format. For MyISAM tables, the prefix length limit is 1000 bytes.



#### Note

Prefix limits are measured in bytes, whereas the prefix length in CREATE TABLE, ALTER TABLE, and CREATE INDEX statements is interpreted as number of characters for nonbinary string types (CHAR, VARCHAR, TEXT) and number of bytes for binary string types (BINARY, VARBINARY, BLOB). Take this into account when specifying a prefix length for a nonbinary string column that uses a multibyte character set.

If a search term exceeds the index prefix length, the index is used to exclude non-matching rows, and the remaining rows are examined for possible matches.

For additional information about index prefixes, see Section 15.1.15, "CREATE INDEX Statement".

### **FULLTEXT Indexes**

FULLTEXT indexes are used for full-text searches. Only the InnoDB and MyISAM storage engines support FULLTEXT indexes and only for CHAR, VARCHAR, and TEXT columns. Indexing always takes place over the entire column and column prefix indexing is not supported. For details, see Section 14.9, "Full-Text Search Functions".

Optimizations are applied to certain kinds of FULLTEXT queries against single InnoDB tables. Queries with these characteristics are particularly efficient:

- FULLTEXT queries that only return the document ID, or the document ID and the search rank.
- FULLTEXT queries that sort the matching rows in descending order of score and apply a LIMIT clause to take the top N matching rows. For this optimization to apply, there must be no WHERE clauses and only a single ORDER BY clause in descending order.
- FULLTEXT queries that retrieve only the COUNT(\*) value of rows matching a search term, with no additional WHERE clauses. Code the WHERE clause as WHERE MATCH(text) AGAINST ('other\_text'), without any > 0 comparison operator.

For queries that contain full-text expressions, MySQL evaluates those expressions during the optimization phase of query execution. The optimizer does not just look at full-text expressions and make estimates, it actually evaluates them in the process of developing an execution plan.

An implication of this behavior is that EXPLAIN for full-text queries is typically slower than for non-full-text queries for which no expression evaluation occurs during the optimization phase.

EXPLAIN for full-text queries may show Select tables optimized away in the Extra column due to matching occurring during optimization; in this case, no table access need occur during later execution.

#### **Spatial Indexes**

You can create indexes on spatial data types. MyISAM and InnoDB support R-tree indexes on spatial types. Other storage engines use B-trees for indexing spatial types (except for ARCHIVE, which does not support spatial type indexing).

### **Indexes in the MEMORY Storage Engine**

The MEMORY storage engine uses HASH indexes by default, but also supports BTREE indexes.

# 10.3.6 Multiple-Column Indexes

MySQL can create composite indexes (that is, indexes on multiple columns). An index may consist of up to 16 columns. For certain data types, you can index a prefix of the column (see Section 10.3.5, "Column Indexes").

MySQL can use multiple-column indexes for queries that test all the columns in the index, or queries that test just the first column, the first two columns, the first three columns, and so on. If you specify the columns in the right order in the index definition, a single composite index can speed up several kinds of queries on the same table.

A multiple-column index can be considered a sorted array, the rows of which contain values that are created by concatenating the values of the indexed columns.



#### Note

As an alternative to a composite index, you can introduce a column that is "hashed" based on information from other columns. If this column is short, reasonably unique, and indexed, it might be faster than a "wide" index on many columns. In MySQL, it is very easy to use this extra column:

```
SELECT * FROM tbl_name
WHERE hash_col=MD5(CONCAT(val1,val2))
AND col1=val1 AND col2=val2;
```

Suppose that a table has the following specification:

The name index is an index over the <code>last\_name</code> and <code>first\_name</code> columns. The index can be used for lookups in queries that specify values in a known range for combinations of <code>last\_name</code> and <code>first\_name</code> values. It can also be used for queries that specify just a <code>last\_name</code> value because that column is a leftmost prefix of the index (as described later in this section). Therefore, the <code>name</code> index is used for lookups in the following queries:

```
SELECT * FROM test WHERE last_name='Jones';

SELECT * FROM test
  WHERE last_name='Jones' AND first_name='John';

SELECT * FROM test
  WHERE last_name='Jones'
  AND (first_name='John' OR first_name='Jon');

SELECT * FROM test
  WHERE last_name='Jones'
  AND first_name >='M' AND first_name < 'N';</pre>
```

However, the name index is *not* used for lookups in the following queries:

```
SELECT * FROM test WHERE first_name='John';

SELECT * FROM test

WHERE last_name='Jones' OR first_name='John';
```

Suppose that you issue the following **SELECT** statement:

```
SELECT * FROM tbl_name
WHERE coll=val1 AND col2=val2;
```

If a multiple-column index exists on coll and col2, the appropriate rows can be fetched directly. If separate single-column indexes exist on coll and col2, the optimizer attempts to use the Index

Merge optimization (see Section 10.2.1.3, "Index Merge Optimization"), or attempts to find the most restrictive index by deciding which index excludes more rows and using that index to fetch the rows.

If the table has a multiple-column index, any leftmost prefix of the index can be used by the optimizer to look up rows. For example, if you have a three-column index on (col1, col2, col3), you have indexed search capabilities on (col1), (col1, col2), and (col1, col2, col3).

MySQL cannot use the index to perform lookups if the columns do not form a leftmost prefix of the index. Suppose that you have the SELECT statements shown here:

```
SELECT * FROM tbl_name WHERE coll=val1;
SELECT * FROM tbl_name WHERE coll=val2;

SELECT * FROM tbl_name WHERE col2=val2;
SELECT * FROM tbl_name WHERE col2=val2 AND col3=val3;
```

If an index exists on (col1, col2, col3), only the first two queries use the index. The third and fourth queries do involve indexed columns, but do not use an index to perform lookups because (col2) and (col2, col3) are not leftmost prefixes of (col1, col2, col3).

## 10.3.7 Verifying Index Usage

Always check whether all your queries really use the indexes that you have created in the tables. Use the EXPLAIN statement, as described in Section 10.8.1, "Optimizing Queries with EXPLAIN".

## 10.3.8 InnoDB and MyISAM Index Statistics Collection

Storage engines collect statistics about tables for use by the optimizer. Table statistics are based on value groups, where a value group is a set of rows with the same key prefix value. For optimizer purposes, an important statistic is the average value group size.

MySQL uses the average value group size in the following ways:

- To estimate how many rows must be read for each ref access
- To estimate how many rows a partial join produces, that is, the number of rows produced by an operation of the form

```
(...) JOIN tbl_name ON tbl_name.key = expr
```

As the average value group size for an index increases, the index is less useful for those two purposes because the average number of rows per lookup increases: For the index to be good for optimization purposes, it is best that each index value target a small number of rows in the table. When a given index value yields a large number of rows, the index is less useful and MySQL is less likely to use it.

The average value group size is related to table cardinality, which is the number of value groups. The SHOW INDEX statement displays a cardinality value based on N/S, where N is the number of rows in the table and S is the average value group size. That ratio yields an approximate number of value groups in the table.

For a join based on the <=> comparison operator, NULL is not treated differently from any other value: NULL <=> NULL, just as N <=> N for any other N.

However, for a join based on the = operator, NULL is different from non-NULL values: expr1 = expr2 is not true when expr1 or expr2 (or both) are NULL. This affects ref accesses for comparisons of the form  $tb1\_name.key = expr$ : MySQL does not access the table if the current value of expr is NULL, because the comparison cannot be true.

For = comparisons, it does not matter how many NULL values are in the table. For optimization purposes, the relevant value is the average size of the non-NULL value groups. However, MySQL does not currently enable that average size to be collected or used.

For InnoDB and MyISAM tables, you have some control over collection of table statistics by means of the innodb\_stats\_method and myisam\_stats\_method system variables, respectively. These variables have three possible values, which differ as follows:

• When the variable is set to nulls\_equal, all NULL values are treated as identical (that is, they all form a single value group).

If the NULL value group size is much higher than the average non-NULL value group size, this method skews the average value group size upward. This makes index appear to the optimizer to be less useful than it really is for joins that look for non-NULL values. Consequently, the nulls\_equal method may cause the optimizer not to use the index for ref accesses when it should.

• When the variable is set to nulls\_unequal, NULL values are not considered the same. Instead, each NULL value forms a separate value group of size 1.

If you have many <code>NULL</code> values, this method skews the average value group size downward. If the average non-<code>NULL</code> value group size is large, counting <code>NULL</code> values each as a group of size 1 causes the optimizer to overestimate the value of the index for joins that look for non-<code>NULL</code> values. Consequently, the <code>nulls\_unequal</code> method may cause the optimizer to use this index for <code>ref</code> lookups when other methods may be better.

• When the variable is set to nulls\_ignored, NULL values are ignored.

If you tend to use many joins that use <=> rather than =, NULL values are not special in comparisons and one NULL is equal to another. In this case, nulls\_equal is the appropriate statistics method.

The innodb\_stats\_method system variable has a global value; the myisam\_stats\_method system variable has both global and session values. Setting the global value affects statistics collection for tables from the corresponding storage engine. Setting the session value affects statistics collection only for the current client connection. This means that you can force a table's statistics to be regenerated with a given method without affecting other clients by setting the session value of myisam\_stats\_method.

To regenerate MyISAM table statistics, you can use any of the following methods:

- Execute myisamchk --stats\_method=method\_name --analyze
- Change the table to cause its statistics to go out of date (for example, insert a row and then delete it), and then set myisam\_stats\_method and issue an ANALYZE TABLE statement

Some caveats regarding the use of innodb\_stats\_method and myisam\_stats\_method:

- You can force table statistics to be collected explicitly, as just described. However, MySQL may also collect statistics automatically. For example, if during the course of executing statements for a table, some of those statements modify the table, MySQL may collect statistics. (This may occur for bulk inserts or deletes, or some ALTER TABLE statements, for example.) If this happens, the statistics are collected using whatever value innodb\_stats\_method or myisam\_stats\_method has at the time. Thus, if you collect statistics using one method, but the system variable is set to the other method when a table's statistics are collected automatically later, the other method is used.
- There is no way to tell which method was used to generate statistics for a given table.
- These variables apply only to InnoDB and MyISAM tables. Other storage engines have only one method for collecting table statistics. Usually it is closer to the nulls\_equal method.

# 10.3.9 Comparison of B-Tree and Hash Indexes

Understanding the B-tree and hash data structures can help predict how different queries perform on different storage engines that use these data structures in their indexes, particularly for the MEMORY storage engine that lets you choose B-tree or hash indexes.

- B-Tree Index Characteristics
- Hash Index Characteristics

#### **B-Tree Index Characteristics**

A B-tree index can be used for column comparisons in expressions that use the =, >, >=, <, <=, or BETWEEN operators. The index also can be used for LIKE comparisons if the argument to LIKE is a constant string that does not start with a wildcard character. For example, the following SELECT statements use indexes:

```
SELECT * FROM tbl_name WHERE key_col LIKE 'Patrick%';
SELECT * FROM tbl_name WHERE key_col LIKE 'Pat%_ck%';
```

In the first statement, only rows with 'Patrick'  $<= key\_col < 'Patricl'$  are considered. In the second statement, only rows with 'Pat'  $<= key\_col < 'Pau'$  are considered.

The following SELECT statements do not use indexes:

```
SELECT * FROM tbl_name WHERE key_col LIKE '%Patrick%';
SELECT * FROM tbl_name WHERE key_col LIKE other_col;
```

In the first statement, the LIKE value begins with a wildcard character. In the second statement, the LIKE value is not a constant.

If you use ... LIKE '%string%' and string is longer than three characters, MySQL uses the *Turbo Boyer-Moore algorithm* to initialize the pattern for the string and then uses this pattern to perform the search more quickly.

A search using col\_name IS NULL employs indexes if col\_name is indexed.

Any index that does not span all AND levels in the WHERE clause is not used to optimize the query. In other words, to be able to use an index, a prefix of the index must be used in every AND group.

The following WHERE clauses use indexes:

```
... WHERE index_part1=1 AND index_part2=2 AND other_column=3

/* index = 1 OR index = 2 */

... WHERE index=1 OR A=10 AND index=2

/* optimized like "index_part1='hello'" */

... WHERE index_part1='hello' AND index_part3=5

/* Can use index on index1 but not on index2 or index3 */

... WHERE index1=1 AND index2=2 OR index1=3 AND index3=3;
```

These WHERE clauses do not use indexes:

```
/* index_part1 is not used */
... WHERE index_part2=1 AND index_part3=2

/* Index is not used in both parts of the WHERE clause */
... WHERE index=1 OR A=10

/* No index spans all rows */
... WHERE index_part1=1 OR index_part2=10
```

Sometimes MySQL does not use an index, even if one is available. One circumstance under which this occurs is when the optimizer estimates that using the index would require MySQL to access a very large percentage of the rows in the table. (In this case, a table scan is likely to be much faster because it requires fewer seeks.) However, if such a query uses LIMIT to retrieve only some of the rows, MySQL uses an index anyway, because it can much more quickly find the few rows to return in the result.

#### **Hash Index Characteristics**

Hash indexes have somewhat different characteristics from those just discussed:

- They are used only for equality comparisons that use the = or <=> operators (but are *very* fast). They are not used for comparison operators such as < that find a range of values. Systems that rely on this type of single-value lookup are known as "key-value stores"; to use MySQL for such applications, use hash indexes wherever possible.
- The optimizer cannot use a hash index to speed up ORDER BY operations. (This type of index cannot be used to search for the next entry in order.)
- MySQL cannot determine approximately how many rows there are between two values (this is used
  by the range optimizer to decide which index to use). This may affect some queries if you change a
  MyISAM or InnobB table to a hash-indexed MEMORY table.
- Only whole keys can be used to search for a row. (With a B-tree index, any leftmost prefix of the key can be used to find rows.)

#### 10.3.10 Use of Index Extensions

InnoDB automatically extends each secondary index by appending the primary key columns to it. Consider this table definition:

```
CREATE TABLE t1 (
   i1 INT NOT NULL DEFAULT 0,
   i2 INT NOT NULL DEFAULT 0,
   d DATE DEFAULT NULL,
   PRIMARY KEY (i1, i2),
   INDEX k_d (d)
) ENGINE = InnoDB;
```

This table defines the primary key on columns (i1, i2). It also defines a secondary index  $k_d$  on column (d), but internally InnoDB extends this index and treats it as columns (d, i1, i2).

The optimizer takes into account the primary key columns of the extended secondary index when determining how and whether to use that index. This can result in more efficient query execution plans and better performance.

The optimizer can use extended secondary indexes for ref, range, and index\_merge index access, for Loose Index Scan access, for join and sorting optimization, and for MIN()/MAX() optimization.

The following example shows how execution plans are affected by whether the optimizer uses extended secondary indexes. Suppose that t1 is populated with these rows:

```
INSERT INTO t1 VALUES
(1, 1, '1998-01-01'), (1, 2, '1999-01-01'),
(1, 3, '2000-01-01'), (1, 4, '2001-01-01'),
(1, 5, '2002-01-01'), (2, 1, '1998-01-01'),
(2, 2, '1999-01-01'), (2, 3, '2000-01-01'),
(2, 4, '2001-01-01'), (2, 5, '2002-01-01'),
(3, 1, '1998-01-01'), (3, 2, '1999-01-01'),
(3, 3, '2000-01-01'), (3, 4, '2001-01-01'),
(3, 5, '2002-01-01'), (4, 1, '1998-01-01'),
(4, 2, '1999-01-01'), (4, 3, '2000-01-01'),
(4, 4, '2001-01-01'), (5, 2, '1999-01-01'),
(5, 1, '1998-01-01'), (5, 2, '1999-01-01'),
(5, 3, '2000-01-01'), (5, 4, '2001-01-01'),
(5, 5, '2002-01-01');
```

Now consider this query:

```
EXPLAIN SELECT COUNT(*) FROM t1 WHERE i1 = 3 AND d = '2000-01-01'
```

The execution plan depends on whether the extended index is used.

When the optimizer does not consider index extensions, it treats the index  $k_d$  as only (d). EXPLAIN for the query produces this result:

```
mysql> EXPLAIN SELECT COUNT(*) FROM t1 WHERE i1 = 3 AND d = '2000-01-01'\G
****************************
        id: 1
        select_type: SIMPLE
            table: t1
            type: ref
possible_keys: PRIMARY,k_d
            key: k_d
        key_len: 4
            ref: const
        rows: 5
        Extra: Using where; Using index
```

When the optimizer takes index extensions into account, it treats  $k_d$  as (d, i1, i2). In this case, it can use the leftmost index prefix (d, i1) to produce a better execution plan:

```
mysql> EXPLAIN SELECT COUNT(*) FROM t1 WHERE i1 = 3 AND d = '2000-01-01'\G
**************************
        id: 1
        select_type: SIMPLE
            table: t1
            type: ref
possible_keys: PRIMARY,k_d
            key: k_d
        key_len: 8
            ref: const,const
        rows: 1
        Extra: Using index
```

In both cases, key indicates that the optimizer uses secondary index k\_d but the EXPLAIN output shows these improvements from using the extended index:

- key\_len goes from 4 bytes to 8 bytes, indicating that key lookups use columns d and i1, not just d.
- The ref value changes from const to const, const because the key lookup uses two key parts, not one.
- The rows count decreases from 5 to 1, indicating that InnoDB should need to examine fewer rows to produce the result.
- The Extra value changes from Using where; Using index to Using index. This means that rows can be read using only the index, without consulting columns in the data row.

Differences in optimizer behavior for use of extended indexes can also be seen with SHOW STATUS:

```
FLUSH TABLE t1;
FLUSH STATUS;
SELECT COUNT(*) FROM t1 WHERE i1 = 3 AND d = '2000-01-01';
SHOW STATUS LIKE 'handler_read%'
```

The preceding statements include FLUSH TABLES and FLUSH STATUS to flush the table cache and clear the status counters.

Without index extensions, SHOW STATUS produces this result:

+----+

With index extensions, SHOW STATUS produces this result. The Handler\_read\_next value decreases from 5 to 1, indicating more efficient use of the index:

The use\_index\_extensions flag of the optimizer\_switch system variable permits control over whether the optimizer takes the primary key columns into account when determining how to use an InnoDB table's secondary indexes. By default, use\_index\_extensions is enabled. To check whether disabling use of index extensions can improve performance, use this statement:

```
SET optimizer_switch = 'use_index_extensions=off';
```

Use of index extensions by the optimizer is subject to the usual limits on the number of key parts in an index (16) and the maximum key length (3072 bytes).

# 10.3.11 Optimizer Use of Generated Column Indexes

MySQL supports indexes on generated columns. For example:

```
CREATE TABLE t1 (f1 INT, gc INT AS (f1 + 1) STORED, INDEX (gc));
```

The generated column, gc, is defined as the expression f1 + 1. The column is also indexed and the optimizer can take that index into account during execution plan construction. In the following query, the WHERE clause refers to gc and the optimizer considers whether the index on that column yields a more efficient plan:

```
SELECT * FROM t1 WHERE gc > 9;
```

The optimizer can use indexes on generated columns to generate execution plans, even in the absence of direct references in queries to those columns by name. This occurs if the WHERE, ORDER BY, or GROUP BY clause refers to an expression that matches the definition of some indexed generated column. The following query does not refer directly to gc but does use an expression that matches the definition of gc:

```
SELECT * FROM t1 WHERE f1 + 1 > 9;
```

The optimizer recognizes that the expression fl+1 matches the definition of gc and that gc is indexed, so it considers that index during execution plan construction. You can see this using EXPLAIN:

```
Extra: Using index condition
```

In effect, the optimizer has replaced the expression £1 + 1 with the name of the generated column that matches the expression. That is also apparent in the rewritten query available in the extended EXPLAIN information displayed by SHOW WARNINGS:

```
mysql> SHOW WARNINGS\G

********************************
Level: Note
   Code: 1003

Message: /* select#1 */ select `test`.`t1`.`f1` AS `f1`,`test`.`t1`.`gc`
        AS `gc` from `test`.`t1` where (`test`.`t1`.`gc` > 9)
```

The following restrictions and conditions apply to the optimizer's use of generated column indexes:

- For a query expression to match a generated column definition, the expression must be identical and it must have the same result type. For example, if the generated column expression is £1 + 1, the optimizer does not recognize a match if the query uses 1 + £1, or if £1 + 1 (an integer expression) is compared with a string.
- The optimization applies to these operators: =, <, <=, >, >=, BETWEEN, and IN().

For operators other than BETWEEN and IN(), either operand can be replaced by a matching generated column. For BETWEEN and IN(), only the first argument can be replaced by a matching generated column, and the other arguments must have the same result type. BETWEEN and IN() are not yet supported for comparisons involving JSON values.

- The generated column must be defined as an expression that contains at least a function call or one of the operators mentioned in the preceding item. The expression cannot consist of a simple reference to another column. For example, gc INT AS (f1) STORED consists only of a column reference, so indexes on gc are not considered.
- For comparisons of strings to indexed generated columns that compute a value from a JSON function that returns a quoted string, JSON\_UNQUOTE() is needed in the column definition to remove the extra quotes from the function value. (For direct comparison of a string to the function result, the JSON comparator handles quote removal, but this does not occur for index lookups.) For example, instead of writing a column definition like this:

```
doc_name TEXT AS (JSON_EXTRACT(jdoc, '$.name')) STORED
```

Write it like this:

```
doc_name TEXT AS (JSON_UNQUOTE(JSON_EXTRACT(jdoc, '$.name'))) STORED
```

With the latter definition, the optimizer can detect a match for both of these comparisons:

```
... WHERE JSON_EXTRACT(jdoc, '$.name') = 'some_string' ...
... WHERE JSON_UNQUOTE(JSON_EXTRACT(jdoc, '$.name')) = 'some_string' ...
```

Without JSON\_UNQUOTE() in the column definition, the optimizer detects a match only for the first of those comparisons.

• If the optimizer picks the wrong index, an index hint can be used to disable it and force the optimizer to make a different choice.

## 10.3.12 Invisible Indexes

MySQL supports invisible indexes; that is, indexes that are not used by the optimizer. The feature applies to indexes other than primary keys (either explicit or implicit).

Indexes are visible by default. To control visibility explicitly for a new index, use a VISIBLE or INVISIBLE keyword as part of the index definition for CREATE TABLE, CREATE INDEX, or ALTER TABLE:

```
CREATE TABLE t1 (
    i INT,
    j INT,
    k INT,
    INDEX i_idx (i) INVISIBLE
) ENGINE = InnoDB;
CREATE INDEX j_idx ON t1 (j) INVISIBLE;
ALTER TABLE t1 ADD INDEX k_idx (k) INVISIBLE;
```

To alter the visibility of an existing index, use a VISIBLE or INVISIBLE keyword with the ALTER TABLE ... ALTER INDEX operation:

```
ALTER TABLE t1 ALTER INDEX i_idx INVISIBLE;
ALTER TABLE t1 ALTER INDEX i_idx VISIBLE;
```

Information about whether an index is visible or invisible is available from the Information Schema STATISTICS table or SHOW INDEX output. For example:

Invisible indexes make it possible to test the effect of removing an index on query performance, without making a destructive change that must be undone should the index turn out to be required. Dropping and re-adding an index can be expensive for a large table, whereas making it invisible and visible are fast, in-place operations.

If an index made invisible actually is needed or used by the optimizer, there are several ways to notice the effect of its absence on queries for the table:

- Errors occur for queries that include index hints that refer to the invisible index.
- Performance Schema data shows an increase in workload for affected queries.
- Queries have different EXPLAIN execution plans.
- Queries appear in the slow query log that did not appear there previously.

The use\_invisible\_indexes flag of the optimizer\_switch system variable controls whether the optimizer uses invisible indexes for query execution plan construction. If the flag is off (the default), the optimizer ignores invisible indexes (the same behavior as prior to the introduction of this flag). If the flag is on, invisible indexes remain invisible but the optimizer takes them into account for execution plan construction.

Using the SET\_VAR optimizer hint to update the value of optimizer\_switch temporarily, you can enable invisible indexes for the duration of a single query only, like this:

```
ref: NULL
       rows: 2
    filtered: 100.00
       Extra: Using index condition
mysql> EXPLAIN SELECT i, j FROM t1 WHERE j >= 50\G
 id: 1
 select_type: SIMPLE
       table: t1
  partitions: NULL
       type: ALL
possible_keys: NULL
        key: NULL
     key_len: NULL
        ref: NULL
       rows: 5
    filtered: 33.33
       Extra: Using where
```

Index visibility does not affect index maintenance. For example, an index continues to be updated per changes to table rows, and a unique index prevents insertion of duplicates into a column, regardless of whether the index is visible or invisible.

A table with no explicit primary key may still have an effective implicit primary key if it has any UNIQUE indexes on NOT NULL columns. In this case, the first such index places the same constraint on table rows as an explicit primary key and that index cannot be made invisible. Consider the following table definition:

```
CREATE TABLE t2 (
i INT NOT NULL,
j INT NOT NULL,
UNIQUE j_idx (j)
) ENGINE = InnoDB;
```

The definition includes no explicit primary key, but the index on NOT NULL column j places the same constraint on rows as a primary key and cannot be made invisible:

```
mysql> ALTER TABLE t2 ALTER INDEX j_idx INVISIBLE;
ERROR 3522 (HY000): A primary key index cannot be invisible.
```

Now suppose that an explicit primary key is added to the table:

```
ALTER TABLE t2 ADD PRIMARY KEY (i);
```

The explicit primary key cannot be made invisible. In addition, the unique index on j no longer acts as an implicit primary key and as a result can be made invisible:

```
mysql> ALTER TABLE t2 ALTER INDEX j_idx INVISIBLE;
Query OK, 0 rows affected (0.03 sec)
```

# 10.3.13 Descending Indexes

MySQL supports descending indexes: DESC in an index definition is no longer ignored but causes storage of key values in descending order. Previously, indexes could be scanned in reverse order but at a performance penalty. A descending index can be scanned in forward order, which is more efficient. Descending indexes also make it possible for the optimizer to use multiple-column indexes when the most efficient scan order mixes ascending order for some columns and descending order for others.

Consider the following table definition, which contains two columns and four two-column index definitions for the various combinations of ascending and descending indexes on the columns:

```
CREATE TABLE t (
cl INT, c2 INT,
INDEX idxl (cl ASC, c2 ASC),
```

```
INDEX idx2 (c1 ASC, c2 DESC),
INDEX idx3 (c1 DESC, c2 ASC),
INDEX idx4 (c1 DESC, c2 DESC)
);
```

The table definition results in four distinct indexes. The optimizer can perform a forward index scan for each of the ORDER BY clauses and need not use a filesort operation:

```
ORDER BY c1 ASC, c2 ASC -- optimizer can use idx1
ORDER BY c1 DESC, c2 DESC -- optimizer can use idx4
ORDER BY c1 ASC, c2 DESC -- optimizer can use idx2
ORDER BY c1 DESC, c2 ASC -- optimizer can use idx3
```

Use of descending indexes is subject to these conditions:

- Descending indexes are supported only for the InnoDB storage engine, with these limitations:
  - Change buffering is not supported for a secondary index if the index contains a descending index key column or if the primary key includes a descending index column.
  - The InnoDB SQL parser does not use descending indexes. For InnoDB full-text search, this means that the index required on the FTS\_DOC\_ID column of the indexed table cannot be defined as a descending index. For more information, see Section 17.6.2.4, "InnoDB Full-Text Indexes".
- Descending indexes are supported for all data types for which ascending indexes are available.
- Descending indexes are supported for ordinary (nongenerated) and generated columns (both VIRTUAL and STORED).
- DISTINCT can use any index containing matching columns, including descending key parts.
- Indexes that have descending key parts are not used for MIN()/MAX() optimization of queries that invoke aggregate functions but do not have a GROUP BY clause.
- Descending indexes are supported for BTREE but not HASH indexes. Descending indexes are not supported for FULLTEXT or SPATIAL indexes.

Explicitly specified ASC and DESC designators for HASH, FULLTEXT, and SPATIAL indexes results in an error.

You can see in the **Extra** column of the output of **EXPLAIN** that the optimizer is able to use a descending index, as shown here:

```
mysql> CREATE TABLE t1 (
   -> a INT,
   -> b INT,
   -> INDEX a desc b asc (a DESC, b ASC)
   -> );
mysql> EXPLAIN SELECT * FROM t1 ORDER BY a ASC\G
  ****** 1. row ********
         id: 1
 select_type: SIMPLE
       table: t1
  partitions: NULL
        type: index
possible_keys: NULL
        key: a_desc_b_asc
     key_len: 10
         ref: NULL
        rows: 1
    filtered: 100.00
       Extra: Backward index scan; Using index
```

In EXPLAIN FORMAT=TREE output, use of a descending index is indicated by the addition of (reverse) following the name of the index, like this:

```
mysql> EXPLAIN FORMAT=TREE SELECT * FROM t1 ORDER BY a ASC\G
**************************
EXPLAIN: -> Index scan on t1 using a_desc_b_asc (reverse) (cost=0.35 rows=1)
```

See also EXPLAIN Extra Information.

# 10.3.14 Indexed Lookups from TIMESTAMP Columns

Temporal values are stored in TIMESTAMP columns as UTC values, and values inserted into and retrieved from TIMESTAMP columns are converted between the session time zone and UTC. (This is the same type of conversion performed by the CONVERT\_TZ() function. If the session time zone is UTC, there is effectively no time zone conversion.)

Due to conventions for local time zone changes such as Daylight Saving Time (DST), conversions between UTC and non-UTC time zones are not one-to-one in both directions. UTC values that are distinct may not be distinct in another time zone. The following example shows distinct UTC values that become identical in a non-UTC time zone:



#### Note

To use named time zones such as <code>'MET'</code> or <code>'Europe/Amsterdam'</code>, the time zone tables must be properly set up. For instructions, see Section 7.1.15, "MySQL Server Time Zone Support".

You can see that the two distinct UTC values are the same when converted to the 'MET' time zone. This phenomenon can lead to different results for a given TIMESTAMP column query, depending on whether the optimizer uses an index to execute the query.

Suppose that a query selects values from the table shown earlier using a WHERE clause to search the ts column for a single specific value such as a user-provided timestamp literal:

```
SELECT ts FROM tstable
WHERE ts = 'literal';
```

Suppose further that the query executes under these conditions:

• The session time zone is not UTC and has a DST shift. For example:

```
SET time_zone = 'MET';
```

- Unique UTC values stored in the TIMESTAMP column are not unique in the session time zone due to DST shifts. (The example shown earlier illustrates how this can occur.)
- The query specifies a search value that is within the hour of entry into DST in the session time zone.

Under those conditions, the comparison in the WHERE clause occurs in different ways for nonindexed and indexed lookups and leads to different results:

If there is no index or the optimizer cannot use it, comparisons occur in the session time zone. The
optimizer performs a table scan in which it retrieves each ts column value, converts it from UTC
to the session time zone, and compares it to the search value (also interpreted in the session time
zone):

Because the stored ts values are converted to the session time zone, it is possible for the query to return two timestamp values that are distinct as UTC values but equal in the session time zone: One value that occurs before the DST shift when clocks are changed, and one value that was occurs after the DST shift.

If there is a usable index, comparisons occur in UTC. The optimizer performs an index scan, first
converting the search value from the session time zone to UTC, then comparing the result to the
UTC index entries:

In this case, the (converted) search value is matched only to index entries, and because the index entries for the distinct stored UTC values are also distinct, the search value can match only one of them.

Due to different optimizer operation for nonindexed and indexed lookups, the query produces different results in each case. The result from the nonindexed lookup returns all values that match in the session time zone. The indexed lookup cannot do so:

- It is performed within the storage engine, which knows only about UTC values.
- For the two distinct session time zone values that map to the same UTC value, the indexed lookup matches only the corresponding UTC index entry and returns only a single row.

In the preceding discussion, the data set stored in tstable happens to consist of distinct UTC values. In such cases, all index-using queries of the form shown match at most one index entry.

If the index is not UNIQUE, it is possible for the table (and the index) to store multiple instances of a given UTC value. For example, the ts column might contain multiple instances of the UTC value '2018-10-28 00:30:00'. In this case, the index-using query would return each of them (converted to the MET value '2018-10-28 02:30:00' in the result set). It remains true that index-using queries match the converted search value to a single value in the UTC index entries, rather than matching multiple UTC values that convert to the search value in the session time zone.

If it is important to return all ts values that match in the session time zone, the workaround is to suppress use of the index with an IGNORE INDEX hint:

```
mysql> SELECT ts FROM tstable
    IGNORE INDEX (ts)
    WHERE ts = '2018-10-28 02:30:00';
+-----+
```

The same lack of one-to-one mapping for time zone conversions in both directions occurs in other contexts as well, such as conversions performed with the FROM\_UNIXTIME() and UNIX TIMESTAMP() functions. See Section 14.7, "Date and Time Functions".

# 10.4 Optimizing Database Structure

In your role as a database designer, look for the most efficient way to organize your schemas, tables, and columns. As when tuning application code, you minimize I/O, keep related items together, and plan ahead so that performance stays high as the data volume increases. Starting with an efficient database design makes it easier for team members to write high-performing application code, and makes the database likely to endure as applications evolve and are rewritten.

# 10.4.1 Optimizing Data Size

Design your tables to minimize their space on the disk. This can result in huge improvements by reducing the amount of data written to and read from disk. Smaller tables normally require less main memory while their contents are being actively processed during query execution. Any space reduction for table data also results in smaller indexes that can be processed faster.

MySQL supports many different storage engines (table types) and row formats. For each table, you can decide which storage and indexing method to use. Choosing the proper table format for your application can give you a big performance gain. See Chapter 17, *The InnoDB Storage Engine*, and Chapter 18, *Alternative Storage Engines*.

You can get better performance for a table and minimize storage space by using the techniques listed here:

- Table Columns
- Row Format
- Indexes
- Joins
- Normalization

#### **Table Columns**

- Use the most efficient (smallest) data types possible. MySQL has many specialized types that save disk space and memory. For example, use the smaller integer types if possible to get smaller tables. MEDIUMINT is often a better choice than INT because a MEDIUMINT column uses 25% less space.
- Declare columns to be NOT NULL if possible. It makes SQL operations faster, by enabling better use of indexes and eliminating overhead for testing whether each value is NULL. You also save some storage space, one bit per column. If you really need NULL values in your tables, use them. Just avoid the default setting that allows NULL values in every column.

#### **Row Format**

• InnoDB tables are created using the DYNAMIC row format by default. To use a row format other than DYNAMIC, configure innodb\_default\_row\_format, or specify the ROW\_FORMAT option explicitly in a CREATE TABLE or ALTER TABLE statement.

The compact family of row formats, which includes COMPACT, DYNAMIC, and COMPRESSED, decreases row storage space at the cost of increasing CPU use for some operations. If your

workload is a typical one that is limited by cache hit rates and disk speed it is likely to be faster. If it is a rare case that is limited by CPU speed, it might be slower.

The compact family of row formats also optimizes CHAR column storage when using a variable-length character set such as  $\mathtt{utf8mb3}$  or  $\mathtt{utf8mb4}$ . With  $\mathtt{ROW\_FORMAT=REDUNDANT}$ ,  $\mathtt{CHAR}(N)$  occupies  $N \times \text{the maximum byte length of the character set.}$  Many languages can be written primarily using single-byte  $\mathtt{utf8mb3}$  or  $\mathtt{utf8mb4}$  characters, so a fixed storage length often wastes space. With the compact family of rows formats,  $\mathtt{InnoDB}$  allocates a variable amount of storage in the range of  $N \times \text{the maximum byte length of the character set for these columns by stripping trailing spaces. The minimum storage length is <math>N \times \text{the maximum byte sto facilitate in-place updates in typical cases. For more information, see Section 17.10, "InnoDB Row Formats".$ 

- To minimize space even further by storing table data in compressed form, specify ROW\_FORMAT=COMPRESSED when creating InnoDB tables, or run the myisampack command on an existing MyISAM table. (InnoDB compressed tables are readable and writable, while MyISAM compressed tables are read-only.)
- For MyISAM tables, if you do not have any variable-length columns (VARCHAR, TEXT, or BLOB columns), a fixed-size row format is used. This is faster but may waste some space. See Section 18.2.3, "MyISAM Table Storage Formats". You can hint that you want to have fixed length rows even if you have VARCHAR columns with the CREATE TABLE option ROW\_FORMAT=FIXED.

#### **Indexes**

- The primary index of a table should be as short as possible. This makes identification of each row easy and efficient. For InnoDB tables, the primary key columns are duplicated in each secondary index entry, so a short primary key saves considerable space if you have many secondary indexes.
- Create only the indexes that you need to improve query performance. Indexes are good for retrieval, but slow down insert and update operations. If you access a table mostly by searching on a combination of columns, create a single composite index on them rather than a separate index for each column. The first part of the index should be the column most used. If you *always* use many columns when selecting from the table, the first column in the index should be the one with the most duplicates, to obtain better compression of the index.
- If it is very likely that a long string column has a unique prefix on the first number of characters, it is better to index only this prefix, using MySQL's support for creating an index on the leftmost part of the column (see Section 15.1.15, "CREATE INDEX Statement"). Shorter indexes are faster, not only because they require less disk space, but because they also give you more hits in the index cache, and thus fewer disk seeks. See Section 7.1.1, "Configuring the Server".

#### **Joins**

- In some circumstances, it can be beneficial to split into two a table that is scanned very often. This is especially true if it is a dynamic-format table and it is possible to use a smaller static format table that can be used to find the relevant rows when scanning the table.
- Declare columns with identical information in different tables with identical data types, to speed up joins based on the corresponding columns.
- Keep column names simple, so that you can use the same name across different tables and simplify
  join queries. For example, in a table named customer, use a column name of name instead of
  customer\_name. To make your names portable to other SQL servers, consider keeping them
  shorter than 18 characters.

### **Normalization**

• Normally, try to keep all data nonredundant (observing what is referred to in database theory as *third normal form*). Instead of repeating lengthy values such as names and addresses, assign them unique IDs, repeat these IDs as needed across multiple smaller tables, and join the tables in queries by referencing the IDs in the join clause.

If speed is more important than disk space and the maintenance costs of keeping multiple copies
of data, for example in a business intelligence scenario where you analyze all the data from large
tables, you can relax the normalization rules, duplicating information or creating summary tables to
gain more speed.

# 10.4.2 Optimizing MySQL Data Types

### 10.4.2.1 Optimizing for Numeric Data

- For unique IDs or other values that can be represented as either strings or numbers, prefer numeric columns to string columns. Since large numeric values can be stored in fewer bytes than the corresponding strings, it is faster and takes less memory to transfer and compare them.
- If you are using numeric data, it is faster in many cases to access information from a database (using
  a live connection) than to access a text file. Information in the database is likely to be stored in a
  more compact format than in the text file, so accessing it involves fewer disk accesses. You also
  save code in your application because you can avoid parsing the text file to find line and column
  boundaries.

## 10.4.2.2 Optimizing for Character and String Types

For character and string columns, follow these guidelines:

- Use binary collation order for fast comparison and sort operations, when you do not need languagespecific collation features. You can use the BINARY operator to use binary collation within a particular query.
- When comparing values from different columns, declare those columns with the same character set and collation wherever possible, to avoid string conversions while running the query.
- For column values less than 8KB in size, use binary VARCHAR instead of BLOB. The GROUP BY and ORDER BY clauses can generate temporary tables, and these temporary tables can use the MEMORY storage engine if the original table does not contain any BLOB columns.
- If a table contains string columns such as name and address, but many queries do not retrieve
  those columns, consider splitting the string columns into a separate table and using join queries
  with a foreign key when necessary. When MySQL retrieves any value from a row, it reads a data
  block containing all the columns of that row (and possibly other adjacent rows). Keeping each row
  small, with only the most frequently used columns, allows more rows to fit in each data block. Such
  compact tables reduce disk I/O and memory usage for common queries.
- When you use a randomly generated value as a primary key in an InnoDB table, prefix it with an ascending value such as the current date and time if possible. When consecutive primary values are physically stored near each other, InnoDB can insert and retrieve them faster.
- See Section 10.4.2.1, "Optimizing for Numeric Data" for reasons why a numeric column is usually preferable to an equivalent string column.

### 10.4.2.3 Optimizing for BLOB Types

- When storing a large blob containing textual data, consider compressing it first. Do not use this
  technique when the entire table is compressed by InnoDB or MyISAM.
- For a table with several columns, to reduce memory requirements for queries that do not use the BLOB column, consider splitting the BLOB column into a separate table and referencing it with a join query when needed.
- Since the performance requirements to retrieve and display a BLOB value might be very different from other data types, you could put the BLOB-specific table on a different storage device or even a

separate database instance. For example, to retrieve a BLOB might require a large sequential disk read that is better suited to a traditional hard drive than to an SSD device.

- See Section 10.4.2.2, "Optimizing for Character and String Types" for reasons why a binary VARCHAR column is sometimes preferable to an equivalent BLOB column.
- Rather than testing for equality against a very long text string, you can store a hash of the column value in a separate column, index that column, and test the hashed value in queries. (Use the MD5() or CRC32() function to produce the hash value.) Since hash functions can produce duplicate results for different inputs, you still include a clause AND blob\_column = long\_string\_value in the query to guard against false matches; the performance benefit comes from the smaller, easily scanned index for the hashed values.

# 10.4.3 Optimizing for Many Tables

Some techniques for keeping individual queries fast involve splitting data across many tables. When the number of tables runs into the thousands or even millions, the overhead of dealing with all these tables becomes a new performance consideration.

## 10.4.3.1 How MySQL Opens and Closes Tables

When you execute a mysqladmin status command, you should see something like this:

```
Uptime: 426 Running threads: 1 Questions: 11082
Reloads: 1 Open tables: 12
```

The Open tables value of 12 can be somewhat puzzling if you have fewer than 12 tables.

MySQL is multithreaded, so there may be many clients issuing queries for a given table simultaneously. To minimize the problem with multiple client sessions having different states on the same table, the table is opened independently by each concurrent session. This uses additional memory but normally increases performance. With MyISAM tables, one extra file descriptor is required for the data file for each client that has the table open. (By contrast, the index file descriptor is shared between all sessions.)

The table\_open\_cache and max\_connections system variables affect the maximum number of files the server keeps open. If you increase one or both of these values, you may run up against a limit imposed by your operating system on the per-process number of open file descriptors. Many operating systems permit you to increase the open-files limit, although the method varies widely from system to system. Consult your operating system documentation to determine whether it is possible to increase the limit and how to do so.

table\_open\_cache is related to max\_connections. For example, for 200 concurrent running connections, specify a table cache size of at least  $200 \, * \, N$ , where N is the maximum number of tables per join in any of the queries which you execute. You must also reserve some extra file descriptors for temporary tables and files.

Make sure that your operating system can handle the number of open file descriptors implied by the table\_open\_cache setting. If table\_open\_cache is set too high, MySQL may run out of file descriptors and exhibit symptoms such as refusing connections or failing to perform queries.

Also take into account that the MyISAM storage engine needs two file descriptors for each unique open table. To increase the number of file descriptors available to MySQL, set the open\_files\_limit system variable. See Section B.3.2.16, "File Not Found and Similar Errors".

The cache of open tables is kept at a level of table\_open\_cache entries. The server autosizes the cache size at startup. To set the size explicitly, set the table\_open\_cache system variable at startup. MySQL may temporarily open more tables than this to execute queries, as described later in this section.

MySQL closes an unused table and removes it from the table cache under the following circumstances:

- When the cache is full and a thread tries to open a table that is not in the cache.
- When the cache contains more than table\_open\_cache entries and a table in the cache is no longer being used by any threads.
- When a table-flushing operation occurs. This happens when someone issues a FLUSH TABLES statement or executes a mysgladmin flush-tables or mysgladmin refresh command.

When the table cache fills up, the server uses the following procedure to locate a cache entry to use:

- Tables not currently in use are released, beginning with the table least recently used.
- If a new table must be opened, but the cache is full and no tables can be released, the cache is temporarily extended as necessary. When the cache is in a temporarily extended state and a table goes from a used to unused state, the table is closed and released from the cache.

A MyISAM table is opened for each concurrent access. This means the table needs to be opened twice if two threads access the same table or if a thread accesses the table twice in the same query (for example, by joining the table to itself). Each concurrent open requires an entry in the table cache. The first open of any MyISAM table takes two file descriptors: one for the data file and one for the index file. Each additional use of the table takes only one file descriptor for the data file. The index file descriptor is shared among all threads.

If you are opening a table with the HANDLER <code>tbl\_name</code> OPEN statement, a dedicated table object is allocated for the thread. This table object is not shared by other threads and is not closed until the thread calls <code>HANDLER tbl\_name</code> CLOSE or the thread terminates. When this happens, the table is put back in the table cache (if the cache is not full). See Section 15.2.5, "HANDLER Statement".

To determine whether your table cache is too small, check the Opened\_tables status variable, which indicates the number of table-opening operations since the server started:

If the value is very large or increases rapidly, even when you have not issued many FLUSH TABLES statements, increase the table\_open\_cache value at server startup.

#### 10.4.3.2 Disadvantages of Creating Many Tables in the Same Database

If you have many MyISAM tables in the same database directory, open, close, and create operations are slow. If you execute SELECT statements on many different tables, there is a little overhead when the table cache is full, because for every table that has to be opened, another must be closed. You can reduce this overhead by increasing the number of entries permitted in the table cache.

# 10.4.4 Internal Temporary Table Use in MySQL

In some cases, the server creates internal temporary tables while processing statements. Users have no direct control over when this occurs.

The server creates temporary tables under conditions such as these:

- Evaluation of UNION statements, with some exceptions described later.
- Evaluation of some views, such those that use the TEMPTABLE algorithm, UNION, or aggregation.
- Evaluation of derived tables (see Section 15.2.15.8, "Derived Tables").

- Evaluation of common table expressions (see Section 15.2.20, "WITH (Common Table Expressions)").
- Tables created for subquery or semijoin materialization (see Section 10.2.2, "Optimizing Subqueries, Derived Tables, View References, and Common Table Expressions").
- Evaluation of statements that contain an ORDER BY clause and a different GROUP BY clause, or for which the ORDER BY or GROUP BY contains columns from tables other than the first table in the join queue.
- Evaluation of DISTINCT combined with ORDER BY may require a temporary table.
- For queries that use the SQL\_SMALL\_RESULT modifier, MySQL uses an in-memory temporary table, unless the query also contains elements (described later) that require on-disk storage.
- To evaluate INSERT ... SELECT statements that select from and insert into the same table, MySQL creates an internal temporary table to hold the rows from the SELECT, then inserts those rows into the target table. See Section 15.2.7.1, "INSERT ... SELECT Statement".
- Evaluation of multiple-table UPDATE statements.
- Evaluation of GROUP\_CONCAT() or COUNT(DISTINCT) expressions.
- Evaluation of window functions (see Section 14.20, "Window Functions") uses temporary tables as necessary.

To determine whether a statement requires a temporary table, use EXPLAIN and check the Extra column to see whether it says Using temporary (see Section 10.8.1, "Optimizing Queries with EXPLAIN"). EXPLAIN does not necessarily say Using temporary for derived or materialized temporary tables. For statements that use window functions, EXPLAIN with FORMAT=JSON always provides information about the windowing steps. If the windowing functions use temporary tables, it is indicated for each step.

Some query conditions prevent the use of an in-memory temporary table, in which case the server uses an on-disk table instead:

- Presence of a BLOB or TEXT column in the table. The TempTable storage engine, which is the
  default storage engine for in-memory internal temporary tables in MySQL 8.4, supports binary large
  object types. See Internal Temporary Table Storage Engine.
- Presence of any string column with a maximum length larger than 512 (bytes for binary strings, characters for nonbinary strings) in the SELECT list, if UNION or UNION ALL is used.
- The SHOW COLUMNS and DESCRIBE statements use BLOB as the type for some columns, thus the temporary table used for the results is an on-disk table.

The server does not use a temporary table for UNION statements that meet certain qualifications. Instead, it retains from temporary table creation only the data structures necessary to perform result column typecasting. The table is not fully instantiated and no rows are written to or read from it; rows are sent directly to the client. The result is reduced memory and disk requirements, and smaller delay before the first row is sent to the client because the server need not wait until the last query block is executed. EXPLAIN and optimizer trace output reflects this execution strategy: The UNION RESULT query block is not present because that block corresponds to the part that reads from the temporary table.

These conditions qualify a UNION for evaluation without a temporary table:

- The union is union all, not union or union distinct.
- There is no global ORDER BY clause.
- The union is not the top-level query block of an {INSERT | REPLACE} ... SELECT ... statement.

### **Internal Temporary Table Storage Engine**

An internal temporary table can be held in memory and processed by the TempTable or MEMORY storage engine, or stored on disk by the InnoDB storage engine.

#### Storage Engine for In-Memory Internal Temporary Tables

The internal\_tmp\_mem\_storage\_engine variable defines the storage engine used for in-memory internal temporary tables. Permitted values are TempTable (the default) and MEMORY.



#### Note

Configuring a session setting for internal\_tmp\_mem\_storage\_engine requires the SESSION\_VARIABLES\_ADMIN or SYSTEM\_VARIABLES\_ADMIN privilege.

The TempTable storage engine provides efficient storage for VARCHAR and VARBINARY columns, and other binary large object types.

The following variables control TempTable storage engine limits and behavior:

• tmp\_table\_size: Defines the maximum size of any individual in-memory internal temporary table created using the TempTable storage engine. When the limit determined by tmp\_table\_size is reached, MySQL automatically converts the in-memory internal temporary table to an InnoDB on-disk internal temporary table. The default value is 16777216 bytes (16 MiB).

The tmp\_table\_size limit is intended to prevent individual queries from consuming an inordinate amount of global TempTable resources, which can affect the performance of concurrent queries that require such resources. Global TempTable resources are controlled by temptable\_max\_ram and temptable max mmap.

If tmp\_table\_size is less than temptable\_max\_ram, it is not possible for an in-memory temporary table to use more than tmp\_table\_size. If tmp\_table\_size is greater than the sum of temptable\_max\_ram and temptable\_max\_mmap, an in-memory temporary table cannot use more than the sum of the temptable max ram and temptable max mmap limits.

• temptable\_max\_ram: Defines the maximum amount of RAM that can be used by the TempTable storage engine before it starts allocating space from memory-mapped files or before MySQL starts using InnoDB on-disk internal temporary tables, depending on your configuration. If not set explicitly, the value of temptable\_max\_ram is 3% of the total memory available on the server, with a minimum of 1 GB and a maximum of 4 GB.



#### Note

temptable\_max\_ram does not account for the thread-local memory block allocated to each thread that uses the TempTable storage engine. The size of the thread-local memory block depends on the size of the thread's first memory allocation request. If the request is less than 1MB, which it is in most cases, the thread-local memory block size is 1MB. If the request is greater than 1MB, the thread-local memory block is approximately the same size as the initial memory request. The thread-local memory block is held in thread-local storage until thread exit.

• temptable\_use\_mmap: Controls whether the TempTable storage engine allocates space from memory-mapped files or MySQL uses InnoDB on-disk internal temporary tables when the limit determined by temptable\_max\_ram is exceeded. The default value is OFF.



#### Note

temptable\_use\_mmap is deprecated; expect support for it to be removed in a future version of MySQL. Setting temptable\_max\_mmap=0 is equivalent to setting temptable use mmap=OFF.

• temptable\_max\_mmap: Sets the maximum amount of memory the TempTable storage engine is permitted to allocate from memory-mapped files before MySQL starts using InnoDB on-disk internal temporary tables. The default value is 0 (disabled). The limit is intended to address the risk of memory mapped files using too much space in the temporary directory (tmpdir). temptable\_max\_mmap = 0 disables allocation from memory-mapped files, effectively disabling their use, regardless of the value of temptable\_use\_mmap.

Use of memory-mapped files by the TempTable storage engine is governed by these rules:

- Temporary files are created in the directory defined by the tmpdir variable.
- Temporary files are deleted immediately after they are created and opened, and therefore do not remain visible in the tmpdir directory. The space occupied by temporary files is held by the operating system while temporary files are open. The space is reclaimed when temporary files are closed by the TempTable storage engine, or when the mysqld process is shut down.
- Data is never moved between RAM and temporary files, within RAM, or between temporary files.
- New data is stored in RAM if space becomes available within the limit defined by temptable\_max\_ram. Otherwise, new data is stored in temporary files.
- If space becomes available in RAM after some of the data for a table is written to temporary files, it is possible for the remaining table data to be stored in RAM.

When using the MEMORY storage engine for in-memory temporary tables (internal\_tmp\_mem\_storage\_engine=MEMORY), MySQL automatically converts an in-memory temporary table to an on-disk table if it becomes too large. The maximum size of an in-memory temporary table is defined by the tmp\_table\_size or max\_heap\_table\_size value, whichever is smaller. This differs from MEMORY tables explicitly created with CREATE TABLE. For such tables, only the max\_heap\_table\_size variable determines how large a table can grow, and there is no conversion to on-disk format.

### Storage Engine for On-Disk Internal Temporary Tables

MySQL 8.4 uses only the InnoDB storage engine for on-disk internal temporary tables. (The MYISAM storage engine is no longer supported for this purpose.)

InnoDB on-disk internal temporary tables are created in session temporary tablespaces that reside in the data directory by default. For more information, see Section 17.6.3.5, "Temporary Tablespaces".

### **Internal Temporary Table Storage Format**

When in-memory internal temporary tables are managed by the TempTable storage engine, rows that include VARCHAR columns, VARBINARY columns, and other binary large object type columns are represented in memory by an array of cells, with each cell containing a NULL flag, the data length, and a data pointer. Column values are placed in consecutive order after the array, in a single region of memory, without padding. Each cell in the array uses 16 bytes of storage. The same storage format applies when the TempTable storage engine allocates space from memory-mapped files.

When in-memory internal temporary tables are managed by the MEMORY storage engine, fixed-length row format is used. VARCHAR and VARBINARY column values are padded to the maximum column length, in effect storing them as CHAR and BINARY columns.

Internal temporary tables on disk are always managed by InnoDB.

When using the MEMORY storage engine, statements can initially create an in-memory internal temporary table and then convert it to an on-disk table if the table becomes too large. In such cases, better performance might be achieved by skipping the conversion and creating the internal temporary table on disk to begin with. The big\_tables variable can be used to force disk storage of internal temporary tables.

### **Monitoring Internal Temporary Table Creation**

When an internal temporary table is created in memory or on disk, the server increments the Created\_tmp\_tables value. When an internal temporary table is created on disk, the server increments the Created\_tmp\_disk\_tables value. If too many internal temporary tables are created on disk, consider adjusting the engine-specific limits described in Internal Temporary Table Storage Engine.



#### Note

Due to a known limitation, Created\_tmp\_disk\_tables does not count on-disk temporary tables created in memory-mapped files. By default, the TempTable storage engine overflow mechanism creates internal temporary tables in memory-mapped files. See Internal Temporary Table Storage Engine.

The memory/temptable/physical\_ram and memory/temptable/physical\_disk Performance Schema instruments can be used to monitor TempTable space allocation from memory and disk. memory/temptable/physical\_ram reports the amount of allocated RAM. memory/temptable/physical\_disk reports the amount of space allocated from disk when memory-mapped files are used as the TempTable overflow mechanism. If the physical\_disk instrument reports a value other than 0 and memory-mapped files are used as the TempTable overflow mechanism, a TempTable memory limit was reached at some point. Data can be queried in Performance Schema memory summary tables such as memory\_summary\_global\_by\_event\_name. See Section 29.12.20.10, "Memory Summary Tables".

### 10.4.5 Limits on Number of Databases and Tables

MySQL has no limit on the number of databases. The underlying file system may have a limit on the number of directories.

MySQL has no limit on the number of tables. The underlying file system may have a limit on the number of files that represent tables. Individual storage engines may impose engine-specific constraints. Innobb permits up to 4 billion tables.

### 10.4.6 Limits on Table Size

The effective maximum table size for MySQL databases is usually determined by operating system constraints on file sizes, not by MySQL internal limits. For up-to-date information operating system file size limits, refer to the documentation specific to your operating system.

Windows users, please note that FAT and VFAT (FAT32) are *not* considered suitable for production use with MySQL. Use NTFS instead.

If you encounter a full-table error, there are several reasons why it might have occurred:

- The disk might be full.
- You are using InnoDB tables and have run out of room in an InnoDB tablespace file. The maximum tablespace size is also the maximum size for a table. For tablespace size limits, see Section 17.21, "InnoDB Limits".

Generally, partitioning of tables into multiple tablespace files is recommended for tables larger than 1TB in size.

- You have hit an operating system file size limit. For example, you are using MyISAM tables on an
  operating system that supports files only up to 2GB in size and you have hit this limit for the data file
  or index file.
- You are using a MyISAM table and the space required for the table exceeds what is permitted by the
  internal pointer size. MyISAM permits data and index files to grow up to 256TB by default, but this
  limit can be changed up to the maximum permissible size of 65,536TB (256<sup>7</sup> 1 bytes).

If you need a MyISAM table that is larger than the default limit and your operating system supports large files, the CREATE TABLE statement supports AVG\_ROW\_LENGTH and MAX\_ROWS options. See Section 15.1.20, "CREATE TABLE Statement". The server uses these options to determine how large a table to permit.

If the pointer size is too small for an existing table, you can change the options with ALTER TABLE to increase a table's maximum permissible size. See Section 15.1.9, "ALTER TABLE Statement".

```
ALTER TABLE tbl_name MAX_ROWS=100000000 AVG_ROW_LENGTH=nnn;
```

You have to specify AVG\_ROW\_LENGTH only for tables with BLOB or TEXT columns; in this case, MySQL cannot optimize the space required based only on the number of rows.

To change the default size limit for MyISAM tables, set the myisam\_data\_pointer\_size, which sets the number of bytes used for internal row pointers. The value is used to set the pointer size for new tables if you do not specify the MAX\_ROWS option. The value of myisam\_data\_pointer\_size can be from 2 to 7. For example, for tables that use the dynamic storage format, a value of 4 permits tables up to 4GB; a value of 6 permits tables up to 256TB. Tables that use the fixed storage format have a larger maximum data length. For storage format characteristics, see Section 18.2.3, "MyISAM Table Storage Formats".

You can check the maximum data and index sizes by using this statement:

```
SHOW TABLE STATUS FROM db_name LIKE 'tbl_name';
```

You also can use myisamchk -dv /path/to/table-index-file. See Section 15.7.7, "SHOW Statements", or Section 6.6.4, "myisamchk — MyISAM Table-Maintenance Utility".

Other ways to work around file-size limits for MyISAM tables are as follows:

- If your large table is read only, you can use myisampack to compress it. myisampack usually compresses a table by at least 50%, so you can have, in effect, much bigger tables. myisampack also can merge multiple tables into a single table. See Section 6.6.6, "myisampack Generate Compressed, Read-Only MyISAM Tables".
- MySQL includes a MERGE library that enables you to handle a collection of MyISAM tables that have identical structure as a single MERGE table. See Section 18.7, "The MERGE Storage Engine".
- You are using the MEMORY (HEAP) storage engine; in this case you need to increase the value of the max\_heap\_table\_size system variable. See Section 7.1.8, "Server System Variables".

### 10.4.7 Limits on Table Column Count and Row Size

This section describes limits on the number of columns in tables and the size of individual rows.

- · Column Count Limits
- Row Size Limits

#### **Column Count Limits**

MySQL has hard limit of 4096 columns per table, but the effective maximum may be less for a given table. The exact column limit depends on several factors:

- The maximum row size for a table constrains the number (and possibly size) of columns because the total length of all columns cannot exceed this size. See Row Size Limits.
- The storage requirements of individual columns constrain the number of columns that fit within
  a given maximum row size. Storage requirements for some data types depend on factors such
  as storage engine, storage format, and character set. See Section 13.7, "Data Type Storage
  Requirements".

- Storage engines may impose additional restrictions that limit table column count. For example, InnoDB has a limit of 1017 columns per table. See Section 17.21, "InnoDB Limits". For information about other storage engines, see Chapter 18, *Alternative Storage Engines*.
- Functional key parts (see Section 15.1.15, "CREATE INDEX Statement") are implemented as hidden
  virtual generated stored columns, so each functional key part in a table index counts against the
  table total column limit.

#### **Row Size Limits**

The maximum row size for a given table is determined by several factors:

- The internal representation of a MySQL table has a maximum row size limit of 65,535 bytes, even if
  the storage engine is capable of supporting larger rows. BLOB and TEXT columns only contribute 9
  to 12 bytes toward the row size limit because their contents are stored separately from the rest of the
  row.
- The maximum row size for an InnoDB table, which applies to data stored locally within a database page, is slightly less than half a page for 4KB, 8KB, 16KB, and 32KB innodb\_page\_size settings. For example, the maximum row size is slightly less than 8KB for the default 16KB InnoDB page size. For 64KB pages, the maximum row size is slightly less than 16KB. See Section 17.21, "InnoDB Limits".

If a row containing variable-length columns exceeds the InnoDB maximum row size, InnoDB selects variable-length columns for external off-page storage until the row fits within the InnoDB row size limit. The amount of data stored locally for variable-length columns that are stored off-page differs by row format. For more information, see Section 17.10, "InnoDB Row Formats".

- Different storage formats use different amounts of page header and trailer data, which affects the amount of storage available for rows.
  - For information about InnoDB row formats, see Section 17.10, "InnoDB Row Formats".
  - For information about MyISAM storage formats, see Section 18.2.3, "MyISAM Table Storage Formats".

### **Row Size Limit Examples**

• The MySQL maximum row size limit of 65,535 bytes is demonstrated in the following InnoDB and MyISAM examples. The limit is enforced regardless of storage engine, even though the storage engine may be capable of supporting larger rows.

In the following MyISAM example, changing a column to TEXT avoids the 65,535-byte row size limit and permits the operation to succeed because BLOB and TEXT columns only contribute 9 to 12 bytes toward the row size.

The operation succeeds for an InnoDB table because changing a column to TEXT avoids the MySQL 65,535-byte row size limit, and InnoDB off-page storage of variable-length columns avoids the InnoDB row size limit.

• Storage for variable-length columns includes length bytes, which are counted toward the row size. For example, a VARCHAR(255) CHARACTER SET utf8mb3 column takes two bytes to store the length of the value, so each value can take up to 767 bytes.

The statement to create table t1 succeeds because the columns require 32,765 + 2 bytes and 32,766 + 2 bytes, which falls within the maximum row size of 65,535 bytes:

```
mysql> CREATE TABLE t1
      (c1 VARCHAR(32765) NOT NULL, c2 VARCHAR(32766) NOT NULL)
      ENGINE = Innobe CHARACTER SET latin1;
Query OK, 0 rows affected (0.02 sec)
```

The statement to create table  $\pm 2$  fails because, although the column length is within the maximum length of 65,535 bytes, two additional bytes are required to record the length, which causes the row size to exceed 65,535 bytes:

Reducing the column length to 65,533 or less permits the statement to succeed.

• For MyISAM tables, NULL columns require additional space in the row to record whether their values are NULL. Each NULL column takes one bit extra, rounded up to the nearest byte.

The statement to create table t3 fails because MyISAM requires space for NULL columns in addition to the space required for variable-length column length bytes, causing the row size to exceed 65,535 bytes:

For information about InnoDB NULL column storage, see Section 17.10, "InnoDB Row Formats".

• InnoDB restricts row size (for data stored locally within the database page) to slightly less than half a database page for 4KB, 8KB, 16KB, and 32KB innodb\_page\_size settings, and to slightly less than 16KB for 64KB pages.

The statement to create table t4 fails because the defined columns exceed the row size limit for a 16KB InnoDB page.

```
mysql> CREATE TABLE t4 (
c1 CHAR(255),c2 CHAR(255),c3 CHAR(255),
c4 CHAR(255),c5 CHAR(255),c6 CHAR(255),
c7 CHAR(255),c8 CHAR(255),c9 CHAR(255),
```

```
c10 CHAR(255),c11 CHAR(255),c12 CHAR(255),
c13 CHAR(255),c14 CHAR(255),c15 CHAR(255),
c16 CHAR(255),c17 CHAR(255),c18 CHAR(255),
c19 CHAR(255),c20 CHAR(255),c21 CHAR(255),
c22 CHAR(255),c23 CHAR(255),c24 CHAR(255),
c25 CHAR(255),c26 CHAR(255),c27 CHAR(255),
c28 CHAR(255),c29 CHAR(255),c30 CHAR(255),
c31 CHAR(255),c32 CHAR(255),c33 CHAR(255)
) ENGINE=InnoDB ROW_FORMAT=DYNAMIC DEFAULT CHARSET latin1;
ERROR 1118 (42000): Row size too large (> 8126). Changing some columns to TEXT or BLOB may help.
In current row format, BLOB prefix of 0 bytes is stored inline.
```

# 10.5 Optimizing for InnoDB Tables

InnoDB is the storage engine that MySQL customers typically use in production databases where reliability and concurrency are important. InnoDB is the default storage engine in MySQL. This section explains how to optimize database operations for InnoDB tables.

# 10.5.1 Optimizing Storage Layout for InnoDB Tables

Once your data reaches a stable size, or a growing table has increased by tens or some hundreds
of megabytes, consider using the OPTIMIZE TABLE statement to reorganize the table and compact
any wasted space. The reorganized tables require less disk I/O to perform full table scans. This is a
straightforward technique that can improve performance when other techniques such as improving
index usage or tuning application code are not practical.

OPTIMIZE TABLE copies the data part of the table and rebuilds the indexes. The benefits come from improved packing of data within indexes, and reduced fragmentation within the tablespaces and on disk. The benefits vary depending on the data in each table. You may find that there are significant gains for some and not for others, or that the gains decrease over time until you next optimize the table. This operation can be slow if the table is large or if the indexes being rebuilt do not fit into the buffer pool. The first run after adding a lot of data to a table is often much slower than later runs.

- In Innobb, having a long PRIMARY KEY (either a single column with a lengthy value, or several columns that form a long composite value) wastes a lot of disk space. The primary key value for a row is duplicated in all the secondary index records that point to the same row. (See Section 17.6.2.1, "Clustered and Secondary Indexes".) Create an AUTO\_INCREMENT column as the primary key if your primary key is long, or index a prefix of a long VARCHAR column instead of the entire column.
- Use the VARCHAR data type instead of CHAR to store variable-length strings or for columns with many NULL values. A CHAR (N) column always takes N characters to store data, even if the string is shorter or its value is NULL. Smaller tables fit better in the buffer pool and reduce disk I/O.

When using COMPACT row format (the default InnoDB format) and variable-length character sets, such as  $\mathtt{utf8mb4}$  or  $\mathtt{sjis}$ ,  $\mathtt{CHAR}(N)$  columns occupy a variable amount of space, but still at least N bytes.

• For tables that are big, or contain lots of repetitive text or numeric data, consider using COMPRESSED row format. Less disk I/O is required to bring data into the buffer pool, or to perform full table scans. Before making a permanent decision, measure the amount of compression you can achieve by using COMPRESSED versus COMPACT row format.

# 10.5.2 Optimizing InnoDB Transaction Management

To optimize InnoDB transaction processing, find the ideal balance between the performance overhead of transactional features and the workload of your server. For example, an application might encounter performance issues if it commits thousands of times per second, and different performance issues if it commits only every 2-3 hours.

• The default MySQL setting AUTOCOMMIT=1 can impose performance limitations on a busy database server. Where practical, wrap several related data change operations into a single transaction, by issuing SET AUTOCOMMIT=0 or a START TRANSACTION statement, followed by a COMMIT statement after making all the changes.

InnoDB must flush the log to disk at each transaction commit if that transaction made modifications to the database. When each change is followed by a commit (as with the default autocommit setting), the I/O throughput of the storage device puts a cap on the number of potential operations per second.

- Alternatively, for transactions that consist only of a single SELECT statement, turning on AUTOCOMMIT helps InnoDB to recognize read-only transactions and optimize them. See Section 10.5.3, "Optimizing InnoDB Read-Only Transactions" for requirements.
- Avoid performing rollbacks after inserting, updating, or deleting huge numbers of rows. If a big
  transaction is slowing down server performance, rolling it back can make the problem worse,
  potentially taking several times as long to perform as the original data change operations. Killing the
  database process does not help, because the rollback starts again on server startup.

To minimize the chance of this issue occurring:

- Increase the size of the buffer pool so that all the data change changes can be cached rather than immediately written to disk.
- Set innodb\_change\_buffering=all so that update and delete operations are buffered in addition to inserts.
- Consider issuing COMMIT statements periodically during the big data change operation, possibly breaking a single delete or update into multiple statements that operate on smaller numbers of rows.

To get rid of a runaway rollback once it occurs, increase the buffer pool so that the rollback becomes CPU-bound and runs fast, or kill the server and restart with <code>innodb\_force\_recovery=3</code>, as explained in Section 17.18.2, "InnoDB Recovery".

- If you can afford the loss of some of the latest committed transactions if an unexpected exit occurs, you can set the innodb\_flush\_log\_at\_trx\_commit parameter to 0. InnoDB tries to flush the log once per second anyway, although the flush is not guaranteed.
- When rows are modified or deleted, the rows and associated undo logs are not physically removed immediately, or even immediately after the transaction commits. The old data is preserved until transactions that started earlier or concurrently are finished, so that those transactions can access the previous state of modified or deleted rows. Thus, a long-running transaction can prevent Innobe from purging data that was changed by a different transaction.
- When rows are modified or deleted within a long-running transaction, other transactions using the READ COMMITTED and REPEATABLE READ isolation levels have to do more work to reconstruct the older data if they read those same rows.
- When a long-running transaction modifies a table, queries against that table from other transactions
  do not make use of the covering index technique. Queries that normally could retrieve all the result
  columns from a secondary index, instead look up the appropriate values from the table data.

If secondary index pages are found to have a PAGE\_MAX\_TRX\_ID that is too new, or if records in the secondary index are delete-marked, InnoDB may need to look up records using a clustered index.

# 10.5.3 Optimizing InnoDB Read-Only Transactions

InnoDB can avoid the overhead associated with setting up the transaction ID (TRX\_ID field) for transactions that are known to be read-only. A transaction ID is only needed for a transaction that

might perform write operations or locking reads such as SELECT ... FOR UPDATE. Eliminating unnecessary transaction IDs reduces the size of internal data structures that are consulted each time a query or data change statement constructs a read view.

InnoDB detects read-only transactions when:

• The transaction is started with the START TRANSACTION READ ONLY statement. In this case, attempting to make changes to the database (for Innodb, MyISAM, or other types of tables) causes an error, and the transaction continues in read-only state:

```
ERROR 1792 (25006): Cannot execute statement in a READ ONLY transaction.
```

You can still make changes to session-specific temporary tables in a read-only transaction, or issue locking queries for them, because those changes and locks are not visible to any other transaction.

- The autocommit setting is turned on, so that the transaction is guaranteed to be a single statement, and the single statement making up the transaction is a "non-locking" SELECT statement. That is, a SELECT that does not use a FOR UPDATE or LOCK IN SHARED MODE clause.
- The transaction is started without the READ ONLY option, but no updates or statements that explicitly lock rows have been executed yet. Until updates or explicit locks are required, a transaction stays in read-only mode.

Thus, for a read-intensive application such as a report generator, you can tune a sequence of InnoDB queries by grouping them inside START TRANSACTION READ ONLY and COMMIT, or by turning on the autocommit setting before running the SELECT statements, or simply by avoiding any data change statements interspersed with the queries.

For information about START TRANSACTION and autocommit, see Section 15.3.1, "START TRANSACTION, COMMIT, and ROLLBACK Statements".



#### Note

Transactions that qualify as auto-commit, non-locking, and read-only (AC-NL-RO) are kept out of certain internal InnoDB data structures and are therefore not listed in SHOW ENGINE INNODB STATUS output.

# 10.5.4 Optimizing InnoDB Redo Logging

Consider the following guidelines for optimizing redo logging:

• Increase the size of your redo log files. When InnoDB has written redo log files full, it must write the modified contents of the buffer pool to disk in a checkpoint. Small redo log files cause many unnecessary disk writes.

The redo log file size is determined by <code>innodb\_redo\_log\_capacity</code>. <code>InnoDB</code> tries to maintain 32 redo log files of the same size, with each file equal to 1/32 \* <code>innodb\_redo\_log\_capacity</code>. Therefore, changing the <code>innodb\_redo\_log\_capacity</code> setting changes the size of the redo log files.

For information about modifying your redo log file configuration, see Section 17.6.5, "Redo Log".

- Consider increasing the size of the log buffer. A large log buffer enables large transactions to run
  without a need to write the log to disk before the transactions commit. Thus, if you have transactions
  that update, insert, or delete many rows, making the log buffer larger saves disk I/O. Log buffer size
  is configured using the innodb\_log\_buffer\_size configuration option, which can be dynamically
  configured.
- Configure the <code>innodb\_log\_write\_ahead\_size</code> configuration option to avoid "read-on-write". This option defines the write-ahead block size for the redo log. Set <code>innodb\_log\_write\_ahead\_size</code>

to match the operating system or file system cache block size. Read-on-write occurs when redo log blocks are not entirely cached to the operating system or file system due to a mismatch between write-ahead block size for the redo log and operating system or file system cache block size.

Valid values for <code>innodb\_log\_write\_ahead\_size</code> are multiples of the <code>InnoDB</code> log file block size (2<sup>n</sup>). The minimum value is the <code>InnoDB</code> log file block size (512). Write-ahead does not occur when the minimum value is specified. The maximum value is equal to the <code>innodb\_page\_size</code> value. If you specify a value for <code>innodb\_log\_write\_ahead\_size</code> that is larger than the <code>innodb\_page\_size</code> value, the <code>innodb\_log\_write\_ahead\_size</code> setting is truncated to the <code>innodb\_page\_size</code> value.

Setting the innodb\_log\_write\_ahead\_size value too low in relation to the operating system or file system cache block size results in read-on-write. Setting the value too high may have a slight impact on fsync performance for log file writes due to several blocks being written at once.

- MySQL provides dedicated log writer threads for writing redo log records from the log buffer to the
  system buffers and flushing the system buffers to the redo log files. You can enable or disable log
  writer threads using the innodb\_log\_writer\_threads variable. Dedicated log writer threads
  can improve performance on high-concurrency systems, but for low-concurrency systems, disabling
  dedicated log writer threads provides better performance.
- Optimize the use of spin delay by user threads waiting for flushed redo. Spin delay helps reduce latency. During periods of low concurrency, reducing latency may be less of a priority, and avoiding the use of spin delay during these periods may reduce energy consumption. During periods of high concurrency, you may want to avoid expending processing power on spin delay so that it can be used for other work. The following system variables permit setting high and low watermark values that define boundaries for the use of spin delay.
  - innodb\_log\_wait\_for\_flush\_spin\_hwm: Defines the maximum average log flush time beyond which user threads no longer spin while waiting for flushed redo. The default value is 400 microseconds.
  - innodb\_log\_spin\_cpu\_abs\_lwm: Defines the minimum amount of CPU usage below which user threads no longer spin while waiting for flushed redo. The value is expressed as a sum of CPU core usage. For example, The default value of 80 is 80% of a single CPU core. On a system with a multi-core processor, a value of 150 represents 100% usage of one CPU core plus 50% usage of a second CPU core.
  - innodb\_log\_spin\_cpu\_pct\_hwm: Defines the maximum amount of CPU usage above which user threads no longer spin while waiting for flushed redo. The value is expressed as a percentage of the combined total processing power of all CPU cores. The default value is 50%. For example, 100% usage of two CPU cores is 50% of the combined CPU processing power on a server with four CPU cores.

The  $innodb\_log\_spin\_cpu\_pct\_hwm$  configuration option respects processor affinity. For example, if a server has 48 cores but the mysqld process is pinned to only four CPU cores, the other 44 CPU cores are ignored.

# 10.5.5 Bulk Data Loading for InnoDB Tables

These performance tips supplement the general guidelines for fast inserts in Section 10.2.5.1, "Optimizing INSERT Statements".

When importing data into InnoDB, turn off autocommit mode, because it performs a log flush to
disk for every insert. To disable autocommit during your import operation, surround it with SET
autocommit and COMMIT statements:

```
SET autocommit=0;
... SQL import statements ...
COMMIT;
```

The mysqldump option --opt creates dump files that are fast to import into an InnoDB table, even without wrapping them with the SET autocommit and COMMIT statements.

• If you have UNIQUE constraints on secondary keys, you can speed up table imports by temporarily turning off the uniqueness checks during the import session:

```
SET unique_checks=0;
... SQL import statements ...
SET unique_checks=1;
```

For big tables, this saves a lot of disk I/O because InnoDB can use its change buffer to write secondary index records in a batch. Be certain that the data contains no duplicate keys.

• If you have FOREIGN KEY constraints in your tables, you can speed up table imports by turning off the foreign key checks for the duration of the import session:

```
SET foreign_key_checks=0;
... SQL import statements ...
SET foreign_key_checks=1;
```

For big tables, this can save a lot of disk I/O.

• Use the multiple-row INSERT syntax to reduce communication overhead between the client and the server if you need to insert many rows:

```
INSERT INTO yourtable VALUES (1,2), (5,5), ...;
```

This tip is valid for inserts into any table, not just InnoDB tables.

- When doing bulk inserts into tables with auto-increment columns, set innodb\_autoinc\_lock\_mode to 2 (interleaved) instead of 1 (consecutive). See Section 17.6.1.6, "AUTO\_INCREMENT Handling in InnoDB" for details.
- When performing bulk inserts, it is faster to insert rows in PRIMARY KEY order. InnoDB tables use a clustered index, which makes it relatively fast to use data in the order of the PRIMARY KEY. Performing bulk inserts in PRIMARY KEY order is particularly important for tables that do not fit entirely within the buffer pool.
- For optimal performance when loading data into an InnoDB FULLTEXT index, follow this set of steps:
  - 1. Define a column FTS\_DOC\_ID at table creation time, of type BIGINT UNSIGNED NOT NULL, with a unique index named FTS\_DOC\_ID\_INDEX. For example:

```
CREATE TABLE t1 (
    FTS_DOC_ID BIGINT unsigned NOT NULL AUTO_INCREMENT,
    title varchar(255) NOT NULL DEFAULT '',
    text mediumtext NOT NULL,
PRIMARY KEY (`FTS_DOC_ID`)
) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4;

CREATE UNIQUE INDEX FTS_DOC_ID_INDEX on t1(FTS_DOC_ID);
```

- 2. Load the data into the table.
- Create the FULLTEXT index after the data is loaded.



#### Note

When adding FTS\_DOC\_ID column at table creation time, ensure that the FTS\_DOC\_ID column is updated when the FULLTEXT indexed column is updated, as the FTS\_DOC\_ID must increase monotonically with each INSERT or UPDATE. If you choose not to add the FTS\_DOC\_ID at table creation time and have InnoDB manage DOC IDs for you, InnoDB adds the FTS\_DOC\_ID

as a hidden column with the next CREATE FULLTEXT INDEX call. This approach, however, requires a table rebuild which can impact performance.

• If loading data into a *new* MySQL instance, consider disabling redo logging using ALTER INSTANCE {ENABLE | DISABLE} INNODB REDO\_LOG syntax. Disabling redo logging helps speed up data loading by avoiding redo log writes. For more information, see Disabling Redo Logging.



#### 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.

• Use MySQL Shell to import data. MySQL Shell's parallel table import utility util.importTable() provides rapid data import to a MySQL relational table for large data files. MySQL Shell's dump loading utility util.loadDump() also offers parallel load capabilities. See MySQL Shell Utilities.

## 10.5.6 Optimizing InnoDB Queries

To tune queries for InnoDB tables, create an appropriate set of indexes on each table. See Section 10.3.1, "How MySQL Uses Indexes" for details. Follow these guidelines for InnoDB indexes:

- Because each InnoDB table has a primary key (whether you request one or not), specify a set of
  primary key columns for each table, columns that are used in the most important and time-critical
  queries.
- Do not specify too many or too long columns in the primary key, because these column values are duplicated in each secondary index. When an index contains unnecessary data, the I/O to read this data and memory to cache it reduce the performance and scalability of the server.
- Do not create a separate secondary index for each column, because each query can only make use of one index. Indexes on rarely tested columns or columns with only a few different values might not be helpful for any queries. If you have many queries for the same table, testing different combinations of columns, try to create a small number of concatenated indexes rather than a large number of single-column indexes. If an index contains all the columns needed for the result set (known as a covering index), the query might be able to avoid reading the table data at all.
- If an indexed column cannot contain any NULL values, declare it as NOT NULL when you create the table. The optimizer can better determine which index is most effective to use for a query, when it knows whether each column contains NULL values.
- You can optimize single-query transactions for InnoDB tables, using the technique in Section 10.5.3, "Optimizing InnoDB Read-Only Transactions".

# 10.5.7 Optimizing InnoDB DDL Operations

- Many DDL operations on tables and indexes (CREATE, ALTER, and DROP statements) can be performed online. See Section 17.12, "InnoDB and Online DDL" for details.
- Online DDL support for adding secondary indexes means that you can generally speed up the
  process of creating and loading a table and associated indexes by creating the table without
  secondary indexes, then adding secondary indexes after the data is loaded.
- Use TRUNCATE TABLE to empty a table, not DELETE FROM tb1\_name. Foreign key constraints can make a TRUNCATE statement work like a regular DELETE statement, in which case a sequence of commands like DROP TABLE and CREATE TABLE might be fastest.
- Because the primary key is integral to the storage layout of each InnoDB table, and changing the definition of the primary key involves reorganizing the whole table, always set up the primary key as

part of the CREATE TABLE statement, and plan ahead so that you do not need to ALTER or DROP the primary key afterward.

## 10.5.8 Optimizing InnoDB Disk I/O

If you follow best practices for database design and tuning techniques for SQL operations, but your database is still slow due to heavy disk I/O activity, consider these disk I/O optimizations. If the Unix top tool or the Windows Task Manager shows that the CPU usage percentage with your workload is less than 70%, your workload is probably disk-bound.

• Increase buffer pool size

When table data is cached in the InnoDB buffer pool, it can be accessed repeatedly by queries without requiring any disk I/O. Specify the size of the buffer pool with the innodb\_buffer\_pool\_size option. This memory area is important enough that it is typically recommended that innodb\_buffer\_pool\_size is configured to 50 to 75 percent of system memory. For more information see, Section 10.12.3.1, "How MySQL Uses Memory".

· Adjust the flush method

In some versions of GNU/Linux and Unix, flushing files to disk with the Unix fsync() call and similar methods is surprisingly slow. If database write performance is an issue, conduct benchmarks with the innodb flush method parameter set to O DSYNC.

Configure a threshold for operating system flushes

By default, when InnoDB creates a new data file, such as a new log file or tablespace file, the file is fully written to the operating system cache before it is flushed to disk, which can cause a large amount of disk write activity to occur at once. To force smaller, periodic flushes of data from the operating system cache, you can use the <code>innodb\_fsync\_threshold</code> variable to define a threshold value, in bytes. When the byte threshold is reached, the contents of the operating system cache are flushed to disk. The default value of 0 forces the default behavior, which is to flush data to disk only after a file is fully written to the cache.

Specifying a threshold to force smaller, periodic flushes may be beneficial in cases where multiple MySQL instances use the same storage devices. For example, creating a new MySQL instance and its associated data files could cause large surges of disk write activity, impeding the performance of other MySQL instances that use the same storage devices. Configuring a threshold helps avoid such surges in write activity.

• Use fdatasync() instead of fsync()

On platforms that support fdatasync() system calls, the innodb\_use\_fdatasync variable permits using fdatasync() instead of fsync() for operating system flushes. An fdatasync() system call does not flush changes to file metadata unless required for subsequent data retrieval, providing a potential performance benefit.

A subset of innodb\_flush\_method settings such as fsync, O\_DSYNC, and O\_DIRECT use fsync() system calls. The innodb\_use\_fdatasync variable is applicable when using those settings.

• Use a noop or deadline I/O scheduler with native AIO 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 <code>innodb\_use\_native\_aio</code> configuration option, which is enabled by default. With native AIO, the type of I/O scheduler has greater influence on I/O performance. Generally, noop and deadline I/O schedulers are recommended. Conduct benchmarks to determine which I/O scheduler provides the best results for your workload and environment. For more information, see Section 17.8.6, "Using Asynchronous I/O on Linux".

• Use direct I/O on Solaris 10 for x86\_64 architecture

When using the InnoDB storage engine on Solaris 10 for x86\_64 architecture (AMD Opteron), use direct I/O for InnoDB-related files to avoid degradation of InnoDB performance. To use direct I/O for an entire UFS file system used for storing InnoDB-related files, mount it with the forcedirectio option; see mount\_ufs(1M). (The default on Solaris 10/x86\_64 is not to use this option.) To apply direct I/O only to InnoDB file operations rather than the whole file system, set innodb\_flush\_method = O\_DIRECT. With this setting, InnoDB calls directio() instead of fcntl() for I/O to data files (not for I/O to log files).

· Use raw storage for data and log files with Solaris 2.6 or later

When using the InnoDB storage engine with a large innodb\_buffer\_pool\_size value on any release of Solaris 2.6 and up and any platform (sparc/x86/x64/amd64), conduct benchmarks with InnoDB data files and log files on raw devices or on a separate direct I/O UFS file system, using the forcedirectio mount option as described previously. (It is necessary to use the mount option rather than setting innodb\_flush\_method if you want direct I/O for the log files.) Users of the Veritas file system VxFS should use the convosync=direct mount option.

Do not place other MySQL data files, such as those for MyISAM tables, on a direct I/O file system. Executables or libraries *must not* be placed on a direct I/O file system.

Use additional storage devices

Additional storage devices could be used to set up a RAID configuration. For related information, see Section 10.12.1, "Optimizing Disk I/O".

Alternatively, InnoDB tablespace data files and log files can be placed on different physical disks. For more information, refer to the following sections:

- Section 17.8.1, "InnoDB Startup Configuration"
- Section 17.6.1.2, "Creating Tables Externally"
- Creating a General Tablespace
- Section 17.6.1.4, "Moving or Copying InnoDB Tables"
- · Consider non-rotational storage

Non-rotational storage generally provides better performance for random I/O operations; and rotational storage for sequential I/O operations. When distributing data and log files across rotational

and non-rotational storage devices, consider the type of I/O operations that are predominantly performed on each file.

Random I/O-oriented files typically include file-per-table and general tablespace data files, undo tablespace files, and temporary tablespace files. Sequential I/O-oriented files include InnoDB system tablespace files, doublewrite files, and log files such as binary log files and redo log files.

Review settings for the following configuration options when using non-rotational storage:

• innodb\_checksum\_algorithm

The crc32 option uses a faster checksum algorithm and is recommended for fast storage systems.

• innodb\_flush\_neighbors

Optimizes I/O for rotational storage devices. Disable it for non-rotational storage or a mix of rotational and non-rotational storage. It is disabled by default.

• innodb\_idle\_flush\_pct

Permits placing a limit on page flushing during idle periods, which can help extend the life of non-rotational storage devices.

• innodb\_io\_capacity

The default setting of 10000 is generally sufficient.

• innodb\_io\_capacity\_max

The default value of (2 \* innodb\_io\_capacity) is intended for most workloads.

• innodb\_log\_compressed\_pages

If redo logs are on non-rotational storage, consider disabling this option to reduce logging. See Disable logging of compressed pages.

• innodb\_log\_file\_size (deprecated)

If redo logs are on non-rotational storage, configure this option to maximize caching and write combining.

• innodb\_redo\_log\_capacity

If redo logs are on non-rotational storage, configure this option to maximize caching and write combining.

• innodb page size

Consider using a page size that matches the internal sector size of the disk. Early-generation SSD devices often have a 4KB sector size. Some newer devices have a 16KB sector size. The default InnobB page size is 16KB. Keeping the page size close to the storage device block size minimizes the amount of unchanged data that is rewritten to disk.

• binlog\_row\_image

If binary logs are on non-rotational storage and all tables have primary keys, consider setting this option to minimal to reduce logging.

Ensure that TRIM support is enabled for your operating system. It is typically enabled by default.

· Increase I/O capacity to avoid backlogs

If throughput drops periodically because of InnoDB checkpoint operations, consider increasing the value of the innodb\_io\_capacity configuration option. Higher values cause more frequent flushing, avoiding the backlog of work that can cause dips in throughput.

· Lower I/O capacity if flushing does not fall behind

If the system is not falling behind with InnoDB flushing operations, consider lowering the value of the innodb\_io\_capacity configuration option. Typically, you keep this option value as low as practical, but not so low that it causes periodic drops in throughput as mentioned in the preceding bullet. In a typical scenario where you could lower the option value, you might see a combination like this in the output from SHOW ENGINE INNODB STATUS:

- · History list length low, below a few thousand.
- · Insert buffer merges close to rows inserted.
- Modified pages in buffer pool consistently well below innodb\_max\_dirty\_pages\_pct of the
  buffer pool. (Measure at a time when the server is not doing bulk inserts; it is normal during bulk
  inserts for the modified pages percentage to rise significantly.)
- Log sequence number Last checkpoint is at less than 7/8 or ideally less than 6/8 of the total size of the InnoDB log files.
- · Store system tablespace files on Fusion-io devices

You can take advantage of a doublewrite buffer-related I/O optimization by storing the files that contain the doublewrite storage area on Fusion-io devices that support atomic writes. (The doublewrite buffer storage area resides in doublewrite files. See Section 17.6.4, "Doublewrite Buffer".) When doublewrite storage area files are placed on Fusion-io devices that support atomic writes, the doublewrite buffer is automatically disabled and Fusion-io atomic writes are used for all data files. 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 <code>innodb\_flush\_method</code> setting of <code>O\_DIRECT</code> is recommended.



#### Note

Because the doublewrite buffer setting is global, the doublewrite buffer is also disabled for data files that do not reside on Fusion-io hardware.

· Disable logging of compressed pages

When using the InnoDB table compression feature, images of re-compressed pages are written to the redo log when changes are made to compressed data. This behavior is controlled by innodb\_log\_compressed\_pages, which is enabled by default to prevent corruption that can occur if a different version of the zlib compression algorithm is used during recovery. If you are certain that the zlib version is not subject to change, disable innodb\_log\_compressed\_pages to reduce redo log generation for workloads that modify compressed data.

# 10.5.9 Optimizing InnoDB Configuration Variables

Different settings work best for servers with light, predictable loads, versus servers that are running near full capacity all the time, or that experience spikes of high activity.

Because the InnoDB storage engine performs many of its optimizations automatically, many performance-tuning tasks involve monitoring to ensure that the database is performing well, and changing configuration options when performance drops. See Section 17.16, "InnoDB Integration with MySQL Performance Schema" for information about detailed InnoDB performance monitoring.

The main configuration steps you can perform include:

- Controlling the types of data change operations for which InnoDB buffers the changed data, to avoid
  frequent small disk writes. See Configuring Change Buffering. Enabling change buffering can give
  better performance on IO-bound workloads, but can cause issues during recovery, bulk load, or
  during buffer pool resizing. Having it disabled (default as of MySQL 8.4) helps ensure stability even if
  it may lower performance.
- Turning the adaptive hash indexing feature on and off using the innodb\_adaptive\_hash\_index option. See Section 17.5.3, "Adaptive Hash Index" for more information. You might change this setting during periods of unusual activity, then restore it to its original setting.
- Setting a limit on the number of concurrent threads that InnoDB processes, if context switching is a bottleneck. See Section 17.8.4, "Configuring Thread Concurrency for InnoDB".
- Controlling the amount of prefetching that InnoDB does with its read-ahead operations. When the system has unused I/O capacity, more read-ahead can improve the performance of queries. Too much read-ahead can cause periodic drops in performance on a heavily loaded system. See Section 17.8.3.4, "Configuring InnoDB Buffer Pool Prefetching (Read-Ahead)".
- Increasing the number of background threads for read or write operations, if you have a high-end I/O subsystem that is not fully utilized by the default values. See Section 17.8.5, "Configuring the Number of Background InnoDB I/O Threads".
- Controlling how much I/O InnoDB performs in the background. See Section 17.8.7, "Configuring InnoDB I/O Capacity". You might scale back this setting if you observe periodic drops in performance.
- Controlling the algorithm that determines when InnoDB performs certain types of background
  writes. See Section 17.8.3.5, "Configuring Buffer Pool Flushing". The algorithm works for some
  types of workloads but not others, so you might disable this feature if you observe periodic drops in
  performance.
- Taking advantage of multicore processors and their cache memory configuration, to minimize delays in context switching. See Section 17.8.8, "Configuring Spin Lock Polling".
- Preventing one-time operations such as table scans from interfering with the frequently accessed data stored in the InnoDB buffer cache. See Section 17.8.3.3, "Making the Buffer Pool Scan Resistant".
- Adjusting log files to a size that makes sense for reliability and crash recovery. InnoDB log files have often been kept small to avoid long startup times after a crash. Optimizations introduced in MySQL 5.5 speed up certain steps of the crash recovery process. In particular, scanning the redo log and applying the redo log are faster due to improved algorithms for memory management. If you have kept your log files artificially small to avoid long startup times, you can now consider increasing log file size to reduce the I/O that occurs due recycling of redo log records.
- Configuring the size and number of instances for the InnoDB buffer pool, especially important for systems with multi-gigabyte buffer pools. See Section 17.8.3.2, "Configuring Multiple Buffer Pool Instances".
- Increasing the maximum number of concurrent transactions, which dramatically improves scalability for the busiest databases. See Section 17.6.6, "Undo Logs".
- Moving purge operations (a type of garbage collection) into a background thread. See Section 17.8.9, "Purge Configuration". To effectively measure the results of this setting, tune the other I/O-related and thread-related configuration settings first.
- Reducing the amount of switching that InnoDB does between concurrent threads, so that SQL operations on a busy server do not queue up and form a "traffic jam". Set a value for the innodb\_thread\_concurrency option, up to approximately 32 for a high-powered modern system. Increase the value for the innodb\_concurrency\_tickets option, typically to 5000 or so. This combination of options sets a cap on the number of threads that InnoDB processes at any one time,

and allows each thread to do substantial work before being swapped out, so that the number of waiting threads stays low and operations can complete without excessive context switching.

## 10.5.10 Optimizing InnoDB for Systems with Many Tables

• If you have configured non-persistent optimizer statistics (a non-default configuration), InnoDB computes index cardinality values for a table the first time that table is accessed after startup, instead of storing such values in the table. This step can take significant time on systems that partition the data into many tables. Since this overhead only applies to the initial table open operation, to "warm up" a table for later use, access it immediately after startup by issuing a statement such as SELECT 1 FROM tbl\_name LIMIT 1.

Optimizer statistics are persisted to disk by default, enabled by the <code>innodb\_stats\_persistent</code> configuration option. For information about persistent optimizer statistics, see Section 17.8.10.1, "Configuring Persistent Optimizer Statistics Parameters".

# 10.6 Optimizing for MyISAM Tables

The MyISAM storage engine performs best with read-mostly data or with low-concurrency operations, because table locks limit the ability to perform simultaneous updates. In MySQL, InnoDB is the default storage engine rather than MyISAM.

# 10.6.1 Optimizing MyISAM Queries

Some general tips for speeding up queries on MyISAM tables:

- To help MySQL better optimize queries, use ANALYZE TABLE or run myisamchk —-analyze on a table after it has been loaded with data. This updates a value for each index part that indicates the average number of rows that have the same value. (For unique indexes, this is always 1.) MySQL uses this to decide which index to choose when you join two tables based on a nonconstant expression. You can check the result from the table analysis by using SHOW INDEX FROM tbl\_name and examining the Cardinality value. myisamchk —-description —-verbose shows index distribution information.
- To sort an index and data according to an index, use myisamchk --sort-index --sort-records=1 (assuming that you want to sort on index 1). This is a good way to make queries faster if you have a unique index from which you want to read all rows in order according to the index. The first time you sort a large table this way, it may take a long time.
- Try to avoid complex SELECT queries on MyISAM tables that are updated frequently, to avoid problems with table locking that occur due to contention between readers and writers.
- MyISAM supports concurrent inserts: If a table has no free blocks in the middle of the data file, you can INSERT new rows into it at the same time that other threads are reading from the table. If it is important to be able to do this, consider using the table in ways that avoid deleting rows. Another possibility is to run OPTIMIZE TABLE to defragment the table after you have deleted a lot of rows from it. This behavior is altered by setting the concurrent\_insert variable. You can force new rows to be appended (and therefore permit concurrent inserts), even in tables that have deleted rows. See Section 10.11.3, "Concurrent Inserts".
- For MyISAM tables that change frequently, try to avoid all variable-length columns (VARCHAR, BLOB, and TEXT). The table uses dynamic row format if it includes even a single variable-length column. See Chapter 18, *Alternative Storage Engines*.
- It is normally not useful to split a table into different tables just because the rows become large. In accessing a row, the biggest performance hit is the disk seek needed to find the first byte of the row. After finding the data, most modern disks can read the entire row fast enough for most applications. The only cases where splitting up a table makes an appreciable difference is if it is a MyISAM table using dynamic row format that you can change to a fixed row size, or if you very often need to scan the table but do not need most of the columns. See Chapter 18, Alternative Storage Engines.

- Use ALTER TABLE ... ORDER BY expr1, expr2, ... if you usually retrieve rows in expr1, expr2, ... order. By using this option after extensive changes to the table, you may be able to get higher performance.
- If you often need to calculate results such as counts based on information from a lot of rows, it may
  be preferable to introduce a new table and update the counter in real time. An update of the following
  form is very fast:

```
UPDATE tbl_name SET count_col=count_col+1 WHERE key_col=constant;
```

This is very important when you use MySQL storage engines such as MyISAM that has only table-level locking (multiple readers with single writers). This also gives better performance with most database systems, because the row locking manager in this case has less to do.

- Use OPTIMIZE TABLE periodically to avoid fragmentation with dynamic-format MyISAM tables. See Section 18.2.3, "MyISAM Table Storage Formats".
- Declaring a MyISAM table with the DELAY\_KEY\_WRITE=1 table option makes index updates faster because they are not flushed to disk until the table is closed. The downside is that if something kills the server while such a table is open, you must ensure that the table is okay by running the server with the myisam\_recover\_options system variable set, or by running myisamchk before restarting the server. (However, even in this case, you should not lose anything by using DELAY\_KEY\_WRITE, because the key information can always be generated from the data rows.)
- Strings are automatically prefix- and end-space compressed in MyISAM indexes. See Section 15.1.15, "CREATE INDEX Statement".
- You can increase performance by caching queries or answers in your application and then executing
  many inserts or updates together. Locking the table during this operation ensures that the index
  cache is only flushed once after all updates.

# 10.6.2 Bulk Data Loading for MyISAM Tables

These performance tips supplement the general guidelines for fast inserts in Section 10.2.5.1, "Optimizing INSERT Statements".

- For a MyISAM table, you can use concurrent inserts to add rows at the same time that SELECT statements are running, if there are no deleted rows in middle of the data file. See Section 10.11.3, "Concurrent Inserts".
- With some extra work, it is possible to make LOAD DATA run even faster for a MyISAM table when the table has many indexes. Use the following procedure:
  - 1. Execute a FLUSH TABLES statement or a mysqladmin flush-tables command.
  - 2. Use myisamchk --keys-used=0 -rq /path/to/db/tbl\_name to remove all use of indexes for the table.
  - 3. Insert data into the table with LOAD DATA. This does not update any indexes and therefore is very fast.
  - 4. If you intend only to read from the table in the future, use myisampack to compress it. See Section 18.2.3.3, "Compressed Table Characteristics".
  - 5. Re-create the indexes with myisamchk -rq /path/to/db/tbl\_name. This creates the index tree in memory before writing it to disk, which is much faster than updating the index during LOAD DATA because it avoids lots of disk seeks. The resulting index tree is also perfectly balanced.
  - 6. Execute a FLUSH TABLES statement or a mysgladmin flush-tables command.

LOAD DATA performs the preceding optimization automatically if the MyISAM table into which you insert data is empty. The main difference between automatic optimization and using the procedure

explicitly is that you can let myisamchk allocate much more temporary memory for the index creation than you might want the server to allocate for index re-creation when it executes the LOAD DATA statement.

You can also disable or enable the nonunique indexes for a MyISAM table by using the following statements rather than myisamchk. If you use these statements, you can skip the FLUSH TABLES operations:

```
ALTER TABLE tbl_name DISABLE KEYS;
ALTER TABLE tbl_name ENABLE KEYS;
```

• To speed up INSERT operations that are performed with multiple statements for nontransactional tables, lock your tables:

```
LOCK TABLES a WRITE;
INSERT INTO a VALUES (1,23),(2,34),(4,33);
INSERT INTO a VALUES (8,26),(6,29);
...
UNLOCK TABLES;
```

This benefits performance because the index buffer is flushed to disk only once, after all INSERT statements have completed. Normally, there would be as many index buffer flushes as there are INSERT statements. Explicit locking statements are not needed if you can insert all rows with a single INSERT.

Locking also lowers the total time for multiple-connection tests, although the maximum wait time for individual connections might go up because they wait for locks. Suppose that five clients attempt to perform inserts simultaneously as follows:

- Connection 1 does 1000 inserts
- · Connections 2, 3, and 4 do 1 insert
- Connection 5 does 1000 inserts

If you do not use locking, connections 2, 3, and 4 finish before 1 and 5. If you use locking, connections 2, 3, and 4 probably do not finish before 1 or 5, but the total time should be about 40% faster.

INSERT, UPDATE, and DELETE operations are very fast in MySQL, but you can obtain better overall performance by adding locks around everything that does more than about five successive inserts or updates. If you do very many successive inserts, you could do a LOCK TABLES followed by an UNLOCK TABLES once in a while (each 1,000 rows or so) to permit other threads to access table. This would still result in a nice performance gain.

INSERT is still much slower for loading data than LOAD DATA, even when using the strategies just outlined.

To increase performance for MyISAM tables, for both LOAD DATA and INSERT, enlarge the key
cache by increasing the key\_buffer\_size system variable. See Section 7.1.1, "Configuring the
Server".

# 10.6.3 Optimizing REPAIR TABLE Statements

REPAIR TABLE for MyISAM tables is similar to using myisamchk for repair operations, and some of the same performance optimizations apply:

 myisamchk has variables that control memory allocation. You may be able to its improve performance by setting these variables, as described in Section 6.6.4.6, "myisamchk Memory Usage". • For REPAIR TABLE, the same principle applies, but because the repair is done by the server, you set server system variables instead of myisamchk variables. Also, in addition to setting memoryallocation variables, increasing the myisam\_max\_sort\_file\_size system variable increases the likelihood that the repair uses the faster filesort method and avoids the slower repair by key cache method. Set the variable to the maximum file size for your system, after checking to be sure that there is enough free space to hold a copy of the table files. The free space must be available in the file system containing the original table files.

Suppose that a myisamchk table-repair operation is done using the following options to set its memory-allocation variables:

```
--key_buffer_size=128M --myisam_sort_buffer_size=256M --read_buffer_size=64M --write_buffer_size=64M
```

Some of those myisamchk variables correspond to server system variables:

myisamchk Variable	System Variable
key_buffer_size	key_buffer_size
myisam_sort_buffer_size	myisam_sort_buffer_size
read_buffer_size	read_buffer_size
write_buffer_size	none

Each of the server system variables can be set at runtime, and some of them (myisam\_sort\_buffer\_size, read\_buffer\_size) have a session value in addition to a global value. Setting a session value limits the effect of the change to your current session and does not affect other users. Changing a global-only variable (key\_buffer\_size, myisam\_max\_sort\_file\_size) affects other users as well. For key\_buffer\_size, you must take into account that the buffer is shared with those users. For example, if you set the myisamchk key\_buffer\_size variable to 128MB, you could set the corresponding key\_buffer\_size system variable larger than that (if it is not already set larger), to permit key buffer use by activity in other sessions. However, changing the global key buffer size invalidates the buffer, causing increased disk I/O and slowdown for other sessions. An alternative that avoids this problem is to use a separate key cache, assign to it the indexes from the table to be repaired, and deallocate it when the repair is complete. See Section 10.10.2.2, "Multiple Key Caches".

Based on the preceding remarks, a REPAIR TABLE operation can be done as follows to use settings similar to the myisamchk command. Here a separate 128MB key buffer is allocated and the file system is assumed to permit a file size of at least 100GB.

```
SET SESSION myisam_sort_buffer_size = 256*1024*1024;

SET SESSION read_buffer_size = 64*1024*1024;

SET GLOBAL myisam_max_sort_file_size = 100*1024*1024*1024;

SET GLOBAL repair_cache.key_buffer_size = 128*1024*1024;

CACHE INDEX tbl_name IN repair_cache;

LOAD INDEX INTO CACHE tbl_name;

REPAIR TABLE tbl_name;

SET GLOBAL repair_cache.key_buffer_size = 0;
```

If you intend to change a global variable but want to do so only for the duration of a REPAIR TABLE operation to minimally affect other users, save its value in a user variable and restore it afterward. For example:

```
SET @old_myisam_sort_buffer_size = @@GLOBAL.myisam_max_sort_file_size;
SET GLOBAL myisam_max_sort_file_size = 100*1024*1024*1024;
REPAIR TABLE tbl_name;
SET GLOBAL myisam_max_sort_file_size = @old_myisam_max_sort_file_size;
```

The system variables that affect REPAIR TABLE can be set globally at server startup if you want the values to be in effect by default. For example, add these lines to the server my.cnf file:

```
[mysqld]
myisam_sort_buffer_size=256M
key_buffer_size=1G
```

myisam\_max\_sort\_file\_size=100G

These settings do not include read\_buffer\_size. Setting read\_buffer\_size globally to a large value does so for all sessions and can cause performance to suffer due to excessive memory allocation for a server with many simultaneous sessions.

# 10.7 Optimizing for MEMORY Tables

Consider using MEMORY tables for noncritical data that is accessed often, and is read-only or rarely updated. Benchmark your application against equivalent InnoDB or MyISAM tables under a realistic workload, to confirm that any additional performance is worth the risk of losing data, or the overhead of copying data from a disk-based table at application start.

For best performance with MEMORY tables, examine the kinds of queries against each table, and specify the type to use for each associated index, either a B-tree index or a hash index. On the CREATE INDEX statement, use the clause USING BTREE or USING HASH. B-tree indexes are fast for queries that do greater-than or less-than comparisons through operators such as > or BETWEEN. Hash indexes are only fast for queries that look up single values through the = operator, or a restricted set of values through the IN operator. For why USING BTREE is often a better choice than the default USING HASH, see Section 10.2.1.23, "Avoiding Full Table Scans". For implementation details of the different types of MEMORY indexes, see Section 10.3.9, "Comparison of B-Tree and Hash Indexes".

# 10.8 Understanding the Query Execution Plan

Depending on the details of your tables, columns, indexes, and the conditions in your WHERE clause, the MySQL optimizer considers many techniques to efficiently perform the lookups involved in an SQL query. A query on a huge table can be performed without reading all the rows; a join involving several tables can be performed without comparing every combination of rows. The set of operations that the optimizer chooses to perform the most efficient query is called the "query execution plan", also known as the EXPLAIN plan. Your goals are to recognize the aspects of the EXPLAIN plan that indicate a query is optimized well, and to learn the SQL syntax and indexing techniques to improve the plan if you see some inefficient operations.

# 10.8.1 Optimizing Queries with EXPLAIN

The EXPLAIN statement provides information about how MySQL executes statements:

- EXPLAIN works with SELECT, DELETE, INSERT, REPLACE, and UPDATE 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 10.8.2, "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 10.8.4, "Obtaining Execution Plan Information for a Named Connection".
- For SELECT statements, EXPLAIN produces additional execution plan information that can be displayed using SHOW WARNINGS. See Section 10.8.3, "Extended EXPLAIN Output Format".
- EXPLAIN is useful for examining queries involving partitioned tables. See Section 26.3.5, "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.

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 15.2.13, "SELECT Statement".) However, STRAIGHT\_JOIN may prevent indexes from being used because it disables semijoin transformations. See Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations.

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 Section 10.15, "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 15.7.3.1, "ANALYZE TABLE Statement".



#### Note

EXPLAIN can also be used to obtain information about the columns in a table. EXPLAIN <code>tbl\_name</code> is synonymous with <code>DESCRIBE tbl\_name</code> and <code>SHOW COLUMNS FROM tbl\_name</code>. For more information, see Section 15.8.1, "DESCRIBE Statement", and Section 15.7.7.6, "SHOW COLUMNS Statement".

# 10.8.2 EXPLAIN Output Format

The EXPLAIN statement provides information about how MySQL executes statements. EXPLAIN works with SELECT, DELETE, INSERT, REPLACE, and UPDATE statements.

EXPLAIN returns a row of information for each table used in the SELECT statement. It lists the tables in the output in the order that MySQL would read them while processing the statement. This means that MySQL reads a row from the first table, then finds a matching row in the second table, and then in the third table, and so on. When all tables are processed, MySQL outputs the selected columns and backtracks through the table list until a table is found for which there are more matching rows. The next row is read from this table and the process continues with the next table.



### Note

MySQL Workbench has a Visual Explain capability that provides a visual representation of  ${\tt EXPLAIN}$  output. See Tutorial: Using Explain to Improve Query Performance.

- EXPLAIN Output Columns
- EXPLAIN Join Types
- EXPLAIN Extra Information
- EXPLAIN Output Interpretation

### **EXPLAIN Output Columns**

This section describes the output columns produced by EXPLAIN. Later sections provide additional information about the type and Extra columns.

Each output row from EXPLAIN provides information about one table. Each row contains the values summarized in Table 10.1, "EXPLAIN Output Columns", and described in more detail following the table. Column names are shown in the table's first column; the second column provides the equivalent property name shown in the output when FORMAT=JSON is used.

**Table 10.1 EXPLAIN Output Columns** 

Column	JSON Name	Meaning
id	select_id	The SELECT identifier
select_type	None	The SELECT type

Column	JSON Name	Meaning
table	table_name	The table for the output row
partitions	partitions	The matching partitions
type	access_type	The join type
possible_keys	possible_keys	The possible indexes to choose
key	key	The index actually chosen
key_len	key_length	The length of the chosen key
ref	ref	The columns compared to the index
rows	rows	Estimate of rows to be examined
filtered	filtered	Percentage of rows filtered by table condition
Extra	None	Additional information



### Note

JSON properties which are  ${\tt NULL}$  are not displayed in JSON-formatted  ${\tt EXPLAIN}$  output.

• id (JSON name: select\_id)

The SELECT identifier. This is the sequential number of the SELECT within the query. The value can be NULL if the row refers to the union result of other rows. In this case, the table column shows a value like < unionM, N> to indicate that the row refers to the union of the rows with id values of M and N.

• select\_type (JSON name: none)

The type of SELECT, which can be any of those shown in the following table. A JSON-formatted EXPLAIN exposes the SELECT type as a property of a query\_block, unless it is SIMPLE or PRIMARY. The JSON names (where applicable) are also shown in the table.

select_type Value	JSON Name	Meaning
SIMPLE	None	Simple SELECT (not using UNION or subqueries)
PRIMARY	None	Outermost SELECT
UNION	None	Second or later SELECT statement in a UNION
DEPENDENT UNION	dependent (true)	Second or later SELECT statement in a UNION, dependent on outer query
UNION RESULT	union_result	Result of a UNION.
SUBQUERY	None	First SELECT in subquery
DEPENDENT SUBQUERY	dependent (true)	First SELECT in subquery, dependent on outer query
DERIVED	None	Derived table
DEPENDENT DERIVED	dependent (true)	Derived table dependent on another table
MATERIALIZED	materialized_from_subque	rMaterialized subquery
UNCACHEABLE SUBQUERY	cacheable (false)	A subquery for which the result cannot be cached and must be

select_type Value	JSON Name	Meaning
		re-evaluated for each row of the outer query
UNCACHEABLE UNION	cacheable (false)	The second or later select in a UNION that belongs to an uncacheable subquery (see UNCACHEABLE SUBQUERY)

DEPENDENT typically signifies the use of a correlated subquery. See Section 15.2.15.7, "Correlated Subqueries".

DEPENDENT SUBQUERY evaluation differs from UNCACHEABLE SUBQUERY evaluation. For DEPENDENT SUBQUERY, the subquery is re-evaluated only once for each set of different values of the variables from its outer context. For UNCACHEABLE SUBQUERY, the subquery is re-evaluated for each row of the outer context.

When you specify FORMAT=JSON with EXPLAIN, the output has no single property directly equivalent to select\_type; the query\_block property corresponds to a given SELECT. Properties equivalent to most of the SELECT subquery types just shown are available (an example being materialized\_from\_subquery for MATERIALIZED), and are displayed when appropriate. There are no JSON equivalents for SIMPLE or PRIMARY.

The select\_type value for non-SELECT statements displays the statement type for affected tables. For example, select\_type is DELETE for DELETE statements.

• table (JSON name: table\_name)

The name of the table to which the row of output refers. This can also be one of the following values:

- <unionM, N>: The row refers to the union of the rows with id values of M and N.
- <derivedN>: The row refers to the derived table result for the row with an id value of N. A derived table may result, for example, from a subquery in the FROM clause.
- <subqueryN>: The row refers to the result of a materialized subquery for the row with an id value of N. See Section 10.2.2.2, "Optimizing Subqueries with Materialization".
- partitions (JSON name: partitions)

The partitions from which records would be matched by the query. The value is NULL for nonpartitioned tables. See Section 26.3.5, "Obtaining Information About Partitions".

• type (JSON name: access type)

The join type. For descriptions of the different types, see EXPLAIN Join Types.

• possible\_keys (JSON name: possible\_keys)

The possible\_keys column indicates the indexes from which MySQL can choose to find the rows in this table. Note that this column is totally independent of the order of the tables as displayed in the output from EXPLAIN. That means that some of the keys in possible\_keys might not be usable in practice with the generated table order.

If this column is NULL (or undefined in JSON-formatted output), there are no relevant indexes. In this case, you may be able to improve the performance of your query by examining the WHERE clause to check whether it refers to some column or columns that would be suitable for indexing. If so, create an appropriate index and check the query with EXPLAIN again. See Section 15.1.9, "ALTER TABLE Statement".

To see what indexes a table has, use SHOW INDEX FROM tbl\_name.

### • key (JSON name: key)

The key column indicates the key (index) that MySQL actually decided to use. If MySQL decides to use one of the possible\_keys indexes to look up rows, that index is listed as the key value.

It is possible that key may name an index that is not present in the possible\_keys value. This can happen if none of the possible\_keys indexes are suitable for looking up rows, but all the columns selected by the query are columns of some other index. That is, the named index covers the selected columns, so although it is not used to determine which rows to retrieve, an index scan is more efficient than a data row scan.

For InnoDB, a secondary index might cover the selected columns even if the query also selects the primary key because InnoDB stores the primary key value with each secondary index. If key is NULL, MySQL found no index to use for executing the query more efficiently.

To force MySQL to use or ignore an index listed in the possible\_keys column, use FORCE INDEX, USE INDEX, or IGNORE INDEX in your query. See Section 10.9.4, "Index Hints".

For MyISAM tables, running ANALYZE TABLE helps the optimizer choose better indexes. For MyISAM tables, myisamchk —analyze does the same. See Section 15.7.3.1, "ANALYZE TABLE Statement", and Section 9.6, "MyISAM Table Maintenance and Crash Recovery".

• key len (JSON name: key length)

The key\_len column indicates the length of the key that MySQL decided to use. The value of key\_len enables you to determine how many parts of a multiple-part key MySQL actually uses. If the key column says NULL, the key\_len column also says NULL.

Due to the key storage format, the key length is one greater for a column that can be NULL than for a NOT NULL column.

• ref (JSON name: ref)

The ref column shows which columns or constants are compared to the index named in the key column to select rows from the table.

If the value is func, the value used is the result of some function. To see which function, use SHOW WARNINGS following EXPLAIN to see the extended EXPLAIN output. The function might actually be an operator such as an arithmetic operator.

• rows (JSON name: rows)

The rows column indicates the number of rows MySQL believes it must examine to execute the query.

For InnoDB tables, this number is an estimate, and may not always be exact.

• filtered (JSON name: filtered)

The filtered column indicates an estimated percentage of table rows that are filtered by the table condition. The maximum value is 100, which means no filtering of rows occurred. Values decreasing from 100 indicate increasing amounts of filtering. rows shows the estimated number of rows examined and rows x filtered shows the number of rows that are joined with the following table. For example, if rows is 1000 and filtered is 50.00 (50%), the number of rows to be joined with the following table is  $1000 \times 50\% = 500$ .

• Extra (JSON name: none)

This column contains additional information about how MySQL resolves the query. For descriptions of the different values, see EXPLAIN Extra Information.

There is no single JSON property corresponding to the Extra column; however, values that can occur in this column are exposed as JSON properties, or as the text of the message property.

### **EXPLAIN Join Types**

The type column of EXPLAIN output describes how tables are joined. In JSON-formatted output, these are found as values of the access\_type property. The following list describes the join types, ordered from the best type to the worst:

• system

The table has only one row (= system table). This is a special case of the const join type.

• const

The table has at most one matching row, which is read at the start of the query. Because there is only one row, values from the column in this row can be regarded as constants by the rest of the optimizer. const tables are very fast because they are read only once.

const is used when you compare all parts of a PRIMARY KEY or UNIQUE index to constant values. In the following queries, tbl\_name can be used as a const table:

```
SELECT * FROM tbl_name WHERE primary_key=1;

SELECT * FROM tbl_name

WHERE primary_key_part1=1 AND primary_key_part2=2;
```

• eq\_ref

One row is read from this table for each combination of rows from the previous tables. Other than the system and const types, this is the best possible join type. It is used when all parts of an index are used by the join and the index is a PRIMARY KEY OF UNIQUE NOT NULL index.

eq\_ref can be used for indexed columns that are compared using the = operator. The comparison value can be a constant or an expression that uses columns from tables that are read before this table. In the following examples, MySQL can use an eq\_ref join to process ref\_table:

```
SELECT * FROM ref_table,other_table
WHERE ref_table.key_column=other_table.column;

SELECT * FROM ref_table,other_table
WHERE ref_table.key_column_part1=other_table.column
AND ref_table.key_column_part2=1;
```

• ref

All rows with matching index values are read from this table for each combination of rows from the previous tables. ref is used if the join uses only a leftmost prefix of the key or if the key is not a PRIMARY KEY OR UNIQUE index (in other words, if the join cannot select a single row based on the key value). If the key that is used matches only a few rows, this is a good join type.

ref can be used for indexed columns that are compared using the = or <=> operator. In the following examples, MySQL can use a ref join to process  $ref\_table$ :

```
SELECT * FROM ref_table WHERE key_column=expr;

SELECT * FROM ref_table,other_table
WHERE ref_table.key_column=other_table.column;
```

```
SELECT * FROM ref_table,other_table
WHERE ref_table.key_column_part1=other_table.column
AND ref_table.key_column_part2=1;
```

• fulltext

The join is performed using a FULLTEXT index.

• ref\_or\_null

This join type is like ref, but with the addition that MySQL does an extra search for rows that contain NULL values. This join type optimization is used most often in resolving subqueries. In the following examples, MySQL can use a ref\_or\_null join to process ref\_table:

```
SELECT * FROM ref_table
WHERE key_column=expr OR key_column IS NULL;
```

See Section 10.2.1.15, "IS NULL Optimization".

• index\_merge

This join type indicates that the Index Merge optimization is used. In this case, the key column in the output row contains a list of indexes used, and key\_len contains a list of the longest key parts for the indexes used. For more information, see Section 10.2.1.3, "Index Merge Optimization".

• unique\_subquery

This type replaces eq\_ref for some IN subqueries of the following form:

```
value IN (SELECT primary_key FROM single_table WHERE some_expr)
```

unique\_subquery is just an index lookup function that replaces the subquery completely for better efficiency.

• index\_subquery

This join type is similar to unique\_subquery. It replaces IN subqueries, but it works for nonunique indexes in subqueries of the following form:

```
value IN (SELECT key_column FROM single_table WHERE some_expr)
```

• range

Only rows that are in a given range are retrieved, using an index to select the rows. The key column in the output row indicates which index is used. The  $key\_len$  contains the longest key part that was used. The ref column is NULL for this type.

range can be used when a key column is compared to a constant using any of the =, <>, >, >=, <, <=, IS NULL, <=>, BETWEEN, LIKE, or IN() operators:

```
SELECT * FROM tbl_name
WHERE key_column = 10;

SELECT * FROM tbl_name
WHERE key_column BETWEEN 10 and 20;

SELECT * FROM tbl_name
WHERE key_column IN (10,20,30);

SELECT * FROM tbl_name
WHERE key_column IN (10,20,30);
```

• index

The index join type is the same as ALL, except that the index tree is scanned. This occurs two ways:

- If the index is a covering index for the queries and can be used to satisfy all data required from the table, only the index tree is scanned. In this case, the Extra column says Using index. An index-only scan usually is faster than ALL because the size of the index usually is smaller than the table data.
- A full table scan is performed using reads from the index to look up data rows in index order. Uses index does not appear in the Extra column.

MySQL can use this join type when the query uses only columns that are part of a single index.

• ALL

A full table scan is done for each combination of rows from the previous tables. This is normally not good if the table is the first table not marked const, and usually *very* bad in all other cases. Normally, you can avoid ALL by adding indexes that enable row retrieval from the table based on constant values or column values from earlier tables.

### **EXPLAIN Extra Information**

The Extra column of EXPLAIN output contains additional information about how MySQL resolves the query. The following list explains the values that can appear in this column. Each item also indicates for JSON-formatted output which property displays the Extra value. For some of these, there is a specific property. The others display as the text of the message property.

If you want to make your queries as fast as possible, look out for Extra column values of Using filesort and Using temporary, or, in JSON-formatted EXPLAIN output, for using\_filesort and using\_temporary\_table properties equal to true.

• Backward index scan (JSON: backward\_index\_scan)

The optimizer is able to use a descending index on an InnoDB table. Shown together with Using index. For more information, see Section 10.3.13, "Descending Indexes".

• Child of 'table' pushed join@1 (JSON: message text)

This table is referenced as the child of *table* in a join that can be pushed down to the NDB kernel. Applies only in NDB Cluster, when pushed-down joins are enabled. See the description of the <a href="mailto:ndb\_join\_pushdown">ndb\_join\_pushdown</a> server system variable for more information and examples.

• const row not found (JSON property: const\_row\_not\_found)

For a query such as SELECT ... FROM tbl\_name, the table was empty.

Deleting all rows (JSON property: message)

For DELETE, some storage engines (such as MyISAM) support a handler method that removes all table rows in a simple and fast way. This Extra value is displayed if the engine uses this optimization.

• Distinct (JSON property: distinct)

MySQL is looking for distinct values, so it stops searching for more rows for the current row combination after it has found the first matching row.

• FirstMatch(tbl\_name) (JSON property: first\_match)

The semijoin FirstMatch join shortcutting strategy is used for tbl name.

• Full scan on NULL key (JSON property: message)

This occurs for subquery optimization as a fallback strategy when the optimizer cannot use an index-lookup access method.

Impossible HAVING (JSON property: message)

The HAVING clause is always false and cannot select any rows.

• Impossible WHERE (JSON property: message)

The WHERE clause is always false and cannot select any rows.

• Impossible WHERE noticed after reading const tables (JSON property: message)

MySQL has read all const (and system) tables and notice that the WHERE clause is always false.

• LooseScan(m..n) (JSON property: message)

The semijoin LooseScan strategy is used. m and n are key part numbers.

• No matching min/max row (JSON property: message)

No row satisfies the condition for a query such as  $\mathtt{SELECT\ MIN}(\ldots)$  FROM  $\ldots$  WHERE condition.

• no matching row in const table (JSON property: message)

For a query with a join, there was an empty table or a table with no rows satisfying a unique index condition.

• No matching rows after partition pruning (JSON property: message)

For DELETE or UPDATE, the optimizer found nothing to delete or update after partition pruning. It is similar in meaning to Impossible WHERE for SELECT statements.

• No tables used (JSON property: message)

The query has no FROM clause, or has a FROM DUAL clause.

For INSERT or REPLACE statements, EXPLAIN displays this value when there is no SELECT part. For example, it appears for EXPLAIN INSERT INTO t VALUES(10) because that is equivalent to EXPLAIN INSERT INTO t SELECT 10 FROM DUAL.

• Not exists (JSON property: message)

MySQL was able to do a LEFT JOIN optimization on the query and does not examine more rows in this table for the previous row combination after it finds one row that matches the LEFT JOIN criteria. Here is an example of the type of query that can be optimized this way:

```
SELECT * FROM t1 LEFT JOIN t2 ON t1.id=t2.id
WHERE t2.id IS NULL;
```

Assume that t2.id is defined as NOT NULL. In this case, MySQL scans t1 and looks up the rows in t2 using the values of t1.id. If MySQL finds a matching row in t2, it knows that t2.id can never be NULL, and does not scan through the rest of the rows in t2 that have the same id value. In other words, for each row in t1, MySQL needs to do only a single lookup in t2, regardless of how many rows actually match in t2.

This can also indicate that a <code>WHERE</code> condition of the form <code>NOT IN (subquery)</code> or <code>NOT EXISTS (subquery)</code> has been transformed internally into an antijoin. This removes the subquery and brings its tables into the plan for the topmost query, providing improved cost planning. By merging semijoins and antijoins, the optimizer can reorder tables in the execution plan more freely, in some cases resulting in a faster plan.

You can see when an antijoin transformation is performed for a given query by checking the Message column from SHOW WARNINGS following execution of EXPLAIN, or in the output of EXPLAIN FORMAT=TREE.



#### Note

An antijoin is the complement of a semijoin table\_a JOIN table\_b ON condition. The antijoin returns all rows from table\_a for which there is no row in table\_b which matches condition.

• Plan isn't ready yet (JSON property: none)

This value occurs with EXPLAIN FOR CONNECTION when the optimizer has not finished creating the execution plan for the statement executing in the named connection. If execution plan output comprises multiple lines, any or all of them could have this Extra value, depending on the progress of the optimizer in determining the full execution plan.

• Range checked for each record (index map: N) (JSON property: message)

MySQL found no good index to use, but found that some of indexes might be used after column values from preceding tables are known. For each row combination in the preceding tables, MySQL checks whether it is possible to use a range or index\_merge access method to retrieve rows. This is not very fast, but is faster than performing a join with no index at all. The applicability criteria are as described in Section 10.2.1.2, "Range Optimization", and Section 10.2.1.3, "Index Merge Optimization", with the exception that all column values for the preceding table are known and considered to be constants.

Indexes are numbered beginning with 1, in the same order as shown by SHOW INDEX for the table. The index map value N is a bitmask value that indicates which indexes are candidates. For example, a value of  $0 \times 19$  (binary 11001) means that indexes 1, 4, and 5 are considered.

• Recursive (JSON property: recursive)

This indicates that the row applies to the recursive SELECT part of a recursive common table expression. See Section 15.2.20, "WITH (Common Table Expressions)".

• Rematerialize (JSON property: rematerialize)

Rematerialize (X, ...) is displayed in the EXPLAIN row for table T, where X is any lateral derived table whose rematerialization is triggered when a new row of T is read. For example:

```
SELECT
...

FROM
t,
LATERAL (derived table that refers to t) AS dt
...
```

The content of the derived table is rematerialized to bring it up to date each time a new row of t is processed by the top query.

• Scanned N databases (JSON property: message)

This indicates how many directory scans the server performs when processing a query for INFORMATION\_SCHEMA tables, as described in Section 10.2.3, "Optimizing INFORMATION\_SCHEMA Queries". The value of N can be 0, 1, or all.

• Select tables optimized away (JSON property: message)

The optimizer determined 1) that at most one row should be returned, and 2) that to produce this row, a deterministic set of rows must be read. When the rows to be read can be read during the

optimization phase (for example, by reading index rows), there is no need to read any tables during query execution.

The first condition is fulfilled when the query is implicitly grouped (contains an aggregate function but no GROUP BY clause). The second condition is fulfilled when one row lookup is performed per index used. The number of indexes read determines the number of rows to read.

Consider the following implicitly grouped query:

```
SELECT MIN(c1), MIN(c2) FROM t1;
```

Suppose that MIN(cl) can be retrieved by reading one index row and MIN(cl) can be retrieved by reading one row from a different index. That is, for each column cl and c2, there exists an index where the column is the first column of the index. In this case, one row is returned, produced by reading two deterministic rows.

This Extra value does not occur if the rows to read are not deterministic. Consider this query:

```
SELECT MIN(c2) FROM t1 WHERE c1 <= 10;
```

Suppose that (c1, c2) is a covering index. Using this index, all rows with  $c1 \le 10$  must be scanned to find the minimum c2 value. By contrast, consider this query:

```
SELECT MIN(c2) FROM t1 WHERE c1 = 10;
```

In this case, the first index row with c1 = 10 contains the minimum c2 value. Only one row must be read to produce the returned row.

For storage engines that maintain an exact row count per table (such as MyISAM, but not InnoDB), this Extra value can occur for COUNT(\*) queries for which the WHERE clause is missing or always true and there is no GROUP BY clause. (This is an instance of an implicitly grouped query where the storage engine influences whether a deterministic number of rows can be read.)

Skip\_open\_table, Open\_frm\_only, Open\_full\_table (JSON property: message)

These values indicate file-opening optimizations that apply to queries for INFORMATION\_SCHEMA tables.

- Skip\_open\_table: Table files do not need to be opened. The information is already available from the data dictionary.
- Open\_frm\_only: Only the data dictionary need be read for table information.
- Open\_full\_table: Unoptimized information lookup. Table information must be read from the data dictionary and by reading table files.
- Start temporary, End temporary (JSON property: message)

This indicates temporary table use for the semijoin Duplicate Weedout strategy.

• unique row not found (JSON property: message)

For a query such as SELECT ... FROM *tbl\_name*, no rows satisfy the condition for a UNIQUE index or PRIMARY KEY on the table.

• Using filesort (JSON property: using\_filesort)

MySQL must do an extra pass to find out how to retrieve the rows in sorted order. The sort is done by going through all rows according to the join type and storing the sort key and pointer to the row for all rows that match the WHERE clause. The keys then are sorted and the rows are retrieved in sorted order. See Section 10.2.1.16, "ORDER BY Optimization".

• Using index (JSON property: using index)

The column information is retrieved from the table using only information in the index tree without having to do an additional seek to read the actual row. This strategy can be used when the query uses only columns that are part of a single index.

For InnoDB tables that have a user-defined clustered index, that index can be used even when Using index is absent from the Extra column. This is the case if type is index and key is PRIMARY.

Information about any covering indexes used is shown for EXPLAIN FORMAT=TRADITIONAL and EXPLAIN FORMAT=JSON. It is also shown for EXPLAIN FORMAT=TREE.

• Using index condition (JSON property: using\_index\_condition)

Tables are read by accessing index tuples and testing them first to determine whether to read full table rows. In this way, index information is used to defer ("push down") reading full table rows unless it is necessary. See Section 10.2.1.6, "Index Condition Pushdown Optimization".

• Using index for group-by (JSON property: using\_index\_for\_group\_by)

Similar to the Using index table access method, Using index for group-by indicates that MySQL found an index that can be used to retrieve all columns of a GROUP BY Or DISTINCT query without any extra disk access to the actual table. Additionally, the index is used in the most efficient way so that for each group, only a few index entries are read. For details, see Section 10.2.1.17, "GROUP BY Optimization".

• Using index for skip scan (JSON property: using\_index\_for\_skip\_scan)

Indicates that the Skip Scan access method is used. See Skip Scan Range Access Method.

• Using join buffer (Block Nested Loop), Using join buffer (Batched Key Access), Using join buffer (hash join) (JSON property: using\_join\_buffer)

Tables from earlier joins are read in portions into the join buffer, and then their rows are used from the buffer to perform the join with the current table. (Block Nested Loop) indicates use of the Block Nested-Loop algorithm, (Batched Key Access) indicates use of the Batched Key Access algorithm, and (hash join) indicates use of a hash join. That is, the keys from the table on the preceding line of the EXPLAIN output are buffered, and the matching rows are fetched in batches from the table represented by the line in which Using join buffer appears.

In JSON-formatted output, the value of using\_join\_buffer is always one of Block Nested Loop, Batched Key Access, or hash join.

For more information about hash joins, see Section 10.2.1.4, "Hash Join Optimization".

See Batched Key Access Joins, for information about the Batched Key Access algorithm.

• Using MRR (JSON property: message)

Tables are read using the Multi-Range Read optimization strategy. See Section 10.2.1.11, "Multi-Range Read Optimization".

Using sort\_union(...), Using union(...), Using intersect(...) (JSON property: message)

These indicate the particular algorithm showing how index scans are merged for the index\_merge join type. See Section 10.2.1.3, "Index Merge Optimization".

• Using temporary (JSON property: using\_temporary\_table)

To resolve the query, MySQL needs to create a temporary table to hold the result. This typically happens if the query contains GROUP BY and ORDER BY clauses that list columns differently.

• Using where (JSON property: attached\_condition)

A WHERE clause is used to restrict which rows to match against the next table or send to the client. Unless you specifically intend to fetch or examine all rows from the table, you may have something wrong in your query if the Extra value is not Using where and the table join type is ALL or index.

Using where has no direct counterpart in JSON-formatted output; the attached\_condition property contains any WHERE condition used.

• Using where with pushed condition (JSON property: message)

This item applies to NDB tables *only*. It means that NDB Cluster is using the Condition Pushdown optimization to improve the efficiency of a direct comparison between a nonindexed column and a constant. In such cases, the condition is "pushed down" to the cluster's data nodes and is evaluated on all data nodes simultaneously. This eliminates the need to send nonmatching rows over the network, and can speed up such queries by a factor of 5 to 10 times over cases where Condition Pushdown could be but is not used. For more information, see Section 10.2.1.5, "Engine Condition Pushdown Optimization".

• Zero limit (JSON property: message)

The query had a LIMIT 0 clause and cannot select any rows.

### **EXPLAIN Output Interpretation**

You can get a good indication of how good a join is by taking the product of the values in the rows column of the EXPLAIN output. This should tell you roughly how many rows MySQL must examine to execute the query. If you restrict queries with the max\_join\_size system variable, this row product also is used to determine which multiple-table SELECT statements to execute and which to abort. See Section 7.1.1, "Configuring the Server".

The following example shows how a multiple-table join can be optimized progressively based on the information provided by EXPLAIN.

Suppose that you have the SELECT statement shown here and that you plan to examine it using EXPLAIN:

For this example, make the following assumptions:

• The columns being compared have been declared as follows.

Table	Column	Data Type
tt	ActualPC	CHAR(10)
tt	AssignedPC	CHAR(10)
tt	ClientID	CHAR(10)
et	EMPLOYID	CHAR (15)
do	CUSTNMBR	CHAR (15)

· The tables have the following indexes.

Table	Index
tt	ActualPC
tt	AssignedPC
tt	ClientID
et	EMPLOYID (primary key)
do	CUSTNMBR (primary key)

• The tt.ActualPC values are not evenly distributed.

Initially, before any optimizations have been performed, the EXPLAIN statement produces the following information:

Because type is ALL for each table, this output indicates that MySQL is generating a Cartesian product of all the tables; that is, every combination of rows. This takes quite a long time, because the product of the number of rows in each table must be examined. For the case at hand, this product is  $74 \times 2135 \times 74 \times 3872 = 45,268,558,720$  rows. If the tables were bigger, you can only imagine how long it would take.

One problem here is that MySQL can use indexes on columns more efficiently if they are declared as the same type and size. In this context, VARCHAR and CHAR are considered the same if they are declared as the same size. tt.ActualPC is declared as CHAR(10) and et.EMPLOYID is CHAR(15), so there is a length mismatch.

To fix this disparity between column lengths, use ALTER TABLE to lengthen ActualPC from 10 characters to 15 characters:

```
mysql> ALTER TABLE tt MODIFY ActualPC VARCHAR(15);
```

Now tt.ActualPC and et.EMPLOYID are both VARCHAR (15). Executing the EXPLAIN statement again produces this result:

```
table type possible_keys key
                              key_len ref
                                                 rows
                                                        Extra
    ALL
           AssignedPC, NULL
                              NULL
                                    NULL
                                                        Using
           ClientID.
                                                        where
          ActualPC
PRIMARY NULL NULL
     ALL
                                                 2135
do
                                     NULL
     Range checked for each record (index map: 0x1)
et_1 ALL PRIMARY NULL NULL NULL
     Range checked for each record (index map: 0x1)
et
     eq_ref PRIMARY PRIMARY 15
```

This is not perfect, but is much better: The product of the rows values is less by a factor of 74. This version executes in a couple of seconds.

A second alteration can be made to eliminate the column length mismatches for the tt.AssignedPC = et\_1.EMPLOYID and tt.ClientID = do.CUSTNMBR comparisons:

```
mysql> ALTER TABLE tt MODIFY AssignedPC VARCHAR(15),

MODIFY ClientID VARCHAR(15);
```

After that modification, **EXPLAIN** produces the output shown here:

```
table type possible_keys key key_len ref rows Extra
```

At this point, the query is optimized almost as well as possible. The remaining problem is that, by default, MySQL assumes that values in the tt.ActualPC column are evenly distributed, and that is not the case for the tt table. Fortunately, it is easy to tell MySQL to analyze the key distribution:

```
mysql> ANALYZE TABLE tt;
```

With the additional index information, the join is perfect and EXPLAIN produces this result:

The rows column in the output from EXPLAIN is an educated guess from the MySQL join optimizer. Check whether the numbers are even close to the truth by comparing the rows product with the actual number of rows that the query returns. If the numbers are quite different, you might get better performance by using STRAIGHT\_JOIN in your SELECT statement and trying to list the tables in a different order in the FROM clause. (However, STRAIGHT\_JOIN may prevent indexes from being used because it disables semijoin transformations. See Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations.)

It is possible in some cases to execute statements that modify data when EXPLAIN SELECT is used with a subquery; for more information, see Section 15.2.15.8, "Derived Tables".

# 10.8.3 Extended EXPLAIN Output Format

The EXPLAIN statement produces extra ("extended") information that is not part of EXPLAIN output but can be viewed by issuing a SHOW WARNINGS statement following EXPLAIN. Extended information is available for SELECT, DELETE, INSERT, REPLACE, and UPDATE statements.

The Message value in SHOW WARNINGS output displays how the optimizer qualifies table and column names in the SELECT statement, what the SELECT looks like after the application of rewriting and optimization rules, and possibly other notes about the optimization process.

The extended information displayable with a SHOW WARNINGS statement following EXPLAIN is produced only for SELECT statements. SHOW WARNINGS displays an empty result for other explainable statements (DELETE, INSERT, REPLACE, and UPDATE).

Here is an example of extended EXPLAIN output:

```
mysql> EXPLAIN
     SELECT t1.a, t1.a IN (SELECT t2.a FROM t2) FROM t1\G
************************* 1. row *****************
         id: 1
 select_type: PRIMARY
       table: t1
        type: index
possible_keys: NULL
        key: PRIMARY
     key_len: 4
        ref: NULL
        rows: 4
    filtered: 100.00
      Extra: Using index
                    ***** 2. row *****************
         id: 2
 select_type: SUBQUERY
```

```
table: t2
       type: index
possible_keys: a
        key: a
     key_len: 5
        ref: NULL
        rows: 3
    filtered: 100.00
      Extra: Using index
2 rows in set, 1 warning (0.00 sec)
mysql> SHOW WARNINGS\G
           ********** 1. row ****************
 Level: Note
  Code: 1003
( <materialize> (/* select#2 */ select `test`.`t2`.`a`
       from `test`.`t2` where 1 having 1 ),
       <primary_index_lookup>(`test`.`t1`.`a` in
        <temporary table> on <auto_key>
        where ((`test`.`t1`.`a` = `materialized-subquery`.`a`)))) AS `t1.a
       IN (SELECT t2.a FROM t2)` from `test`.`t1`
1 row in set (0.00 sec)
```

Because the statement displayed by SHOW WARNINGS may contain special markers to provide information about query rewriting or optimizer actions, the statement is not necessarily valid SQL and is not intended to be executed. The output may also include rows with Message values that provide additional non-SQL explanatory notes about actions taken by the optimizer.

The following list describes special markers that can appear in the extended output displayed by SHOW WARNINGS:

• <auto key>

An automatically generated key for a temporary table.

• <cache>(expr)

The expression (such as a scalar subquery) is executed once and the resulting value is saved in memory for later use. For results consisting of multiple values, a temporary table may be created and <temporary table> is shown instead.

<exists>(query fragment)

The subquery predicate is converted to an EXISTS predicate and the subquery is transformed so that it can be used together with the EXISTS predicate.

• <in optimizer>(query fragment)

This is an internal optimizer object with no user significance.

• <index\_lookup>(query fragment)

The query fragment is processed using an index lookup to find qualifying rows.

• <if>(condition, expr1, expr2)

If the condition is true, evaluate to *expr1*, otherwise *expr2*.

• <is\_not\_null\_test>(expr)

A test to verify that the expression does not evaluate to NULL.

• <materialize>(query fragment)

Subquery materialization is used.

• `materialized-subquery`.col\_name

A reference to the column col\_name in an internal temporary table materialized to hold the result from evaluating a subquery.

• rimary\_index\_lookup>(query fragment)

The query fragment is processed using a primary key lookup to find qualifying rows.

• <ref null helper>(expr)

This is an internal optimizer object with no user significance.

• /\* select#N \*/ select stmt

The SELECT is associated with the row in non-extended EXPLAIN output that has an id value of N.

• outer\_tables semi join (inner\_tables)

A semijoin operation. *inner\_tables* shows the tables that were not pulled out. See Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations.

• <temporary table>

This represents an internal temporary table created to cache an intermediate result.

When some tables are of const or system type, expressions involving columns from these tables are evaluated early by the optimizer and are not part of the displayed statement. However, with FORMAT=JSON, some const table accesses are displayed as a ref access that uses a const value.

## 10.8.4 Obtaining Execution Plan Information for a Named Connection

To obtain the execution plan for an explainable statement executing in a named connection, use this statement:

```
EXPLAIN [options] FOR CONNECTION connection_id;
```

EXPLAIN FOR CONNECTION returns the EXPLAIN information that is currently being used to execute a query in a given connection. Because of changes to data (and supporting statistics) it may produce a different result from running EXPLAIN on the equivalent query text. This difference in behavior can be useful in diagnosing more transient performance problems. For example, if you are running a statement in one session that is taking a long time to complete, using EXPLAIN FOR CONNECTION in another session may yield useful information about the cause of the delay.

connection\_id is the connection identifier, as obtained from the INFORMATION\_SCHEMA PROCESSLIST table or the SHOW PROCESSLIST statement. If you have the PROCESS privilege, you can specify the identifier for any connection. Otherwise, you can specify the identifier only for your own connections. In all cases, you must have sufficient privileges to explain the query on the specified connection.

If the named connection is not executing a statement, the result is empty. Otherwise, EXPLAIN FOR CONNECTION applies only if the statement being executed in the named connection is explainable. This includes SELECT, DELETE, INSERT, REPLACE, and UPDATE. (However, EXPLAIN FOR CONNECTION does not work for prepared statements, even prepared statements of those types.)

If the named connection is executing an explainable statement, the output is what you would obtain by using EXPLAIN on the statement itself.

If the named connection is executing a statement that is not explainable, an error occurs. For example, you cannot name the connection identifier for your current session because EXPLAIN is not explainable:

```
mysql> SELECT CONNECTION_ID();
+-----+
```

```
| CONNECTION_ID() |
+-----+
| 373 |
+-----+
1 row in set (0.00 sec)

mysql> EXPLAIN FOR CONNECTION 373;
ERROR 1889 (HY000): EXPLAIN FOR CONNECTION command is supported only for SELECT/UPDATE/INSERT/DELETE/REPLACE
```

The Com\_explain\_other status variable indicates the number of EXPLAIN FOR CONNECTION statements executed.

## 10.8.5 Estimating Query Performance

In most cases, you can estimate query performance by counting disk seeks. For small tables, you can usually find a row in one disk seek (because the index is probably cached). For bigger tables, you can estimate that, using B-tree indexes, you need this many seeks to find a row:  $log(row\_count) / log(index\_block\_length / 3 * 2 / (index\_length + data\_pointer\_length)) + 1.$ 

In MySQL, an index block is usually 1,024 bytes and the data pointer is usually four bytes. For a 500,000-row table with a key value length of three bytes (the size of MEDIUMINT), the formula indicates log(500,000)/log(1024/3\*2/(3+4)) + 1 = 4 seeks.

This index would require storage of about 500,000 \* 7 \* 3/2 = 5.2MB (assuming a typical index buffer fill ratio of 2/3), so you probably have much of the index in memory and so need only one or two calls to read data to find the row.

For writes, however, you need four seek requests to find where to place a new index value and normally two seeks to update the index and write the row.

The preceding discussion does not mean that your application performance slowly degenerates by  $\log N$ . As long as everything is cached by the OS or the MySQL server, things become only marginally slower as the table gets bigger. After the data gets too big to be cached, things start to go much slower until your applications are bound only by disk seeks (which increase by  $\log N$ ). To avoid this, increase the key cache size as the data grows. For MyISAM tables, the key cache size is controlled by the key\_buffer\_size system variable. See Section 7.1.1, "Configuring the Server".

# 10.9 Controlling the Query Optimizer

MySQL provides optimizer control through system variables that affect how query plans are evaluated, switchable optimizations, optimizer and index hints, and the optimizer cost model.

The server maintains histogram statistics about column values in the column\_statistics data dictionary table (see Section 10.9.6, "Optimizer Statistics"). Like other data dictionary tables, this table is not directly accessible by users. Instead, you can obtain histogram information by querying INFORMATION\_SCHEMA.COLUMN\_STATISTICS, which is implemented as a view on the data dictionary table. You can also perform histogram management using the ANALYZE TABLE statement.

# 10.9.1 Controlling Query Plan Evaluation

The task of the query optimizer is to find an optimal plan for executing an SQL query. Because the difference in performance between "good" and "bad" plans can be orders of magnitude (that is, seconds versus hours or even days), most query optimizers, including that of MySQL, perform a more or less exhaustive search for an optimal plan among all possible query evaluation plans. For join queries, the number of possible plans investigated by the MySQL optimizer grows exponentially with the number of tables referenced in a query. For small numbers of tables (typically less than 7 to 10) this is not a problem. However, when larger queries are submitted, the time spent in query optimization may easily become the major bottleneck in the server's performance.

A more flexible method for query optimization enables the user to control how exhaustive the optimizer is in its search for an optimal query evaluation plan. The general idea is that the fewer plans that are

investigated by the optimizer, the less time it spends in compiling a query. On the other hand, because the optimizer skips some plans, it may miss finding an optimal plan.

The behavior of the optimizer with respect to the number of plans it evaluates can be controlled using two system variables:

- The optimizer\_prune\_level variable tells the optimizer to skip certain plans based on estimates of the number of rows accessed for each table. Our experience shows that this kind of "educated guess" rarely misses optimal plans, and may dramatically reduce query compilation times. That is why this option is on (optimizer\_prune\_level=1) by default. However, if you believe that the optimizer missed a better query plan, this option can be switched off (optimizer\_prune\_level=0) with the risk that query compilation may take much longer. Note that, even with the use of this heuristic, the optimizer still explores a roughly exponential number of plans.
- The optimizer\_search\_depth variable tells how far into the "future" of each incomplete plan the optimizer should look to evaluate whether it should be expanded further. Smaller values of optimizer\_search\_depth may result in orders of magnitude smaller query compilation times. For example, queries with 12, 13, or more tables may easily require hours and even days to compile if optimizer\_search\_depth is close to the number of tables in the query. At the same time, if compiled with optimizer\_search\_depth equal to 3 or 4, the optimizer may compile in less than a minute for the same query. If you are unsure of what a reasonable value is for optimizer\_search\_depth, this variable can be set to 0 to tell the optimizer to determine the value automatically.

## 10.9.2 Switchable Optimizations

The optimizer\_switch system variable enables control over optimizer behavior. Its value is a set of flags, each of which has a value of on or off to indicate whether the corresponding optimizer behavior is enabled or disabled. This variable has global and session values and can be changed at runtime. The global default can be set at server startup.

To see the current set of optimizer flags, select the variable value:

To change the value of <code>optimizer\_switch</code>, assign a value consisting of a comma-separated list of one or more commands:

```
SET [GLOBAL|SESSION] optimizer_switch='command[,command]...';
```

Each command value should have one of the forms shown in the following table.

Command Syntax	Meaning
default	Reset every optimization to its default value
opt_name=default	Set the named optimization to its default value
opt_name=off	Disable the named optimization
opt_name=on	Enable the named optimization

The order of the commands in the value does not matter, although the default command is executed first if present. Setting an <code>opt\_name</code> flag to default sets it to whichever of on or off is its default value. Specifying any given <code>opt\_name</code> more than once in the value is not permitted and causes an error. Any errors in the value cause the assignment to fail with an error, leaving the value of <code>optimizer\_switch</code> unchanged.

The following list describes the permissible opt\_name flag names, grouped by optimization strategy:

- Batched Key Access Flags
  - batched\_key\_access (default off)

Controls use of BKA join algorithm.

For batched\_key\_access to have any effect when set to on, the mrr flag must also be on. Currently, the cost estimation for MRR is too pessimistic. Hence, it is also necessary for mrr\_cost\_based to be off for BKA to be used.

For more information, see Section 10.2.1.12, "Block Nested-Loop and Batched Key Access Joins".

- Block Nested-Loop Flags
  - block\_nested\_loop (default on)

Controls use of hash joins, as do the BNL and NO BNL optimizer hints.

For more information, see Section 10.2.1.12, "Block Nested-Loop and Batched Key Access Joins".

- Condition Filtering Flags
  - condition\_fanout\_filter (default on)

Controls use of condition filtering.

For more information, see Section 10.2.1.13, "Condition Filtering".

- · Derived Condition Pushdown Flags
  - derived\_condition\_pushdown (default on)

Controls derived condition pushdown.

For more information, see Section 10.2.2.5, "Derived Condition Pushdown Optimization"

- · Derived Table Merging Flags
  - derived\_merge (default on)

Controls merging of derived tables and views into outer query block.

The derived\_merge flag controls whether the optimizer attempts to merge derived tables, view references, and common table expressions into the outer query block, assuming that no other rule prevents merging; for example, an ALGORITHM directive for a view takes precedence over the derived\_merge setting. By default, the flag is on to enable merging.

For more information, see Section 10.2.2.4, "Optimizing Derived Tables, View References, and Common Table Expressions with Merging or Materialization".

- Engine Condition Pushdown Flags
  - engine\_condition\_pushdown (default on)

Controls engine condition pushdown.

For more information, see Section 10.2.1.5, "Engine Condition Pushdown Optimization".

- · Hash Join Flags
  - hash\_join (default on)

Has no effect in MySQL 8.4. Use the block\_nested\_loop flag, instead.

For more information, see Section 10.2.1.4, "Hash Join Optimization".

- Index Condition Pushdown Flags
  - index\_condition\_pushdown (default on)

Controls index condition pushdown.

For more information, see Section 10.2.1.6, "Index Condition Pushdown Optimization".

- Index Extensions Flags
  - use\_index\_extensions (default on)

Controls use of index extensions.

For more information, see Section 10.3.10, "Use of Index Extensions".

- Index Merge Flags
  - index\_merge (default on)

Controls all Index Merge optimizations.

• index\_merge\_intersection (default on)

Controls the Index Merge Intersection Access optimization.

• index\_merge\_sort\_union (default on)

Controls the Index Merge Sort-Union Access optimization.

• index\_merge\_union (default on)

Controls the Index Merge Union Access optimization.

For more information, see Section 10.2.1.3, "Index Merge Optimization".

- Index Visibility Flags
  - use\_invisible\_indexes (default off)

Controls use of invisible indexes.

For more information, see Section 10.3.12, "Invisible Indexes".

- · Limit Optimization Flags
  - prefer\_ordering\_index (default on)

Controls whether, in the case of a query having an ORDER BY or GROUP BY with a LIMIT clause, the optimizer tries to use an ordered index instead of an unordered index, a filesort, or some other

optimization. This optimization is performed by default whenever the optimizer determines that using it would allow for faster execution of the query.

Because the algorithm that makes this determination cannot handle every conceivable case (due in part to the assumption that the distribution of data is always more or less uniform), there are cases in which this optimization may not be desirable. This optimization can be disabled by setting the prefer\_ordering\_index flag to off.

For more information and examples, see Section 10.2.1.19, "LIMIT Query Optimization".

- Multi-Range Read Flags
  - mrr (default on)

Controls the Multi-Range Read strategy.

• mrr\_cost\_based (default on)

Controls use of cost-based MRR if mrr=on.

For more information, see Section 10.2.1.11, "Multi-Range Read Optimization".

- · Semijoin Flags
  - duplicateweedout (default on)

Controls the semijoin Duplicate Weedout strategy.

• firstmatch (default on)

Controls the semijoin FirstMatch strategy.

• loosescan (default on)

Controls the semijoin LooseScan strategy (not to be confused with Loose Index Scan for GROUP BY).

semijoin (default on)

Controls all semijoin strategies.

This also applies to the antijoin optimization.

The semijoin, firstmatch, loosescan, and duplicateweedout flags enable control over semijoin strategies. The semijoin flag controls whether semijoins are used. If it is set to on, the firstmatch and loosescan flags enable finer control over the permitted semijoin strategies.

If the duplicateweedout semijoin strategy is disabled, it is not used unless all other applicable strategies are also disabled.

If semijoin and materialization are both on, semijoins also use materialization where applicable. These flags are on by default.

For more information, see Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations.

- Set Operations Flags
  - hash\_set\_operations (default on)

Enables the hash table optimization for set operations involving EXCEPT and INTERSECT); enabled by default. Otherwise, temporary table based de-duplication is used, as in previous versions of MySQL.

The amount of memory used for hashing by this optimization can be controlled using the set\_operations\_buffer\_size system variable; increasing this generally results in faster execution times for statements using these operations.

- · Skip Scan Flags
  - skip\_scan (default on)

Controls use of Skip Scan access method.

For more information, see Skip Scan Range Access Method.

- Subquery Materialization Flags
  - materialization (default on)

Controls materialization (including semijoin materialization).

• subquery\_materialization\_cost\_based (default on)

Use cost-based materialization choice.

The materialization flag controls whether subquery materialization is used. If semijoin and materialization are both on, semijoins also use materialization where applicable. These flags are on by default.

The subquery\_materialization\_cost\_based flag enables control over the choice between subquery materialization and IN-to-EXISTS subquery transformation. If the flag is on (the default), the optimizer performs a cost-based choice between subquery materialization and IN-to-EXISTS subquery transformation if either method could be used. If the flag is off, the optimizer chooses subquery materialization over IN-to-EXISTS subquery transformation.

For more information, see Section 10.2.2, "Optimizing Subqueries, Derived Tables, View References, and Common Table Expressions".

- Subquery Transformation Flags
  - subquery\_to\_derived (default off)

The optimizer is able in many cases to transform a scalar subquery in a SELECT, WHERE, JOIN, or HAVING clause into a left outer join on a derived table. (Depending on the nullability of the derived table, this can sometimes be simplified further to an inner join.) This can be done for a subquery which meets the following conditions:

- The subquery does not make use of any nondeterministic functions, such as RAND().
- The subquery is not an ANY or ALL subquery which can be rewritten to use MIN() or MAX().
- The parent query does not set a user variable, since rewriting it may affect the order of execution, which could lead to unexpected results if the variable is accessed more than once in the same query.
- The subquery should not be correlated, that is, it should not reference a column from a table in the outer query, or contain an aggregate that is evaluated in the outer query.

This optimization can also be applied to a table subquery which is the argument to IN, NOT IN, EXISTS, or NOT EXISTS, that does not contain a GROUP BY.

The default value for this flag is off, since, in most cases, enabling this optimization does not produce any noticeable improvement in performance (and in many cases can even make queries run more slowly), but you can enable the optimization by setting the subquery\_to\_derived flag to on. It is primarily intended for use in testing.

Example, using a scalar subquery:

```
mysql> CREATE TABLE t1(a INT);
mysql> CREATE TABLE t2(a INT);
mysql> INSERT INTO t1 VALUES ROW(1), ROW(2), ROW(3), ROW(4);
mysql> INSERT INTO t2 VALUES ROW(1), ROW(2);
mysql> SELECT * FROM t1
          WHERE tl.a > (SELECT COUNT(a) FROM t2);
| a |
    3
    4
mysql> SELECT @@optimizer_switch LIKE '%subquery_to_derived=off%';
| @@optimizer_switch LIKE '%subquery_to_derived=off%' |
mysql> EXPLAIN SELECT * FROM t1 WHERE t1.a > (SELECT COUNT(a) FROM t2)\G
               ********** 1. row *******
       id: 1
  select_type: PRIMARY
       table: t1
  partitions: NULL
        type: ALL
possible_keys: NULL
         key: NULL
     key_len: NULL
         ref: NULL
        rows: 4
    filtered: 33.33
      Extra: Using where
************************ 2. row ******************
          id: 2
  select_type: SUBQUERY
       table: t2
  partitions: NULL
        type: ALL
possible_keys: NULL
         key: NULL
     key_len: NULL
         ref: NULL
        rows: 2
     filtered: 100.00
       Extra: NULL
mysql> SET @@optimizer_switch='subquery_to_derived=on';
mysql> SELECT @@optimizer_switch LIKE '%subquery_to_derived=off%';
| @@optimizer_switch LIKE '%subquery_to_derived=off%' |
```

```
_____
mysql> SELECT @@optimizer_switch LIKE '%subquery_to_derived=on%';
| @@optimizer_switch LIKE '%subquery_to_derived=on%' |
                                          1 |
mysql> EXPLAIN SELECT * FROM t1 WHERE t1.a > (SELECT COUNT(a) FROM t2)\G
 id: 1
 select_type: PRIMARY
      table: <derived2>
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
     key_len: NULL
       ref: NULL
       rows: 1
    filtered: 100.00
     Extra: NULL
************************ 2. row *****************
       id: 1
 select_type: PRIMARY
      table: t1
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
    key_len: NULL
        ref: NULL
       rows: 4
    filtered: 33.33
      Extra: Using where; Using join buffer (hash join)
****** 3. row *****
        id: 2
 select_type: DERIVED
      table: t2
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
     key_len: NULL
        ref: NULL
       rows: 2
    filtered: 100.00
      Extra: NULL
```

As can be seen from executing SHOW WARNINGS immediately following the second EXPLAIN statement, with the optimization enabled, the query SELECT \* FROM t1 WHERE t1.a > (SELECT COUNT(a) FROM t2) is rewritten in a form similar to what is shown here:

```
SELECT t1.a FROM t1

JOIN ( SELECT COUNT(t2.a) AS c FROM t2 ) AS d

WHERE t1.a > d.c;
```

Example, using a query with IN (subquery):

```
mysql> SELECT * FROM t1
   -> WHERE t1.b < 0
               OR
               t1.a IN (SELECT t2.a + 1 FROM t2);
  ->
| a | b |
   2 | 20 |
  3 | 30 |
mysql> SET @@optimizer_switch="subquery_to_derived=off";
mysql> EXPLAIN SELECT * FROM t1
       WHERE t1.b < 0
   ->
                     OR
                     t1.a IN (SELECT t2.a + 1 FROM t2)\G
id: 1
 select_type: PRIMARY
      table: t1
  partitions: NULL
      type: ALL
possible_keys: NULL
       key: NULL
    key_len: NULL
       ref: NULL
       rows: 3
    filtered: 100.00
     Extra: Using where
id: 2
 select_type: DEPENDENT SUBQUERY
      table: t2
  partitions: NULL
      type: ALL
possible_keys: NULL
       key: NULL
    key_len: NULL
       ref: NULL
       rows: 6
    filtered: 100.00
      Extra: Using where
mysql> SET @@optimizer_switch="subquery_to_derived=on";
mysql> EXPLAIN SELECT * FROM t1
       WHERE t1.b < 0
                     OR
   ->
                    t1.a IN (SELECT t2.a + 1 FROM t2)\G
id: 1
 select_type: PRIMARY
     table: t1
  partitions: NULL
      type: ALL
possible_keys: NULL
       key: NULL
    key_len: NULL
       ref: NULL
       rows: 3
   filtered: 100.00
     Extra: NULL
************************* 2. row *******************
      id: 1
 select_type: PRIMARY
     table: <derived2>
  partitions: NULL
       type: ref
possible_keys: <auto_key0>
       key: <auto_key0>
    key_len: 9
```

```
ref: std2.t1.a
      rows: 2
   filtered: 100.00
     Extra: Using where; Using index
id: 2
 select_type: DERIVED
     table: t2
  partitions: NULL
      type: ALL
possible_keys: NULL
      key: NULL
    key_len: NULL
       ref: NULL
      rows: 6
   filtered: 100.00
      Extra: Using temporary
```

Checking and simplifying the result of SHOW WARNINGS after executing EXPLAIN on this query shows that, when the subquery\_to\_derived flag enabled, SELECT \* FROM t1 WHERE t1.b < 0 OR t1.a IN (SELECT t2.a + 1 FROM t2) is rewritten in a form similar to what is shown here:

```
SELECT a, b FROM t1

LEFT JOIN (SELECT DISTINCT a + 1 AS e FROM t2) d

ON t1.a = d.e

WHERE t1.b < 0

OR

d.e IS NOT NULL;
```

Example, using a query with EXISTS (subquery) and the same tables and data as in the previous example:

```
mysql> SELECT * FROM t1
   -> WHERE t1.b < 0
               OR
   ->
               EXISTS(SELECT * FROM t2 WHERE t2.a = t1.a + 1);
| a | b |
   1 | 10 |
2 | 20 |
mysql> SET @@optimizer_switch="subquery_to_derived=off";
mysql> EXPLAIN SELECT * FROM t1
      WHERE t1.b < 0
  ->
              OR
EXISTS(SELECT * FROM t2 WHERE t2.a = t1.a + 1)\G
   ->
id: 1
 select_type: PRIMARY
      table: t1
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
     key_len: NULL
       ref: NULL
       rows: 3
    filtered: 100.00
     Extra: Using where
************************ 2. row ******************
        id: 2
 select_type: DEPENDENT SUBQUERY
      table: t2
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
```

```
key_len: NULL
       ref: NULL
       rows: 6
    filtered: 16.67
       Extra: Using where
mysql> SET @@optimizer_switch="subquery_to_derived=on";
mysql> EXPLAIN SELECT * FROM t1
         WHERE t1.b < 0
   ->
                       OR
                      EXISTS(SELECT * FROM t2 WHERE t2.a = t1.a + 1)\G
id: 1
 select_type: PRIMARY
      table: t1
  partitions: NULL
       type: ALL
possible_keys: NULL
        key: NULL
     key_len: NULL
       ref: NULL
       rows: 3
    filtered: 100.00
      Extra: NULL
************************ 2. row *****************
        id: 1
 select_type: PRIMARY
      table: <derived2>
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
     key_len: NULL
        ref: NULL
       rows: 6
    filtered: 100.00
      Extra: Using where; Using join buffer (hash join)
****** 3. row ******
        id: 2
 select_type: DERIVED
      table: t2
  partitions: NULL
       type: ALL
possible_keys: NULL
       key: NULL
     key_len: NULL
        ref: NULL
       rows: 6
    filtered: 100.00
       Extra: Using temporary
```

If we execute SHOW WARNINGS after running EXPLAIN on the query SELECT \* FROM t1 WHERE t1.b < 0 OR EXISTS(SELECT \* FROM t2 WHERE t2.a = t1.a + 1) when subquery\_to\_derived has been enabled, and simplify the second row of the result, we see that it has been rewritten in a form which resembles this:

```
SELECT a, b FROM t1

LEFT JOIN (SELECT DISTINCT 1 AS e1, t2.a AS e2 FROM t2) d

ON t1.a + 1 = d.e2

WHERE t1.b < 0

OR
d.e1 IS NOT NULL;
```

For more information, see Section 10.2.2.4, "Optimizing Derived Tables, View References, and Common Table Expressions with Merging or Materialization", as well as Section 10.2.1.19, "LIMIT Query Optimization", and Optimizing IN and EXISTS Subquery Predicates with Semijoin Transformations.

When you assign a value to <code>optimizer\_switch</code>, flags that are not mentioned keep their current values. This makes it possible to enable or disable specific optimizer behaviors in a single statement without affecting other behaviors. The statement does not depend on what other optimizer flags exist and what their values are. Suppose that all Index Merge optimizations are enabled:

If the server is using the Index Merge Union or Index Merge Sort-Union access methods for certain queries and you want to check whether the optimizer can perform better without them, set the variable value like this:

# 10.9.3 Optimizer Hints

One means of control over optimizer strategies is to set the <code>optimizer\_switch</code> system variable (see Section 10.9.2, "Switchable Optimizations"). Changes to this variable affect execution of all subsequent queries; to affect one query differently from another, it is necessary to change <code>optimizer\_switch</code> before each one.

Another way to control the optimizer is by using optimizer hints, which can be specified within individual statements. Because optimizer hints apply on a per-statement basis, they provide finer control over statement execution plans than can be achieved using optimizer\_switch. For example, you can enable an optimization for one table in a statement and disable the optimization for a different table. Hints within a statement take precedence over optimizer switch flags.

#### Examples:

```
SELECT /*+ NO_RANGE_OPTIMIZATION(t3 PRIMARY, f2_idx) */ f1
FROM t3 WHERE f1 > 30 AND f1 < 33;

SELECT /*+ BKA(t1) NO_BKA(t2) */ * FROM t1 INNER JOIN t2 WHERE ...;

SELECT /*+ NO_ICP(t1, t2) */ * FROM t1 INNER JOIN t2 WHERE ...;

SELECT /*+ SEMIJOIN(FIRSTMATCH, LOOSESCAN) */ * FROM t1 ...;

EXPLAIN SELECT /*+ NO_ICP(t1) */ * FROM t1 WHERE ...;

SELECT /*+ MERGE(dt) */ * FROM (SELECT * FROM t1) AS dt;

INSERT /*+ SET_VAR(foreign_key_checks=OFF) */ INTO t2 VALUES(2);
```

Optimizer hints, described here, differ from index hints, described in Section 10.9.4, "Index Hints". Optimizer and index hints may be used separately or together.

Optimizer Hint Overview

- Optimizer Hint Syntax
- Join-Order Optimizer Hints
- Table-Level Optimizer Hints
- Index-Level Optimizer Hints
- Subquery Optimizer Hints
- Statement Execution Time Optimizer Hints
- Variable-Setting Hint Syntax
- · Resource Group Hint Syntax
- · Optimizer Hints for Naming Query Blocks

## **Optimizer Hint Overview**

Optimizer hints apply at different scope levels:

- · Global: The hint affects the entire statement
- Query block: The hint affects a particular query block within a statement
- · Table-level: The hint affects a particular table within a query block
- Index-level: The hint affects a particular index within a table

The following table summarizes the available optimizer hints, the optimizer strategies they affect, and the scope or scopes at which they apply. More details are given later.

**Table 10.2 Optimizer Hints Available** 

Hint Name	Description	Applicable Scopes
BKA, NO_BKA	Affects Batched Key Access join processing	Query block, table
BNL, NO_BNL	Affects hash join optimization	Query block, table
DERIVED_CONDITION_PUSHDOW NO_DERIVED_CONDITION_PUSH	_	Query block, table
GROUP_INDEX, NO_GROUP_INDEX	Use or ignore the specified index or indexes for index scans in GROUP BY operations	Index
HASH_JOIN, NO_HASH_JOIN	Affects Hash Join optimization (No effect in MySQL 8.4)	Query block, table
INDEX, NO_INDEX	Acts as the combination of JOIN_INDEX, GROUP_INDEX, and ORDER_INDEX, or as the combination of NO_JOIN_INDEX, NO_GROUP_INDEX, and NO_ORDER_INDEX	Index
INDEX_MERGE, NO_INDEX_MERGE	Affects Index Merge optimization	Table, index
JOIN_FIXED_ORDER	Use table order specified in FROM clause for join order	Query block

Hint Name	Description	Applicable Scopes
JOIN_INDEX, NO_JOIN_INDEX	Use or ignore the specified index or indexes for any access method	Index
JOIN_ORDER	Use table order specified in hint for join order	Query block
JOIN_PREFIX	Use table order specified in hint for first tables of join order	Query block
JOIN_SUFFIX	Use table order specified in hint for last tables of join order	Query block
MAX_EXECUTION_TIME	Limits statement execution time	Global
MERGE, NO_MERGE	Affects derived table/view merging into outer query block	Table
MRR, NO_MRR	Affects Multi-Range Read optimization	Table, index
NO_ICP	Affects Index Condition Pushdown optimization	Table, index
NO_RANGE_OPTIMIZATION	Affects range optimization	Table, index
ORDER_INDEX, NO_ORDER_INDEX	Use or ignore the specified index or indexes for sorting rows	Index
QB_NAME	Assigns name to query block	Query block
RESOURCE_GROUP	Set resource group during statement execution	Global
SEMIJOIN, NO_SEMIJOIN	Affects semijoin and antijoin strategies	Query block
SKIP_SCAN, NO_SKIP_SCAN	Affects Skip Scan optimization	Table, index
SET_VAR	Set variable during statement execution	Global
SUBQUERY	Affects materialization, IN- to-EXISTS subquery strategies	Query block

Disabling an optimization prevents the optimizer from using it. Enabling an optimization means the optimizer is free to use the strategy if it applies to statement execution, not that the optimizer necessarily uses it.

## **Optimizer Hint Syntax**

MySQL supports comments in SQL statements as described in Section 11.7, "Comments". Optimizer hints must be specified within /\*+ . . . \*/ comments. That is, optimizer hints use a variant of / \* . . . \*/ C-style comment syntax, with a + character following the /\* comment opening sequence. Examples:

```
/*+ BKA(t1) */
/*+ BNL(t1, t2) */
/*+ NO_RANGE_OPTIMIZATION(t4 PRIMARY) */
/*+ QB_NAME(qb2) */
```

Whitespace is permitted after the + character.

The parser recognizes optimizer hint comments after the initial keyword of SELECT, UPDATE, INSERT, REPLACE, and DELETE statements. Hints are permitted in these contexts:

At the beginning of query and data change statements:

```
SELECT /*+ ... */ ...
INSERT /*+ ... */ ...
REPLACE /*+ ... */ ...
UPDATE /*+ ... */ ...
DELETE /*+ ... */ ...
```

· At the beginning of query blocks:

```
(SELECT /*+ ... */ ... )
(SELECT ... ) UNION (SELECT /*+ ... */ ... )
(SELECT /*+ ... */ ... ) UNION (SELECT /*+ ... */ ... )
UPDATE ... WHERE x IN (SELECT /*+ ... */ ...)
INSERT ... SELECT /*+ ... */ ...
```

• In hintable statements prefaced by EXPLAIN. For example:

```
EXPLAIN SELECT /*+ ... */ ...
EXPLAIN UPDATE ... WHERE x IN (SELECT /*+ ... */ ...)
```

The implication is that you can use EXPLAIN to see how optimizer hints affect execution plans. Use SHOW WARNINGS immediately after EXPLAIN to see how hints are used. The extended EXPLAIN output displayed by a following SHOW WARNINGS indicates which hints were used. Ignored hints are not displayed.

A hint comment may contain multiple hints, but a query block cannot contain multiple hint comments. This is valid:

```
SELECT /*+ BNL(t1) BKA(t2) */ ...
```

But this is invalid:

```
SELECT /*+ BNL(t1) */ /* BKA(t2) */ ...
```

When a hint comment contains multiple hints, the possibility of duplicates and conflicts exists. The following general guidelines apply. For specific hint types, additional rules may apply, as indicated in the hint descriptions.

- Duplicate hints: For a hint such as /\*+ MRR(idx1) MRR(idx1) \*/, MySQL uses the first hint and issues a warning about the duplicate hint.
- Conflicting hints: For a hint such as /\*+ MRR(idx1) NO\_MRR(idx1) \*/, MySQL uses the first hint and issues a warning about the second conflicting hint.

Query block names are identifiers and follow the usual rules about what names are valid and how to quote them (see Section 11.2, "Schema Object Names").

Hint names, query block names, and strategy names are not case-sensitive. References to table and index names follow the usual identifier case-sensitivity rules (see Section 11.2.3, "Identifier Case Sensitivity").

## **Join-Order Optimizer Hints**

Join-order hints affect the order in which the optimizer joins tables.

Syntax of the JOIN FIXED ORDER hint:

```
hint_name([@query_block_name])
```

Syntax of other join-order hints:

```
hint_name([@query_block_name] tbl_name [, tbl_name] ...)
hint_name(tbl_name[@query_block_name] [, tbl_name[@query_block_name]] ...)
```

The syntax refers to these terms:

- hint\_name: These hint names are permitted:
  - JOIN\_FIXED\_ORDER: Force the optimizer to join tables using the order in which they appear in the FROM clause. This is the same as specifying SELECT STRAIGHT\_JOIN.
  - JOIN\_ORDER: Instruct the optimizer to join tables using the specified table order. The hint applies to the named tables. The optimizer may place tables that are not named anywhere in the join order, including between specified tables.
  - JOIN\_PREFIX: Instruct the optimizer to join tables using the specified table order for the first tables of the join execution plan. The hint applies to the named tables. The optimizer places all other tables after the named tables.
  - JOIN\_SUFFIX: Instruct the optimizer to join tables using the specified table order for the last tables of the join execution plan. The hint applies to the named tables. The optimizer places all other tables before the named tables.
- tb1\_name: The name of a table used in the statement. A hint that names tables applies to all tables that it names. The JOIN\_FIXED\_ORDER hint names no tables and applies to all tables in the FROM clause of the query block in which it occurs.

If a table has an alias, hints must refer to the alias, not the table name.

Table names in hints cannot be qualified with schema names.

query\_block\_name: The query block to which the hint applies. If the hint includes no leading @query\_block\_name, the hint applies to the query block in which it occurs. For tbl\_name@query\_block\_name syntax, the hint applies to the named table in the named query block. To assign a name to a query block, see Optimizer Hints for Naming Query Blocks.

#### Example:

Hints control the behavior of semijoin tables that are merged to the outer query block. If subqueries subq1 and subq2 are converted to semijoins, tables t4@subq1 and t5@subq2 are merged to the outer query block. In this case, the hint specified in the outer query block controls the behavior of t4@subq1, t5@subq2 tables.

The optimizer resolves join-order hints according to these principles:

· Multiple hint instances

Only one JOIN\_PREFIX and JOIN\_SUFFIX hint of each type are applied. Any later hints of the same type are ignored with a warning. JOIN\_ORDER can be specified several times.

#### Examples:

```
/*+ JOIN_PREFIX(t1) JOIN_PREFIX(t2) */
```

The second JOIN\_PREFIX hint is ignored with a warning.

```
/*+ JOIN_PREFIX(t1) JOIN_SUFFIX(t2) */
```

Both hints are applicable. No warning occurs.

```
/*+ JOIN_ORDER(t1, t2) JOIN_ORDER(t2, t3) */
```

Both hints are applicable. No warning occurs.

Conflicting hints

In some cases hints can conflict, such as when JOIN\_ORDER and JOIN\_PREFIX have table orders that are impossible to apply at the same time:

```
SELECT /*+ JOIN_ORDER(t1, t2) JOIN_PREFIX(t2, t1) */ ... FROM t1, t2;
```

In this case, the first specified hint is applied and subsequent conflicting hints are ignored with no warning. A valid hint that is impossible to apply is silently ignored with no warning.

· Ignored hints

A hint is ignored if a table specified in the hint has a circular dependency.

#### Example:

```
/*+ JOIN_ORDER(t1, t2) JOIN_PREFIX(t2, t1) */
```

The JOIN\_ORDER hint sets table t2 dependent on t1. The JOIN\_PREFIX hint is ignored because table t1 cannot be dependent on t2. Ignored hints are not displayed in extended EXPLAIN output.

Interaction with const tables

The MySQL optimizer places const tables first in the join order, and the position of a const table cannot be affected by hints. References to const tables in join-order hints are ignored, although the hint is still applicable. For example, these are equivalent:

```
JOIN_ORDER(t1, const_tbl, t2)
JOIN_ORDER(t1, t2)
```

Accepted hints shown in extended EXPLAIN output include const tables as they were specified.

· Interaction with types of join operations

MySQL supports several type of joins: LEFT, RIGHT, INNER, CROSS, STRAIGHT\_JOIN. A hint that conflicts with the specified type of join is ignored with no warning.

#### Example:

```
SELECT /*+ JOIN_PREFIX(t1, t2) */FROM t2 LEFT JOIN t1;
```

Here a conflict occurs between the requested join order in the hint and the order required by the LEFT JOIN. The hint is ignored with no warning.

### **Table-Level Optimizer Hints**

Table-level hints affect:

- Use of the Block Nested-Loop (BNL) and Batched Key Access (BKA) join-processing algorithms (see Section 10.2.1.12, "Block Nested-Loop and Batched Key Access Joins").
- Whether derived tables, view references, or common table expressions should be merged into the outer query block, or materialized using an internal temporary table.
- Use of the derived table condition pushdown optimization. See Section 10.2.2.5, "Derived Condition Pushdown Optimization".

These hint types apply to specific tables, or all tables in a guery block.

#### Syntax of table-level hints:

```
hint_name([@query_block_name] [tbl_name [, tbl_name] ...])
hint_name([tbl_name@query_block_name [, tbl_name@query_block_name] ...])
```

The syntax refers to these terms:

- hint name: These hint names are permitted:
  - BKA, NO\_BKA: Enable or disable batched key access for the specified tables.
  - BNL, NO\_BNL: Enable and disable the hash join optimization.
  - DERIVED\_CONDITION\_PUSHDOWN, NO\_DERIVED\_CONDITION\_PUSHDOWN: Enable or disable use of derived table condition pushdown for the specified tables. For more information, see Section 10.2.2.5, "Derived Condition Pushdown Optimization".
  - HASH\_JOIN, NO\_HASH\_JOIN: These hints have no effect in MySQL 8.4; use BNL or NO\_BNL instead.
  - MERGE, NO\_MERGE: Enable merging for the specified tables, view references or common table expressions; or disable merging and use materialization instead.



#### Note

To use a block nested loop or batched key access hint to enable join buffering for any inner table of an outer join, join buffering must be enabled for all inner tables of the outer join.

• *tbl\_name*: The name of a table used in the statement. The hint applies to all tables that it names. If the hint names no tables, it applies to all tables of the query block in which it occurs.

If a table has an alias, hints must refer to the alias, not the table name.

Table names in hints cannot be qualified with schema names.

• query\_block\_name: The query block to which the hint applies. If the hint includes no leading @query\_block\_name, the hint applies to the query block in which it occurs. For tbl\_name@query\_block\_name syntax, the hint applies to the named table in the named query block. To assign a name to a query block, see Optimizer Hints for Naming Query Blocks.

### Examples:

```
SELECT /*+ NO_BKA(t1, t2) */ t1.* FROM t1 INNER JOIN t2 INNER JOIN t3;

SELECT /*+ NO_BNL() BKA(t1) */ t1.* FROM t1 INNER JOIN t2 INNER JOIN t3;

SELECT /*+ NO_MERGE(dt) */ * FROM (SELECT * FROM t1) AS dt;
```

A table-level hint applies to tables that receive records from previous tables, not sender tables. Consider this statement:

```
SELECT /*+ BNL(t2) */ FROM t1, t2;
```

If the optimizer chooses to process t1 first, it applies a Block Nested-Loop join to t2 by buffering the rows from t1 before starting to read from t2. If the optimizer instead chooses to process t2 first, the hint has no effect because t2 is a sender table.

For the MERGE and NO MERGE hints, these precedence rules apply:

- A hint takes precedence over any optimizer heuristic that is not a technical constraint. (If providing a hint as a suggestion has no effect, the optimizer has a reason for ignoring it.)
- A hint takes precedence over the derived\_merge flag of the optimizer\_switch system variable.

• For view references, an ALGORITHM= {MERGE | TEMPTABLE} clause in the view definition takes precedence over a hint specified in the query referencing the view.

## **Index-Level Optimizer Hints**

Index-level hints affect which index-processing strategies the optimizer uses for particular tables or indexes. These hint types affect use of Index Condition Pushdown (ICP), Multi-Range Read (MRR), Index Merge, and range optimizations (see Section 10.2.1, "Optimizing SELECT Statements").

Syntax of index-level hints:

```
hint_name([@query_block_name] tbl_name [index_name [, index_name] ...])
hint_name(tbl_name@query_block_name [index_name [, index_name] ...])
```

The syntax refers to these terms:

- hint\_name: These hint names are permitted:
  - GROUP\_INDEX, NO\_GROUP\_INDEX: Enable or disable the specified index or indexes for index scans for GROUP BY operations. Equivalent to the index hints FORCE INDEX FOR GROUP BY, IGNORE INDEX FOR GROUP BY.
  - INDEX, NO\_INDEX: Acts as the combination of JOIN\_INDEX, GROUP\_INDEX, and ORDER\_INDEX, forcing the server to use the specified index or indexes for any and all scopes, or as the combination of NO\_JOIN\_INDEX, NO\_GROUP\_INDEX, and NO\_ORDER\_INDEX, which causes the server to ignore the specified index or indexes for any and all scopes. Equivalent to FORCE\_INDEX, IGNORE\_INDEX.
  - INDEX\_MERGE, NO\_INDEX\_MERGE: Enable or disable the Index Merge access method for the specified table or indexes. For information about this access method, see Section 10.2.1.3, "Index Merge Optimization". These hints apply to all three Index Merge algorithms.

The INDEX\_MERGE hint forces the optimizer to use Index Merge for the specified table using the specified set of indexes. If no index is specified, the optimizer considers all possible index combinations and selects the least expensive one. The hint may be ignored if the index combination is inapplicable to the given statement.

The NO\_INDEX\_MERGE hint disables Index Merge combinations that involve any of the specified indexes. If the hint specifies no indexes, Index Merge is not permitted for the table.

- JOIN\_INDEX, NO\_JOIN\_INDEX: Forces MySQL to use or ignore the specified index or indexes for any access method, such as ref, range, index\_merge, and so on. Equivalent to FORCE INDEX FOR JOIN, IGNORE INDEX FOR JOIN.
- MRR, NO\_MRR: Enable or disable MRR for the specified table or indexes. MRR hints apply only to InnoDB and MyISAM tables. For information about this access method, see Section 10.2.1.11, "Multi-Range Read Optimization".
- NO\_ICP: Disable ICP for the specified table or indexes. By default, ICP is a candidate optimization strategy, so there is no hint for enabling it. For information about this access method, see Section 10.2.1.6, "Index Condition Pushdown Optimization".
- NO\_RANGE\_OPTIMIZATION: Disable index range access for the specified table or indexes. This hint also disables Index Merge and Loose Index Scan for the table or indexes. By default, range access is a candidate optimization strategy, so there is no hint for enabling it.

This hint may be useful when the number of ranges may be high and range optimization would require many resources.

• ORDER\_INDEX, NO\_ORDER\_INDEX: Cause MySQL to use or to ignore the specified index or indexes for sorting rows. Equivalent to FORCE INDEX FOR ORDER BY, IGNORE INDEX FOR ORDER BY.

 SKIP\_SCAN, NO\_SKIP\_SCAN: Enable or disable the Skip Scan access method for the specified table or indexes. For information about this access method, see Skip Scan Range Access Method.

The SKIP\_SCAN hint forces the optimizer to use Skip Scan for the specified table using the specified set of indexes. If no index is specified, the optimizer considers all possible indexes and selects the least expensive one. The hint may be ignored if the index is inapplicable to the given statement.

The NO\_SKIP\_SCAN hint disables Skip Scan for the specified indexes. If the hint specifies no indexes, Skip Scan is not permitted for the table.

- tbl\_name: The table to which the hint applies.
- index\_name: The name of an index in the named table. The hint applies to all indexes that it names. If the hint names no indexes, it applies to all indexes in the table.

To refer to a primary key, use the name PRIMARY. To see the index names for a table, use SHOW INDEX.

• query\_block\_name: The query block to which the hint applies. If the hint includes no leading @query\_block\_name, the hint applies to the query block in which it occurs. For tbl\_name@query\_block\_name syntax, the hint applies to the named table in the named query block. To assign a name to a query block, see Optimizer Hints for Naming Query Blocks.

#### Examples:

```
SELECT /*+ INDEX_MERGE(t1 f3, PRIMARY) */ f2 FROM t1
    WHERE f1 = 'o' AND f2 = f3 AND f3 <= 4;
SELECT /*+ MRR(t1) */ * FROM t1 WHERE f2 <= 3 AND 3 <= f3;
SELECT /*+ NO_RANGE_OPTIMIZATION(t3 PRIMARY, f2_idx) */ f1
    FROM t3 WHERE f1 > 30 AND f1 < 33;
INSERT INTO t3(f1, f2, f3)
    (SELECT /*+ NO_ICP(t2) */ t2.f1, t2.f2, t2.f3 FROM t1,t2
    WHERE t1.f1=t2.f1 AND t2.f2 BETWEEN t1.f1
    AND t1.f2 AND t2.f2 + 1 >= t1.f1 + 1);
SELECT /*+ SKIP_SCAN(t1 PRIMARY) */ f1, f2
    FROM t1 WHERE f2 > 40;
```

The following examples use the Index Merge hints, but other index-level hints follow the same principles regarding hint ignoring and precedence of optimizer hints in relation to the optimizer\_switch system variable or index hints.

Assume that table t1 has columns a, b, c, and d; and that indexes named  $i_a$ ,  $i_b$ , and  $i_c$  exist on a, b, and c, respectively:

```
SELECT /*+ INDEX_MERGE(t1 i_a, i_b, i_c)*/ * FROM t1
WHERE a = 1 AND b = 2 AND c = 3 AND d = 4;
```

Index Merge is used for (i\_a, i\_b, i\_c) in this case.

```
SELECT /*+ INDEX_MERGE(t1 i_a, i_b, i_c)*/ * FROM t1
WHERE b = 1 AND c = 2 AND d = 3;
```

Index Merge is used for (i\_b, i\_c) in this case.

```
/*+ INDEX_MERGE(t1 i_a, i_b) NO_INDEX_MERGE(t1 i_b) */
```

NO\_INDEX\_MERGE is ignored because there is a preceding hint for the same table.

```
/*+ NO_INDEX_MERGE(t1 i_a, i_b) INDEX_MERGE(t1 i_b) */
```

INDEX\_MERGE is ignored because there is a preceding hint for the same table.

For the INDEX\_MERGE and NO\_INDEX\_MERGE optimizer hints, these precedence rules apply:

• If an optimizer hint is specified and is applicable, it takes precedence over the Index Merge-related flags of the optimizer\_switch system variable.

```
SET optimizer_switch='index_merge_intersection=off';
SELECT /*+ INDEX_MERGE(t1 i_b, i_c) */ * FROM t1
WHERE b = 1 AND c = 2 AND d = 3;
```

The hint takes precedence over  $optimizer\_switch$ . Index Merge is used for  $(i\_b, i\_c)$  in this case.

```
SET optimizer_switch='index_merge_intersection=on';
SELECT /*+ INDEX_MERGE(t1 i_b) */ * FROM t1
WHERE b = 1 AND c = 2 AND d = 3;
```

The hint specifies only one index, so it is inapplicable, and the optimizer\_switch flag (on) applies. Index Merge is used if the optimizer assesses it to be cost efficient.

```
SET optimizer_switch='index_merge_intersection=off';
SELECT /*+ INDEX_MERGE(t1 i_b) */ * FROM t1
WHERE b = 1 AND c = 2 AND d = 3;
```

The hint specifies only one index, so it is inapplicable, and the optimizer\_switch flag (off) applies. Index Merge is not used.

• The index-level optimizer hints <code>GROUP\_INDEX</code>, <code>INDEX</code>, <code>JOIN\_INDEX</code>, and <code>ORDER\_INDEX</code> all take precedence over the equivalent <code>FORCE INDEX</code> hints; that is, they cause the <code>FORCE INDEX</code> hints to be ignored. Likewise, the <code>NO\_GROUP\_INDEX</code>, <code>NO\_INDEX</code>, <code>NO\_JOIN\_INDEX</code>, and <code>NO\_ORDER\_INDEX</code> hints all take precedence over any <code>IGNORE INDEX</code> equivalents, also causing them to be ignored.

The index-level optimizer hints <code>GROUP\_INDEX</code>, <code>NO\_GROUP\_INDEX</code>, <code>INDEX</code>, <code>NO\_INDEX</code>, <code>JOIN\_INDEX</code>, <code>ORDER\_INDEX</code>, and <code>NO\_ORDER\_INDEX</code> hints all take precedence over all other optimizer hints, including other index-level optimizer hints. Any other optimizer hints are applied only to the indexes permitted by these.

The GROUP\_INDEX, INDEX, JOIN\_INDEX, and ORDER\_INDEX hints are all equivalent to FORCE INDEX and not to USE INDEX. This is because using one or more of these hints means that a table scan is used only if there is no way to use one of the named indexes to find rows in the table. To cause MySQL to use the same index or set of indexes as with a given instance of USE INDEX, you can use NO\_INDEX, NO\_JOIN\_INDEX, NO\_GROUP\_INDEX, NO\_ORDER\_INDEX, or some combination of these.

To replicate the effect that USE INDEX has in the query SELECT a,c FROM t1 USE INDEX FOR ORDER BY  $(i_a)$  ORDER BY a, you can use the NO\_ORDER\_INDEX optimizer hint to cover all indexes on the table except the one that is desired like this:

```
SELECT /*+ NO_ORDER_INDEX(t1 i_b,i_c) */ a,c

FROM t1

ORDER BY a;
```

Attempting to combine NO\_ORDER\_INDEX for the table as a whole with USE INDEX FOR ORDER BY does not work to do this, because NO\_ORDER\_BY causes USE INDEX to be ignored, as shown here:

```
mysql> EXPLAIN SELECT /*+ NO_ORDER_INDEX(t1) */ a,c FROM t1
    -> USE INDEX FOR ORDER BY (i_a) ORDER BY a\G
*****************************
    id: 1
select_type: SIMPLE
    table: t1
partitions: NULL
    type: ALL
possible_keys: NULL
    key: NULL
```

```
key_len: NULL
ref: NULL
rows: 256
filtered: 100.00
Extra: Using filesort
```

• The USE INDEX, FORCE INDEX, and IGNORE INDEX index hints have higher priority than the INDEX\_MERGE and NO\_INDEX\_MERGE optimizer hints.

```
/*+ INDEX_MERGE(t1 i_a, i_b, i_c) */ ... IGNORE INDEX i_a
```

IGNORE INDEX takes precedence over INDEX\_MERGE, so index i\_a is excluded from the possible ranges for Index Merge.

```
/*+ NO_INDEX_MERGE(t1 i_a, i_b) */ ... FORCE INDEX i_a, i_b
```

Index Merge is disallowed for i\_a, i\_b because of FORCE INDEX, but the optimizer is forced to use either i\_a or i\_b for range or ref access. There are no conflicts; both hints are applicable.

- If an IGNORE INDEX hint names multiple indexes, those indexes are unavailable for Index Merge.
- The FORCE INDEX and USE INDEX hints make only the named indexes to be available for Index Merge.

```
SELECT /*+ INDEX_MERGE(t1 i_a, i_b, i_c) */ a FROM t1
FORCE INDEX (i_a, i_b) WHERE c = 'h' AND a = 2 AND b = 'b';
```

The Index Merge intersection access algorithm is used for  $(i_a, i_b)$ . The same is true if FORCE INDEX is changed to USE INDEX.

## **Subquery Optimizer Hints**

Subquery hints affect whether to use semijoin transformations and which semijoin strategies to permit, and, when semijoins are not used, whether to use subquery materialization or IN-to-EXISTS transformations. For more information about these optimizations, see Section 10.2.2, "Optimizing Subqueries, Derived Tables, View References, and Common Table Expressions".

Syntax of hints that affect semijoin strategies:

```
hint_name([@query_block_name] [strategy [, strategy] ...])
```

The syntax refers to these terms:

- hint name: These hint names are permitted:
  - SEMIJOIN, NO\_SEMIJOIN: Enable or disable the named semijoin strategies.
- strategy: A semijoin strategy to be enabled or disabled. These strategy names are permitted: DUPSWEEDOUT, FIRSTMATCH, LOOSESCAN, MATERIALIZATION.

For SEMIJOIN hints, if no strategies are named, semijoin is used if possible based on the strategies enabled according to the <code>optimizer\_switch</code> system variable. If strategies are named but inapplicable for the statement, <code>DUPSWEEDOUT</code> is used.

For NO\_SEMIJOIN hints, if no strategies are named, semijoin is not used. If strategies are named that rule out all applicable strategies for the statement, DUPSWEEDOUT is used.

If one subquery is nested within another and both are merged into a semijoin of an outer query, any specification of semijoin strategies for the innermost query are ignored. SEMIJOIN and NO\_SEMIJOIN hints can still be used to enable or disable semijoin transformations for such nested subqueries.

If DUPSWEEDOUT is disabled, on occasion the optimizer may generate a query plan that is far from optimal. This occurs due to heuristic pruning during greedy search, which can be avoided by setting optimizer\_prune\_level=0.

#### Examples:

```
SELECT /*+ NO_SEMIJOIN(@subq1 FIRSTMATCH, LOOSESCAN) */ * FROM t2
WHERE t2.a IN (SELECT /*+ QB_NAME(subq1) */ a FROM t3);
SELECT /*+ SEMIJOIN(@subq1 MATERIALIZATION, DUPSWEEDOUT) */ * FROM t2
WHERE t2.a IN (SELECT /*+ QB_NAME(subq1) */ a FROM t3);
```

Syntax of hints that affect whether to use subquery materialization or IN-to-EXISTS transformations:

```
SUBQUERY([@query_block_name] strategy)
```

The hint name is always SUBQUERY.

For SUBQUERY hints, these strategy values are permitted: INTOEXISTS, MATERIALIZATION.

#### Examples:

```
SELECT id, a IN (SELECT /*+ SUBQUERY(MATERIALIZATION) */ a FROM t1) FROM t2;
SELECT * FROM t2 WHERE t2.a IN (SELECT /*+ SUBQUERY(INTOEXISTS) */ a FROM t1);
```

For semijoin and SUBQUERY hints, a leading <code>@query\_block\_name</code> specifies the query block to which the hint applies. If the hint includes no leading <code>@query\_block\_name</code>, the hint applies to the query block in which it occurs. To assign a name to a query block, see Optimizer Hints for Naming Query Blocks.

If a hint comment contains multiple subquery hints, the first is used. If there are other following hints of that type, they produce a warning. Following hints of other types are silently ignored.

## **Statement Execution Time Optimizer Hints**

The MAX\_EXECUTION\_TIME hint is permitted only for SELECT statements. It places a limit N (a timeout value in milliseconds) on how long a statement is permitted to execute before the server terminates it:

```
	exttt{MAX\_EXECUTION\_TIME}\left(N
ight)
```

Example with a timeout of 1 second (1000 milliseconds):

```
SELECT /*+ MAX_EXECUTION_TIME(1000) */ * FROM t1 INNER JOIN t2 WHERE ...
```

The MAX\_EXECUTION\_TIME (N) hint sets a statement execution timeout of N milliseconds. If this option is absent or N is 0, the statement timeout established by the max\_execution\_time system variable applies.

The MAX\_EXECUTION\_TIME hint is applicable as follows:

- For statements with multiple SELECT keywords, such as unions or statements with subqueries, MAX\_EXECUTION\_TIME applies to the entire statement and must appear after the first SELECT.
- It applies to read-only SELECT statements. Statements that are not read only are those that invoke a stored function that modifies data as a side effect.
- It does not apply to SELECT statements in stored programs and is ignored.

## Variable-Setting Hint Syntax

The SET\_VAR hint sets the session value of a system variable temporarily (for the duration of a single statement). Examples:

```
SELECT /*+ SET_VAR(sort_buffer_size = 16M) */ name FROM people ORDER BY name;
INSERT /*+ SET_VAR(foreign_key_checks=OFF) */ INTO t2 VALUES(2);
SELECT /*+ SET_VAR(optimizer_switch = 'mrr_cost_based=off') */ 1;
```

Syntax of the SET\_VAR hint:

```
SET_VAR(var_name = value)
```

var\_name names a system variable that has a session value (although not all such variables can be named, as explained later). value is the value to assign to the variable; the value must be a scalar.

SET\_VAR makes a temporary variable change, as demonstrated by these statements:

With SET\_VAR, there is no need to save and restore the variable value. This enables you to replace multiple statements by a single statement. Consider this sequence of statements:

```
SET @saved_val = @@SESSION.var_name;
SET @@SESSION.var_name = value;
SELECT ...
SET @@SESSION.var_name = @saved_val;
```

The sequence can be replaced by this single statement:

```
SELECT /*+ SET_VAR(var_name = value) ...
```

Standalone SET statements permit any of these syntaxes for naming session variables:

```
SET SESSION var_name = value;
SET @@SESSION.var_name = value;
SET @@.var_name = value;
```

Because the SET\_VAR hint applies only to session variables, session scope is implicit, and SESSION, @@SESSION., and @@ are neither needed nor permitted. Including explicit session-indicator syntax results in the SET\_VAR hint being ignored with a warning.

Not all session variables are permitted for use with SET\_VAR. Individual system variable descriptions indicate whether each variable is hintable; see Section 7.1.8, "Server System Variables". You can also check a system variable at runtime by attempting to use it with SET\_VAR. If the variable is not hintable, a warning occurs:

```
mysql> SELECT /*+ SET_VAR(collation_server = 'utf8mb4') */ 1;
+---+
| 1 |
+---+
| 1 |
+---+
1 row in set, 1 warning (0.00 sec)

mysql> SHOW WARNINGS\G
********************************
Level: Warning
    Code: 4537
Message: Variable 'collation_server' cannot be set using SET_VAR hint.
```

SET\_VAR syntax permits setting only a single variable, but multiple hints can be given to set multiple variables:

```
SELECT /*+ SET_VAR(optimizer_switch = 'mrr_cost_based=off')
```

```
SET_VAR(max_heap_table_size = 1G) */ 1;
```

If several hints with the same variable name appear in the same statement, the first one is applied and the others are ignored with a warning:

In this case, the second hint is ignored with a warning that it is conflicting.

A SET\_VAR hint is ignored with a warning if no system variable has the specified name or the variable value is incorrect:

```
SELECT /*+ SET_VAR(max_size = 1G) */ 1;
SELECT /*+ SET_VAR(optimizer_switch = 'mrr_cost_based=yes') */ 1;
```

For the first statement, there is no max\_size variable. For the second statement, mrr\_cost\_based takes values of on or off, so attempting to set it to yes is incorrect. In each case, the hint is ignored with a warning.

The SET\_VAR hint is permitted only at the statement level. If used in a subquery, the hint is ignored with a warning.

Replicas ignore SET VAR hints in replicated statements to avoid the potential for security issues.

## **Resource Group Hint Syntax**

The RESOURCE\_GROUP optimizer hint is used for resource group management (see Section 7.1.16, "Resource Groups"). This hint assigns the thread that executes a statement to the named resource group temporarily (for the duration of the statement). It requires the RESOURCE\_GROUP\_ADMIN or RESOURCE\_GROUP\_USER privilege.

### Examples:

```
SELECT /*+ RESOURCE_GROUP(USR_default) */ name FROM people ORDER BY name;
INSERT /*+ RESOURCE_GROUP(Batch) */ INTO t2 VALUES(2);
```

Syntax of the RESOURCE\_GROUP hint:

```
RESOURCE_GROUP(group_name)
```

group\_name indicates the resource group to which the thread should be assigned for the duration of statement execution. If the group is nonexistent, a warning occurs and the hint is ignored.

The RESOURCE\_GROUP hint must appear after the initial statement keyword (SELECT, INSERT, REPLACE, UPDATE, or DELETE).

An alternative to RESOURCE\_GROUP is the SET RESOURCE GROUP statement, which nontemporarily assigns threads to a resource group. See Section 15.7.2.4, "SET RESOURCE GROUP Statement".

### **Optimizer Hints for Naming Query Blocks**

Table-level, index-level, and subquery optimizer hints permit specific query blocks to be named as part of their argument syntax. To create these names, use the <code>QB\_NAME</code> hint, which assigns a name to the query block in which it occurs:

```
QB_NAME(name)
```

QB\_NAME hints can be used to make explicit in a clear way which query blocks other hints apply to. They also permit all non-query block name hints to be specified within a single hint comment for easier understanding of complex statements. Consider the following statement:

```
SELECT ...
```

```
FROM (SELECT ... FROM ...)) ...
```

QB\_NAME hints assign names to query blocks in the statement:

```
SELECT /*+ QB_NAME(qb1) */ ...
FROM (SELECT /*+ QB_NAME(qb2) */ ...
FROM (SELECT /*+ QB_NAME(qb3) */ ... FROM ...)) ...
```

Then other hints can use those names to refer to the appropriate query blocks:

```
SELECT /*+ QB_NAME(qb1) MRR(@qb1 t1) BKA(@qb2) NO_MRR(@qb3t1 idx1, id2) */ ...
FROM (SELECT /*+ QB_NAME(qb2) */ ...
FROM (SELECT /*+ QB_NAME(qb3) */ ... FROM ...)) ...
```

The resulting effect is as follows:

- MRR(@qb1 t1) applies to table t1 in query block qb1.
- BKA(@qb2) applies to query block qb2.
- NO\_MRR(@qb3 t1 idx1, id2) applies to indexes idx1 and idx2 in table t1 in query block qb3.

Query block names are identifiers and follow the usual rules about what names are valid and how to quote them (see Section 11.2, "Schema Object Names"). For example, a query block name that contains spaces must be quoted, which can be done using backticks:

```
SELECT /*+ BKA(@`my hint name`) */ ...
FROM (SELECT /*+ QB_NAME(`my hint name`) */ ...) ...
```

If the ANSI\_QUOTES SQL mode is enabled, it is also possible to quote query block names within double quotation marks:

```
SELECT /*+ BKA(@"my hint name") */ ...
FROM (SELECT /*+ QB_NAME("my hint name") */ ...) ...
```

### 10.9.4 Index Hints

Index hints give the optimizer information about how to choose indexes during query processing. Index hints, described here, differ from optimizer hints, described in Section 10.9.3, "Optimizer Hints". Index and optimizer hints may be used separately or together.

Index hints apply to SELECT and UPDATE statements. They also work with multi-table DELETE statements, but not with single-table DELETE, as shown later in this section.

Index hints are specified following a table name. (For the general syntax for specifying tables in a SELECT statement, see Section 15.2.13.2, "JOIN Clause".) The syntax for referring to an individual table, including index hints, looks like this:

```
tbl_name [[AS] alias] [index_hint_list]
index_hint_list:
    index_hint [index_hint] ...

index_hint:
    USE {INDEX|KEY}
        [FOR {JOIN|ORDER BY|GROUP BY}] ([index_list])
| {IGNORE|FORCE} {INDEX|KEY}
        [FOR {JOIN|ORDER BY|GROUP BY}] (index_list)

index_list:
    index_name [, index_name] ...
```

The USE INDEX (index\_list) hint tells MySQL to use only one of the named indexes to find rows in the table. The alternative syntax IGNORE INDEX (index\_list) tells MySQL to not use some

particular index or indexes. These hints are useful if EXPLAIN shows that MySQL is using the wrong index from the list of possible indexes.

The FORCE INDEX hint acts like USE INDEX (index\_list), with the addition that a table scan is assumed to be *very* expensive. In other words, a table scan is used only if there is no way to use one of the named indexes to find rows in the table.



#### Note

MySQL 8.4 supports the index-level optimizer hints JOIN\_INDEX, GROUP\_INDEX, ORDER\_INDEX, and INDEX, which are equivalent to and intended to supersede FORCE INDEX index hints, as well as the NO\_JOIN\_INDEX, NO\_GROUP\_INDEX, NO\_ORDER\_INDEX, and NO\_INDEX optimizer hints, which are equivalent to and intended to supersede IGNORE INDEX index hints. Thus, you should expect USE INDEX, FORCE INDEX, and IGNORE INDEX to be deprecated in a future release of MySQL, and at some time thereafter to be removed altogether.

These index-level optimizer hints are supported with both single-table and multitable <code>DELETE</code> statements.

For more information, see Index-Level Optimizer Hints.

Each hint requires index names, not column names. To refer to a primary key, use the name PRIMARY. To see the index names for a table, use the SHOW INDEX statement or the Information Schema STATISTICS table.

An *index\_name* value need not be a full index name. It can be an unambiguous prefix of an index name. If a prefix is ambiguous, an error occurs.

#### Examples:

```
SELECT * FROM table1 USE INDEX (col1_index,col2_index)
WHERE col1=1 AND col2=2 AND col3=3;

SELECT * FROM table1 IGNORE INDEX (col3_index)
WHERE col1=1 AND col2=2 AND col3=3;
```

The syntax for index hints has the following characteristics:

- It is syntactically valid to omit index\_list for USE INDEX, which means "use no indexes." Omitting index\_list for FORCE INDEX or IGNORE INDEX is a syntax error.
- You can specify the scope of an index hint by adding a FOR clause to the hint. This provides more
  fine-grained control over optimizer selection of an execution plan for various phases of query
  processing. To affect only the indexes used when MySQL decides how to find rows in the table and
  how to process joins, use FOR JOIN. To influence index usage for sorting or grouping rows, use FOR
  ORDER BY OF FOR GROUP BY.
- · You can specify multiple index hints:

```
SELECT * FROM t1 USE INDEX (i1) IGNORE INDEX FOR ORDER BY (i2) ORDER BY a;
```

It is not an error to name the same index in several hints (even within the same hint):

```
SELECT * FROM t1 USE INDEX (i1) USE INDEX (i1,i1);

However, it is an error to mix USE INDEX and FORCE INDEX for the same table:
```

```
SELECT * FROM t1 USE INDEX FOR JOIN (i1) FORCE INDEX FOR JOIN (i2);
```

If an index hint includes no FOR clause, the scope of the hint is to apply to all parts of the statement. For example, this hint:

```
IGNORE INDEX (i1)
```

is equivalent to this combination of hints:

```
IGNORE INDEX FOR JOIN (i1)
IGNORE INDEX FOR ORDER BY (i1)
IGNORE INDEX FOR GROUP BY (i1)
```

When index hints are processed, they are collected in a single list by type (USE, FORCE, IGNORE) and by scope (FOR JOIN, FOR ORDER BY, FOR GROUP BY). For example:

```
SELECT * FROM t1
USE INDEX () IGNORE INDEX (i2) USE INDEX (i1) USE INDEX (i2);
```

is equivalent to:

```
SELECT * FROM t1
USE INDEX (i1,i2) IGNORE INDEX (i2);
```

The index hints then are applied for each scope in the following order:

- 1. {USE | FORCE} INDEX is applied if present. (If not, the optimizer-determined set of indexes is used.)
- 2. IGNORE INDEX is applied over the result of the previous step. For example, the following two queries are equivalent:

```
SELECT * FROM tl USE INDEX (i1) IGNORE INDEX (i2) USE INDEX (i2);
SELECT * FROM tl USE INDEX (i1);
```

For FULLTEXT searches, index hints work as follows:

- For natural language mode searches, index hints are silently ignored. For example, IGNORE INDEX(i1) is ignored with no warning and the index is still used.
- For boolean mode searches, index hints with FOR ORDER BY OF FOR GROUP BY are silently ignored. Index hints with FOR JOIN or no FOR modifier are honored. In contrast to how hints apply for non-FULLTEXT searches, the hint is used for all phases of query execution (finding rows and retrieval, grouping, and ordering). This is true even if the hint is given for a non-FULLTEXT index.

For example, the following two queries are equivalent:

```
SELECT * FROM t

USE INDEX (index1)

IGNORE INDEX FOR ORDER BY (index1)

IGNORE INDEX FOR GROUP BY (index1)

WHERE ... IN BOOLEAN MODE ...;

SELECT * FROM t

USE INDEX (index1)

WHERE ... IN BOOLEAN MODE ...;
```

Index hints work with DELETE statements, but only if you use multi-table DELETE syntax, as shown here:

```
mysql> EXPLAIN DELETE FROM t1 USE INDEX(col2)
    -> WHERE col1 BETWEEN 1 AND 100 AND COL2 BETWEEN 1 AND 100\G
ERROR 1064 (42000): You have an error in your SQL syntax; check the manual that
corresponds to your MySQL server version for the right syntax to use near 'use
index(col2) where col1 between 1 and 100 and col2 between 1 and 100' at line 1

mysql> EXPLAIN DELETE t1.* FROM t1 USE INDEX(col2)
    -> WHERE col1 BETWEEN 1 AND 100 AND COL2 BETWEEN 1 AND 100\G
*********************
    id: 1
    select_type: DELETE
        table: t1
```

```
partitions: NULL
    type: range
possible_keys: col2
    key: col2
    key_len: 5
       ref: NULL
    rows: 72
    filtered: 11.11
       Extra: Using where
1 row in set, 1 warning (0.00 sec)
```

## 10.9.5 The Optimizer Cost Model

To generate execution plans, the optimizer uses a cost model that is based on estimates of the cost of various operations that occur during query execution. The optimizer has a set of compiled-in default "cost constants" available to it to make decisions regarding execution plans.

The optimizer also has a database of cost estimates to use during execution plan construction. These estimates are stored in the server\_cost and engine\_cost tables in the mysql system database and are configurable at any time. The intent of these tables is to make it possible to easily adjust the cost estimates that the optimizer uses when it attempts to arrive at query execution plans.

- · Cost Model General Operation
- · The Cost Model Database
- · Making Changes to the Cost Model Database

## **Cost Model General Operation**

The configurable optimizer cost model works like this:

- The server reads the cost model tables into memory at startup and uses the in-memory values
  at runtime. Any non-NULL cost estimate specified in the tables takes precedence over the
  corresponding compiled-in default cost constant. Any NULL estimate indicates to the optimizer to use
  the compiled-in default.
- At runtime, the server may re-read the cost tables. This occurs when a storage engine is dynamically loaded or when a FLUSH OPTIMIZER COSTS statement is executed.
- Cost tables enable server administrators to easily adjust cost estimates by changing entries in the
  tables. It is also easy to revert to a default by setting an entry's cost to NULL. The optimizer uses the
  in-memory cost values, so changes to the tables should be followed by FLUSH OPTIMIZER\_COSTS
  to take effect.
- The in-memory cost estimates that are current when a client session begins apply throughout that session until it ends. In particular, if the server re-reads the cost tables, any changed estimates apply only to subsequently started sessions. Existing sessions are unaffected.
- Cost tables are specific to a given server instance. The server does not replicate cost table changes to replicas.

#### The Cost Model Database

The optimizer cost model database consists of two tables in the mysql system database that contain cost estimate information for operations that occur during query execution:

- server\_cost: Optimizer cost estimates for general server operations
- engine\_cost: Optimizer cost estimates for operations specific to particular storage engines

The server\_cost table contains these columns:

• cost\_name

The name of a cost estimate used in the cost model. The name is not case-sensitive. If the server does not recognize the cost name when it reads this table, it writes a warning to the error log.

#### • cost\_value

The cost estimate value. If the value is non-NULL, the server uses it as the cost. Otherwise, it uses the default estimate (the compiled-in value). DBAs can change a cost estimate by updating this column. If the server finds that the cost value is invalid (nonpositive) when it reads this table, it writes a warning to the error log.

To override a default cost estimate (for an entry that specifies NULL), set the cost to a non-NULL value. To revert to the default, set the value to NULL. Then execute FLUSH OPTIMIZER\_COSTS to tell the server to re-read the cost tables.

#### • last\_update

The time of the last row update.

#### • comment

A descriptive comment associated with the cost estimate. DBAs can use this column to provide information about why a cost estimate row stores a particular value.

#### • default\_value

The default (compiled-in) value for the cost estimate. This column is a read-only generated column that retains its value even if the associated cost estimate is changed. For rows added to the table at runtime, the value of this column is NULL.

The primary key for the server\_cost table is the cost\_name column, so it is not possible to create multiple entries for any cost estimate.

The server recognizes these cost\_name values for the server\_cost table:

• disk\_temptable\_create\_cost, disk\_temptable\_row\_cost

The cost estimates for internally created temporary tables stored in a disk-based storage engine (either InnoDB or MyISAM). Increasing these values increases the cost estimate of using internal temporary tables and makes the optimizer prefer query plans with less use of them. For information about such tables, see Section 10.4.4, "Internal Temporary Table Use in MySQL".

The larger default values for these disk parameters compared to the default values for the corresponding memory parameters (memory\_temptable\_create\_cost, memory\_temptable\_row\_cost) reflects the greater cost of processing disk-based tables.

#### • key compare cost

The cost of comparing record keys. Increasing this value causes a query plan that compares many keys to become more expensive. For example, a query plan that performs a filesort becomes relatively more expensive compared to a query plan that avoids sorting by using an index.

• memory\_temptable\_create\_cost, memory\_temptable\_row\_cost

The cost estimates for internally created temporary tables stored in the MEMORY storage engine. Increasing these values increases the cost estimate of using internal temporary tables and makes the optimizer prefer query plans with less use of them. For information about such tables, see Section 10.4.4, "Internal Temporary Table Use in MySQL".

The smaller default values for these memory parameters compared to the default values for the corresponding disk parameters (disk\_temptable\_create\_cost, disk\_temptable\_row\_cost) reflects the lesser cost of processing memory-based tables.

#### • row\_evaluate\_cost

The cost of evaluating record conditions. Increasing this value causes a query plan that examines many rows to become more expensive compared to a query plan that examines fewer rows. For example, a table scan becomes relatively more expensive compared to a range scan that reads fewer rows.

The engine cost table contains these columns:

• engine name

The name of the storage engine to which this cost estimate applies. The name is not case-sensitive. If the value is default, it applies to all storage engines that have no named entry of their own. If the server does not recognize the engine name when it reads this table, it writes a warning to the error log.

• device\_type

The device type to which this cost estimate applies. The column is intended for specifying different cost estimates for different storage device types, such as hard disk drives versus solid state drives. Currently, this information is not used and 0 is the only permitted value.

• cost\_name

Same as in the server\_cost table.

cost\_value

Same as in the server cost table.

• last\_update

Same as in the server\_cost table.

• comment

Same as in the server\_cost table.

• default value

The default (compiled-in) value for the cost estimate. This column is a read-only generated column that retains its value even if the associated cost estimate is changed. For rows added to the table at runtime, the value of this column is NULL, with the exception that if the row has the same cost\_name value as one of the original rows, the default\_value column has the same value as that row.

The primary key for the engine\_cost table is a tuple comprising the (cost\_name, engine\_name, device\_type) columns, so it is not possible to create multiple entries for any combination of values in those columns.

The server recognizes these <code>cost\_name</code> values for the <code>engine\_cost</code> table:

• io\_block\_read\_cost

The cost of reading an index or data block from disk. Increasing this value causes a query plan that reads many disk blocks to become more expensive compared to a query plan that reads fewer disk blocks. For example, a table scan becomes relatively more expensive compared to a range scan that reads fewer blocks.

• memory\_block\_read\_cost

Similar to io\_block\_read\_cost, but represents the cost of reading an index or data block from an in-memory database buffer.

If the io\_block\_read\_cost and memory\_block\_read\_cost values differ, the execution plan may change between two runs of the same query. Suppose that the cost for memory access is less than the cost for disk access. In that case, at server startup before data has been read into the buffer pool, you may get a different plan than after the query has been run because then the data is in memory.

## **Making Changes to the Cost Model Database**

For DBAs who wish to change the cost model parameters from their defaults, try doubling or halving the value and measuring the effect.

Changes to the io\_block\_read\_cost and memory\_block\_read\_cost parameters are most likely to yield worthwhile results. These parameter values enable cost models for data access methods to take into account the costs of reading information from different sources; that is, the cost of reading information from disk versus reading information already in a memory buffer. For example, all other things being equal, setting io\_block\_read\_cost to a value larger than memory\_block\_read\_cost causes the optimizer to prefer query plans that read information already held in memory to plans that must read from disk.

This example shows how to change the default value for io\_block\_read\_cost:

```
UPDATE mysql.engine_cost
   SET cost_value = 2.0
   WHERE cost_name = 'io_block_read_cost';
FLUSH OPTIMIZER_COSTS;
```

This example shows how to change the value of io\_block\_read\_cost only for the InnoDB storage engine:

```
INSERT INTO mysql.engine_cost
  VALUES ('InnoDB', 0, 'io_block_read_cost', 3.0,
  CURRENT_TIMESTAMP, 'Using a slower disk for InnoDB');
FLUSH OPTIMIZER_COSTS;
```

## 10.9.6 Optimizer Statistics

The column\_statistics data dictionary table stores histogram statistics about column values, for use by the optimizer in constructing query execution plans. To perform histogram management, use the ANALYZE TABLE statement.

The column statistics table has these characteristics:

- The table contains statistics for columns of all data types except geometry types (spatial data) and JSON.
- The table is persistent so that column statistics need not be created each time the server starts.
- The server performs updates to the table; users do not.

The column\_statistics table is not directly accessible by users because it is part of the data dictionary. Histogram information is available using INFORMATION\_SCHEMA.COLUMN\_STATISTICS, which is implemented as a view on the data dictionary table. COLUMN\_STATISTICS has these columns:

- SCHEMA\_NAME, TABLE\_NAME, COLUMN\_NAME: The names of the schema, table, and column for which the statistics apply.
- HISTOGRAM: A JSON value describing the column statistics, stored as a histogram.

Column histograms contain buckets for parts of the range of values stored in the column. Histograms are JSON objects to permit flexibility in the representation of column statistics. Here is a sample histogram object:

```
{
    "buckets": [
```

Histogram objects have these keys:

• buckets: The histogram buckets. Bucket structure depends on the histogram type.

For singleton histograms, buckets contain two values:

- Value 1: The value for the bucket. The type depends on the column data type.
- Value 2: A double representing the cumulative frequency for the value. For example, .25 and .75 indicate that 25% and 75% of the values in the column are less than or equal to the bucket value.

For equi-height histograms, buckets contain four values:

- Values 1, 2: The lower and upper inclusive values for the bucket. The type depends on the column data type.
- Value 3: A double representing the cumulative frequency for the value. For example, .25 and .75 indicate that 25% and 75% of the values in the column are less than or equal to the bucket upper value.
- Value 4: The number of distinct values in the range from the bucket lower value to its upper value.
- null-values: A number between 0.0 and 1.0 indicating the fraction of column values that are SQL NULL values. If 0, the column contains no NULL values.
- last-updated: When the histogram was generated, as a UTC value in YYYY-MM-DD hh:mm:ss.uuuuuu format.
- sampling-rate: A number between 0.0 and 1.0 indicating the fraction of data that was sampled to create the histogram. A value of 1 means that all of the data was read (no sampling).
- histogram-type: The histogram type:
  - singleton: One bucket represents one single value in the column. This histogram type is created when the number of distinct values in the column is less than or equal to the number of buckets specified in the ANALYZE TABLE statement that generated the histogram.
  - equi-height: One bucket represents a range of values. This histogram type is created when the number of distinct values in the column is greater than the number of buckets specified in the ANALYZE TABLE statement that generated the histogram.
- number-of-buckets-specified: The number of buckets specified in the ANALYZE TABLE statement that generated the histogram.

- data-type: The type of data this histogram contains. This is needed when reading and parsing
  histograms from persistent storage into memory. The value is one of int, uint (unsigned integer),
  double, decimal, datetime, or string (includes character and binary strings).
- collation-id: The collation ID for the histogram data. It is mostly meaningful when the datatype value is string. Values correspond to ID column values in the Information Schema COLLATIONS table.

To extract particular values from the histogram objects, you can use JSON operations. For example:

The optimizer uses histogram statistics, if applicable, for columns of any data type for which statistics are collected. The optimizer applies histogram statistics to determine row estimates based on the selectivity (filtering effect) of column value comparisons against constant values. Predicates of these forms qualify for histogram use:

```
col_name = constant
col_name <> constant
col_name != constant
col_name > constant
col_name << constant
col_name >= constant
col_name >= constant
col_name is null
col_name is null
col_name is not null
col_name Between constant AND constant
col_name NOT Between constant AND constant
col_name in (constant[, constant] ...)
col_name NOT in (constant[, constant] ...)
```

For example, these statements contain predicates that qualify for histogram use:

```
SELECT * FROM orders WHERE amount BETWEEN 100.0 AND 300.0;
SELECT * FROM tbl WHERE col1 = 15 AND col2 > 100;
```

The requirement for comparison against a constant value includes functions that are constant, such as ABS() and FLOOR():

```
SELECT * FROM tbl WHERE col1 < ABS(-34);
```

Histogram statistics are useful primarily for nonindexed columns. Adding an index to a column for which histogram statistics are applicable might also help the optimizer make row estimates. The tradeoffs are:

- · An index must be updated when table data is modified.
- A histogram is created or updated only on demand, so it adds no overhead when table data is modified. On the other hand, the statistics become progressively more out of date when table modifications occur, until the next time they are updated.

The optimizer prefers range optimizer row estimates to those obtained from histogram statistics. If the optimizer determines that the range optimizer applies, it does not use histogram statistics.

For columns that are indexed, row estimates can be obtained for equality comparisons using index dives (see Section 10.2.1.2, "Range Optimization"). In this case, histogram statistics are not necessarily useful because index dives can yield better estimates.

In some cases, use of histogram statistics may not improve query execution (for example, if the statistics are out of date). To check whether this is the case, use ANALYZE TABLE to regenerate the histogram statistics, then run the query again.

Alternatively, to disable histogram statistics, use ANALYZE TABLE to drop them. A different method of disabling histogram statistics is to turn off the condition\_fanout\_filter flag of the optimizer\_switch system variable (although this may disable other optimizations as well):

```
SET optimizer_switch='condition_fanout_filter=off';
```

If histogram statistics are used, the resulting effect is visible using EXPLAIN. Consider the following query, where no index is available for column col1:

```
SELECT * FROM t1 WHERE col1 < 24;
```

If histogram statistics indicate that 57% of the rows in t1 satisfy the col1 < 24 predicate, filtering can occur even in the absence of an index, and EXPLAIN shows 57.00 in the filtered column.

# 10.10 Buffering and Caching

MySQL uses several strategies that cache information in memory buffers to increase performance.

## 10.10.1 InnoDB Buffer Pool Optimization

InnoDB maintains a storage area called the buffer pool for caching data and indexes in memory. Knowing how the InnoDB buffer pool works, and taking advantage of it to keep frequently accessed data in memory, is an important aspect of MySQL tuning.

For an explanation of the inner workings of the InnoDB buffer pool, an overview of its LRU replacement algorithm, and general configuration information, see Section 17.5.1, "Buffer Pool".

For additional Innobe buffer pool configuration and tuning information, see these sections:

- Section 17.8.3.4, "Configuring InnoDB Buffer Pool Prefetching (Read-Ahead)"
- Section 17.8.3.5, "Configuring Buffer Pool Flushing"
- Section 17.8.3.3, "Making the Buffer Pool Scan Resistant"
- Section 17.8.3.2, "Configuring Multiple Buffer Pool Instances"
- Section 17.8.3.6, "Saving and Restoring the Buffer Pool State"
- Section 17.8.3.1, "Configuring InnoDB Buffer Pool Size"

# 10.10.2 The MyISAM Key Cache

To minimize disk I/O, the MyISAM storage engine exploits a strategy that is used by many database management systems. It employs a cache mechanism to keep the most frequently accessed table blocks in memory:

- For index blocks, a special structure called the *key cache* (or *key buffer*) is maintained. The structure contains a number of block buffers where the most-used index blocks are placed.
- For data blocks, MySQL uses no special cache. Instead it relies on the native operating system file system cache.

This section first describes the basic operation of the MyISAM key cache. Then it discusses features that improve key cache performance and that enable you to better control cache operation:

Multiple sessions can access the cache concurrently.

You can set up multiple key caches and assign table indexes to specific caches.

To control the size of the key cache, use the key\_buffer\_size system variable. If this variable is set equal to zero, no key cache is used. The key cache also is not used if the key\_buffer\_size value is too small to allocate the minimal number of block buffers (8).

When the key cache is not operational, index files are accessed using only the native file system buffering provided by the operating system. (In other words, table index blocks are accessed using the same strategy as that employed for table data blocks.)

An index block is a contiguous unit of access to the MyISAM index files. Usually the size of an index block is equal to the size of nodes of the index B-tree. (Indexes are represented on disk using a B-tree data structure. Nodes at the bottom of the tree are leaf nodes. Nodes above the leaf nodes are nonleaf nodes.)

All block buffers in a key cache structure are the same size. This size can be equal to, greater than, or less than the size of a table index block. Usually one these two values is a multiple of the other.

When data from any table index block must be accessed, the server first checks whether it is available in some block buffer of the key cache. If it is, the server accesses data in the key cache rather than on disk. That is, it reads from the cache or writes into it rather than reading from or writing to disk. Otherwise, the server chooses a cache block buffer containing a different table index block (or blocks) and replaces the data there by a copy of required table index block. As soon as the new index block is in the cache, the index data can be accessed.

If it happens that a block selected for replacement has been modified, the block is considered "dirty." In this case, prior to being replaced, its contents are flushed to the table index from which it came.

Usually the server follows an *LRU* (*Least Recently Used*) strategy: When choosing a block for replacement, it selects the least recently used index block. To make this choice easier, the key cache module maintains all used blocks in a special list (*LRU chain*) ordered by time of use. When a block is accessed, it is the most recently used and is placed at the end of the list. When blocks need to be replaced, blocks at the beginning of the list are the least recently used and become the first candidates for eviction.

The InnoDB storage engine also uses an LRU algorithm, to manage its buffer pool. See Section 17.5.1, "Buffer Pool".

## 10.10.2.1 Shared Key Cache Access

Threads can access key cache buffers simultaneously, subject to the following conditions:

- A buffer that is not being updated can be accessed by multiple sessions.
- A buffer that is being updated causes sessions that need to use it to wait until the update is complete.
- Multiple sessions can initiate requests that result in cache block replacements, as long as they do not
  interfere with each other (that is, as long as they need different index blocks, and thus cause different
  cache blocks to be replaced).

Shared access to the key cache enables the server to improve throughput significantly.

## 10.10.2.2 Multiple Key Caches



#### Note

As of MySQL 8.4, the compound-part structured-variable syntax discussed here for referring to multiple MyISAM key caches is deprecated.

Shared access to the key cache improves performance but does not eliminate contention among sessions entirely. They still compete for control structures that manage access to the key cache

buffers. To reduce key cache access contention further, MySQL also provides multiple key caches. This feature enables you to assign different table indexes to different key caches.

Where there are multiple key caches, the server must know which cache to use when processing queries for a given MyISAM table. By default, all MyISAM table indexes are cached in the default key cache. To assign table indexes to a specific key cache, use the CACHE INDEX statement (see Section 15.7.8.2, "CACHE INDEX Statement"). For example, the following statement assigns indexes from the tables t1, t2, and t3 to the key cache named hot\_cache:

The key cache referred to in a CACHE INDEX statement can be created by setting its size with a SET GLOBAL parameter setting statement or by using server startup options. For example:

```
mysql> SET GLOBAL keycachel.key_buffer_size=128*1024;
```

To destroy a key cache, set its size to zero:

```
mysql> SET GLOBAL keycache1.key_buffer_size=0;
```

You cannot destroy the default key cache. Any attempt to do this is ignored:

```
mysql> SET GLOBAL key_buffer_size = 0;

mysql> SHOW VARIABLES LIKE 'key_buffer_size';

+-----+

| Variable_name | Value |

+-----+

| key_buffer_size | 8384512 |

+------+
```

Key cache variables are structured system variables that have a name and components. For keycachel.key\_buffer\_size, keycachel is the cache variable name and key\_buffer\_size is the cache component. See Section 7.1.9.5, "Structured System Variables", for a description of the syntax used for referring to structured key cache system variables.

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.

For a busy server, you can use a strategy that involves three key caches:

- A "hot" key cache that takes up 20% of the space allocated for all key caches. Use this for tables that are heavily used for searches but that are not updated.
- A "cold" key cache that takes up 20% of the space allocated for all key caches. Use this cache for medium-sized, intensively modified tables, such as temporary tables.
- A "warm" key cache that takes up 60% of the key cache space. Employ this as the default key cache, to be used by default for all other tables.

One reason the use of three key caches is beneficial is that access to one key cache structure does not block access to the others. Statements that access tables assigned to one cache do not compete with statements that access tables assigned to another cache. Performance gains occur for other reasons as well:

 The hot cache is used only for retrieval queries, so its contents are never modified. Consequently, whenever an index block needs to be pulled in from disk, the contents of the cache block chosen for replacement need not be flushed first.

- For an index assigned to the hot cache, if there are no queries requiring an index scan, there is a
  high probability that the index blocks corresponding to nonleaf nodes of the index B-tree remain in
  the cache.
- An update operation most frequently executed for temporary tables is performed much faster when
  the updated node is in the cache and need not be read from disk first. If the size of the indexes of the
  temporary tables are comparable with the size of cold key cache, the probability is very high that the
  updated node is in the cache.

The CACHE INDEX statement sets up an association between a table and a key cache, but the association is lost each time the server restarts. If you want the association to take effect each time the server starts, one way to accomplish this is to use an option file: Include variable settings that configure your key caches, and an init\_file system variable that names a file containing CACHE INDEX statements to be executed. For example:

```
key_buffer_size = 4G
hot_cache.key_buffer_size = 2G
cold_cache.key_buffer_size = 2G
init_file=/path/to/data-directory/mysqld_init.sql
```

The statements in <code>mysqld\_init.sql</code> are executed each time the server starts. The file should contain one SQL statement per line. The following example assigns several tables each to <code>hot\_cache</code> and <code>cold\_cache</code>:

```
CACHE INDEX db1.t1, db1.t2, db2.t3 IN hot_cache
CACHE INDEX db1.t4, db2.t5, db2.t6 IN cold_cache
```

## 10.10.2.3 Midpoint Insertion Strategy

By default, the key cache management system uses a simple LRU strategy for choosing key cache blocks to be evicted, but it also supports a more sophisticated method called the *midpoint insertion strategy*.

When using the midpoint insertion strategy, the LRU chain is divided into two parts: a hot sublist and a warm sublist. The division point between two parts is not fixed, but the key cache management system takes care that the warm part is not "too short," always containing at least key\_cache\_division\_limit percent of the key cache blocks. key\_cache\_division\_limit is a component of structured key cache variables, so its value is a parameter that can be set per cache.

When an index block is read from a table into the key cache, it is placed at the end of the warm sublist. After a certain number of hits (accesses of the block), it is promoted to the hot sublist. At present, the number of hits required to promote a block (3) is the same for all index blocks.

A block promoted into the hot sublist is placed at the end of the list. The block then circulates within this sublist. If the block stays at the beginning of the sublist for a long enough time, it is demoted to the warm sublist. This time is determined by the value of the key\_cache\_age\_threshold component of the key cache.

The threshold value prescribes that, for a key cache containing N blocks, the block at the beginning of the hot sublist not accessed within the last N \* key\_cache\_age\_threshold / 100 hits is to be moved to the beginning of the warm sublist. It then becomes the first candidate for eviction, because blocks for replacement always are taken from the beginning of the warm sublist.

The midpoint insertion strategy enables you to keep more-valued blocks always in the cache. If you prefer to use the plain LRU strategy, leave the key\_cache\_division\_limit value set to its default of 100.

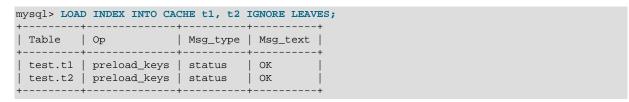
The midpoint insertion strategy helps to improve performance when execution of a query that requires an index scan effectively pushes out of the cache all the index blocks corresponding to valuable high-level B-tree nodes. To avoid this, you must use a midpoint insertion strategy with the key\_cache\_division\_limit set to much less than 100. Then valuable frequently hit nodes are preserved in the hot sublist during an index scan operation as well.

## 10.10.2.4 Index Preloading

If there are enough blocks in a key cache to hold blocks of an entire index, or at least the blocks corresponding to its nonleaf nodes, it makes sense to preload the key cache with index blocks before starting to use it. Preloading enables you to put the table index blocks into a key cache buffer in the most efficient way: by reading the index blocks from disk sequentially.

Without preloading, the blocks are still placed into the key cache as needed by queries. Although the blocks stay in the cache, because there are enough buffers for all of them, they are fetched from disk in random order, and not sequentially.

To preload an index into a cache, use the LOAD INDEX INTO CACHE statement. For example, the following statement preloads nodes (index blocks) of indexes of the tables t1 and t2:



The IGNORE LEAVES modifier causes only blocks for the nonleaf nodes of the index to be preloaded. Thus, the statement shown preloads all index blocks from t1, but only blocks for the nonleaf nodes from t2.

If an index has been assigned to a key cache using a CACHE INDEX statement, preloading places index blocks into that cache. Otherwise, the index is loaded into the default key cache.

## 10.10.2.5 Key Cache Block Size

It is possible to specify the size of the block buffers for an individual key cache using the key\_cache\_block\_size variable. This permits tuning of the performance of I/O operations for index files.

The best performance for I/O operations is achieved when the size of read buffers is equal to the size of the native operating system I/O buffers. But setting the size of key nodes equal to the size of the I/O buffer does not always ensure the best overall performance. When reading the big leaf nodes, the server pulls in a lot of unnecessary data, effectively preventing reading other leaf nodes.

To control the size of blocks in the .MYI index file of MyISAM tables, use the --myisam-block-size option at server startup.

### 10.10.2.6 Restructuring a Key Cache

A key cache can be restructured at any time by updating its parameter values. For example:

```
mysql> SET GLOBAL cold_cache.key_buffer_size=4*1024*1024;
```

If you assign to either the key\_buffer\_size or key\_cache\_block\_size key cache component a value that differs from the component's current value, the server destroys the cache's old structure and creates a new one based on the new values. If the cache contains any dirty blocks, the server saves them to disk before destroying and re-creating the cache. Restructuring does not occur if you change other key cache parameters.

When restructuring a key cache, the server first flushes the contents of any dirty buffers to disk. After that, the cache contents become unavailable. However, restructuring does not block queries that need to use indexes assigned to the cache. Instead, the server directly accesses the table indexes using native file system caching. File system caching is not as efficient as using a key cache, so although queries execute, a slowdown can be anticipated. After the cache has been restructured, it becomes available again for caching indexes assigned to it, and the use of file system caching for the indexes ceases.

## 10.10.3 Caching of Prepared Statements and Stored Programs

For certain statements that a client might execute multiple times during a session, the server converts the statement to an internal structure and caches that structure to be used during execution. Caching enables the server to perform more efficiently because it avoids the overhead of reconverting the statement should it be needed again during the session. Conversion and caching occurs for these statements:

- Prepared statements, both those processed at the SQL level (using the PREPARE statement) and
  those processed using the binary client/server protocol (using the mysql\_stmt\_prepare() C
  API function). The max\_prepared\_stmt\_count system variable controls the total number of
  statements the server caches. (The sum of the number of prepared statements across all sessions.)
- Stored programs (stored procedures and functions, triggers, and events). In this case, the server converts and caches the entire program body. The <a href="mailto:stored\_program\_cache">stored\_program\_cache</a> system variable indicates the approximate number of stored programs the server caches per session.

The server maintains caches for prepared statements and stored programs on a per-session basis. Statements cached for one session are not accessible to other sessions. When a session ends, the server discards any statements cached for it.

When the server uses a cached internal statement structure, it must take care that the structure does not go out of date. Metadata changes can occur for an object used by the statement, causing a mismatch between the current object definition and the definition as represented in the internal statement structure. Metadata changes occur for DDL statements such as those that create, drop, alter, rename, or truncate tables, or that analyze, optimize, or repair tables. Table content changes (for example, with INSERT or UPDATE) do not change metadata, nor do SELECT statements.

Here is an illustration of the problem. Suppose that a client prepares this statement:

```
PREPARE s1 FROM 'SELECT * FROM t1';
```

The SELECT \* expands in the internal structure to the list of columns in the table. If the set of columns in the table is modified with ALTER TABLE, the prepared statement goes out of date. If the server does not detect this change the next time the client executes s1, the prepared statement returns incorrect results.

To avoid problems caused by metadata changes to tables or views referred to by the prepared statement, the server detects these changes and automatically reprepares the statement when it is next executed. That is, the server reparses the statement and rebuilds the internal structure. Reparsing also occurs after referenced tables or views are flushed from the table definition cache, either implicitly to make room for new entries in the cache, or explicitly due to FLUSH TABLES.

Similarly, if changes occur to objects used by a stored program, the server reparses affected statements within the program.

The server also detects metadata changes for objects in expressions. These might be used in statements specific to stored programs, such as DECLARE CURSOR or flow-control statements such as IF, CASE, and RETURN.

To avoid reparsing entire stored programs, the server reparses affected statements or expressions within a program only as needed. Examples:

- Suppose that metadata for a table or view is changed. Reparsing occurs for a SELECT \* within the program that accesses the table or view, but not for a SELECT \* that does not access the table or view.
- When a statement is affected, the server reparses it only partially if possible. Consider this CASE statement:

CASE case\_expr

```
WHEN when_expr1 ...
WHEN when_expr2 ...
WHEN when_expr3 ...
...
END CASE
```

If a metadata change affects only WHEN when\_expr3, that expression is reparsed. case\_expr and the other WHEN expressions are not reparsed.

Reparsing uses the default database and SQL mode that were in effect for the original conversion to internal form.

The server attempts reparsing up to three times. An error occurs if all attempts fail.

Reparsing is automatic, but to the extent that it occurs, diminishes prepared statement and stored program performance.

For prepared statements, the Com\_stmt\_reprepare status variable tracks the number of repreparations.

# 10.11 Optimizing Locking Operations

MySQL manages contention for table contents using locking:

- Internal locking is performed within the MySQL server itself to manage contention for table contents by multiple threads. This type of locking is internal because it is performed entirely by the server and involves no other programs. See Section 10.11.1, "Internal Locking Methods".
- External locking occurs when the server and other programs lock MyISAM table files to coordinate
  among themselves which program can access the tables at which time. See Section 10.11.5,
  "External Locking".

# 10.11.1 Internal Locking Methods

This section discusses internal locking; that is, locking performed within the MySQL server itself to manage contention for table contents by multiple sessions. This type of locking is internal because it is performed entirely by the server and involves no other programs. For locking performed on MySQL files by other programs, see Section 10.11.5, "External Locking".

- Row-Level Locking
- Table-Level Locking
- · Choosing the Type of Locking

## **Row-Level Locking**

MySQL uses row-level locking for InnoDB tables to support simultaneous write access by multiple sessions, making them suitable for multi-user, highly concurrent, and OLTP applications.

To avoid deadlocks when performing multiple concurrent write operations on a single InnoDB table, acquire necessary locks at the start of the transaction by issuing a SELECT ... FOR UPDATE statement for each group of rows expected to be modified, even if the data change statements come later in the transaction. If transactions modify or lock more than one table, issue the applicable statements in the same order within each transaction. Deadlocks affect performance rather than representing a serious error, because InnoDB automatically detects deadlock conditions by default and rolls back one of the affected transactions.

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 configuration option.

Advantages of row-level locking:

- · Fewer lock conflicts when different sessions access different rows.
- · Fewer changes for rollbacks.
- Possible to lock a single row for a long time.

## Table-Level Locking

MySQL uses table-level locking for MyISAM, MEMORY, and MERGE tables, permitting only one session to update those tables at a time. This locking level makes these storage engines more suitable for read-only, read-mostly, or single-user applications.

These storage engines avoid deadlocks by always requesting all needed locks at once at the beginning of a query and always locking the tables in the same order. The tradeoff is that this strategy reduces concurrency; other sessions that want to modify the table must wait until the current data change statement finishes.

Advantages of table-level locking:

- Relatively little memory required (row locking requires memory per row or group of rows locked)
- Fast when used on a large part of the table because only a single lock is involved.
- Fast if you often do GROUP BY operations on a large part of the data or must scan the entire table frequently.

MySQL grants table write locks as follows:

- 1. If there are no locks on the table, put a write lock on it.
- 2. Otherwise, put the lock request in the write lock queue.

MySQL grants table read locks as follows:

- 1. If there are no write locks on the table, put a read lock on it.
- 2. Otherwise, put the lock request in the read lock queue.

Table updates are given higher priority than table retrievals. Therefore, when a lock is released, the lock is made available to the requests in the write lock queue and then to the requests in the read lock queue. This ensures that updates to a table are not "starved" even when there is heavy SELECT activity for the table. However, if there are many updates for a table, SELECT statements wait until there are no more updates.

For information on altering the priority of reads and writes, see Section 10.11.2, "Table Locking Issues".

You can analyze the table lock contention on your system by checking the Table\_locks\_immediate and Table\_locks\_waited status variables, which indicate the number of times that requests for table locks could be granted immediately and the number that had to wait, respectively:

mysql> SHOW STATUS LIKE	
Variable_name	Value
	1151552   15324

The Performance Schema lock tables also provide locking information. See Section 29.12.13, "Performance Schema Lock Tables".

The MyISAM storage engine supports concurrent inserts to reduce contention between readers and writers for a given table: If a MyISAM table has no free blocks in the middle of the data file, rows are

always inserted at the end of the data file. In this case, you can freely mix concurrent INSERT and SELECT statements for a MyISAM table without locks. That is, you can insert rows into a MyISAM table at the same time other clients are reading from it. Holes can result from rows having been deleted from or updated in the middle of the table. If there are holes, concurrent inserts are disabled but are enabled again automatically when all holes have been filled with new data. To control this behavior, use the concurrent\_insert system variable. See Section 10.11.3, "Concurrent Inserts".

If you acquire a table lock explicitly with LOCK TABLES, you can request a READ LOCAL lock rather than a READ lock to enable other sessions to perform concurrent inserts while you have the table locked.

To perform many INSERT and SELECT operations on a table t1 when concurrent inserts are not possible, you can insert rows into a temporary table  $temp_t1$  and update the real table with the rows from the temporary table:

```
mysql> LOCK TABLES t1 WRITE, temp_t1 WRITE;
mysql> INSERT INTO t1 SELECT * FROM temp_t1;
mysql> DELETE FROM temp_t1;
mysql> UNLOCK TABLES;
```

## **Choosing the Type of Locking**

Generally, table locks are superior to row-level locks in the following cases:

- Most statements for the table are reads.
- Statements for the table are a mix of reads and writes, where writes are updates or deletes for a single row that can be fetched with one key read:

```
UPDATE tbl_name SET column=value WHERE unique_key_col=key_value;

DELETE FROM tbl_name WHERE unique_key_col=key_value;
```

- SELECT combined with concurrent INSERT statements, and very few UPDATE or DELETE statements.
- Many scans or GROUP BY operations on the entire table without any writers.

With higher-level locks, you can more easily tune applications by supporting locks of different types, because the lock overhead is less than for row-level locks.

Options other than row-level locking:

- Versioning (such as that used in MySQL for concurrent inserts) where it is possible to have one
  writer at the same time as many readers. This means that the database or table supports different
  views for the data depending on when access begins. Other common terms for this are "time travel,"
  "copy on write," or "copy on demand."
- Copy on demand is in many cases superior to row-level locking. However, in the worst case, it can use much more memory than using normal locks.
- Instead of using row-level locks, you can employ application-level locks, such as those provided by GET\_LOCK() and RELEASE\_LOCK() in MySQL. These are advisory locks, so they work only with applications that cooperate with each other. See Section 14.14, "Locking Functions".

# 10.11.2 Table Locking Issues

InnoDB tables use row-level locking so that multiple sessions and applications can read from and write to the same table simultaneously, without making each other wait or producing inconsistent results. For this storage engine, avoid using the LOCK TABLES statement, because it does not offer any extra protection, but instead reduces concurrency. The automatic row-level locking makes these tables suitable for your busiest databases with your most important data, while also simplifying application logic since you do not need to lock and unlock tables. Consequently, the InnoDB storage engine is the default in MySQL.

MySQL uses table locking (instead of page, row, or column locking) for all storage engines except InnoDB. The locking operations themselves do not have much overhead. But because only one session can write to a table at any one time, for best performance with these other storage engines, use them primarily for tables that are queried often and rarely inserted into or updated.

- Performance Considerations Favoring InnoDB
- Workarounds for Locking Performance Issues

## **Performance Considerations Favoring InnoDB**

When choosing whether to create a table using InnoDB or a different storage engine, keep in mind the following disadvantages of table locking:

- Table locking enables many sessions to read from a table at the same time, but if a session wants to
  write to a table, it must first get exclusive access, meaning it might have to wait for other sessions to
  finish with the table first. During the update, all other sessions that want to access this particular table
  must wait until the update is done.
- Table locking causes problems when a session is waiting because the disk is full and free space
  needs to become available before the session can proceed. In this case, all sessions that want to
  access the problem table are also put in a waiting state until more disk space is made available.
- A SELECT statement that takes a long time to run prevents other sessions from updating the table in
  the meantime, making the other sessions appear slow or unresponsive. While a session is waiting to
  get exclusive access to the table for updates, other sessions that issue SELECT statements queue
  up behind it, reducing concurrency even for read-only sessions.

## **Workarounds for Locking Performance Issues**

The following items describe some ways to avoid or reduce contention caused by table locking:

- Consider switching the table to the Innobs storage engine, either using CREATE TABLE ... ENGINE=INNOBS during setup, or using ALTER TABLE ... ENGINE=INNOBS for an existing table. See Chapter 17, *The Innobs Storage Engine* for more details about this storage engine.
- Optimize SELECT statements to run faster so that they lock tables for a shorter time. You might have to create some summary tables to do this.
- Start mysqld with --low-priority-updates. For storage engines that use only table-level locking (such as MyISAM, MEMORY, and MERGE), this gives all statements that update (modify) a table lower priority than SELECT statements. In this case, the second SELECT statement in the preceding scenario would execute before the UPDATE statement, and would not wait for the first SELECT to finish.
- To specify that all updates issued in a specific connection should be done with low priority, set the low\_priority\_updates server system variable equal to 1.
- To give a specific INSERT, UPDATE, or DELETE statement lower priority, use the LOW\_PRIORITY attribute.
- To give a specific SELECT statement higher priority, use the HIGH\_PRIORITY attribute. See Section 15.2.13, "SELECT Statement".
- Start mysqld with a low value for the max\_write\_lock\_count system variable to force MySQL to temporarily elevate the priority of all SELECT statements that are waiting for a table after a specific number of write locks to the table occur (for example, for insert operations). This permits read locks after a certain number of write locks.
- If you have problems with mixed SELECT and DELETE statements, the LIMIT option to DELETE may help. See Section 15.2.2, "DELETE Statement".

- Using SQL\_BUFFER\_RESULT with SELECT statements can help to make the duration of table locks shorter. See Section 15.2.13, "SELECT Statement".
- Splitting table contents into separate tables may help, by allowing queries to run against columns in one table, while updates are confined to columns in a different table.
- You could change the locking code in mysys/thr\_lock.c to use a single queue. In this case, write locks and read locks would have the same priority, which might help some applications.

## 10.11.3 Concurrent Inserts

The MyISAM storage engine supports concurrent inserts to reduce contention between readers and writers for a given table: If a MyISAM table has no holes in the data file (deleted rows in the middle), an INSERT statement can be executed to add rows to the end of the table at the same time that SELECT statements are reading rows from the table. If there are multiple INSERT statements, they are queued and performed in sequence, concurrently with the SELECT statements. The results of a concurrent INSERT may not be visible immediately.

The concurrent\_insert system variable can be set to modify the concurrent-insert processing. By default, the variable is set to AUTO (or 1) and concurrent inserts are handled as just described. If concurrent\_insert is set to NEVER (or 0), concurrent inserts are disabled. If the variable is set to ALWAYS (or 2), concurrent inserts at the end of the table are permitted even for tables that have deleted rows. See also the description of the concurrent insert system variable.

If you are using the binary log, concurrent inserts are converted to normal inserts for CREATE ... SELECT or INSERT ... SELECT statements. This is done to ensure that you can re-create an exact copy of your tables by applying the log during a backup operation. See Section 7.4.4, "The Binary Log". In addition, for those statements a read lock is placed on the selected-from table such that inserts into that table are blocked. The effect is that concurrent inserts for that table must wait as well.

With LOAD DATA, if you specify CONCURRENT with a MyISAM table that satisfies the condition for concurrent inserts (that is, it contains no free blocks in the middle), other sessions can retrieve data from the table while LOAD DATA is executing. Use of the CONCURRENT option affects the performance of LOAD DATA a bit, even if no other session is using the table at the same time.

If you specify <code>HIGH\_PRIORITY</code>, it overrides the effect of the <code>--low-priority-updates</code> option if the server was started with that option. It also causes concurrent inserts not to be used.

For LOCK TABLE, the difference between READ LOCAL and READ is that READ LOCAL permits nonconflicting INSERT statements (concurrent inserts) to execute while the lock is held. However, this cannot be used if you are going to manipulate the database using processes external to the server while you hold the lock.

# 10.11.4 Metadata Locking

MySQL uses metadata locking to manage concurrent access to database objects and to ensure data consistency. Metadata locking applies not just to tables, but also to schemas, stored programs (procedures, functions, triggers, scheduled events), tablespaces, user locks acquired with the GET\_LOCK() function (see Section 14.14, "Locking Functions"), and locks acquired with the locking service described in Section 7.6.9.1, "The Locking Service".

The Performance Schema metadata\_locks table exposes metadata lock information, which can be useful for seeing which sessions hold locks, are blocked waiting for locks, and so forth. For details, see Section 29.12.13.3, "The metadata\_locks Table".

Metadata locking does involve some overhead, which increases as query volume increases. Metadata contention increases the more that multiple queries attempt to access the same objects.

Metadata locking is not a replacement for the table definition cache, and its mutexes and locks differ from the LOCK\_open mutex. The following discussion provides some information about how metadata locking works.

- · Metadata Lock Acquisition
- Metadata Lock Release

## **Metadata Lock Acquisition**

If there are multiple waiters for a given lock, the highest-priority lock request is satisfied first, with an exception related to the max\_write\_lock\_count system variable. Write lock requests have higher priority than read lock requests. However, if max\_write\_lock\_count is set to some low value (say, 10), read lock requests may be preferred over pending write lock requests if the read lock requests have already been passed over in favor of 10 write lock requests. Normally this behavior does not occur because max\_write\_lock\_count by default has a very large value.

Statements acquire metadata locks one by one, not simultaneously, and perform deadlock detection in the process.

DML statements normally acquire locks in the order in which tables are mentioned in the statement.

DDL statements, LOCK TABLES, and other similar statements try to reduce the number of possible deadlocks between concurrent DDL statements by acquiring locks on explicitly named tables in name order. Locks might be acquired in a different order for implicitly used tables (such as tables in foreign key relationships that also must be locked).

For example, RENAME TABLE is a DDL statement that acquires locks in name order:

This RENAME TABLE statement renames tbla to something else, and renames tblc to tbla:

```
RENAME TABLE tbla TO tbld, tblc TO tbla;
```

The statement acquires metadata locks, in order, on tbla, tblc, and tbld (because tbld follows tblc in name order):

• This slightly different statement also renames tbla to something else, and renames tblc to tbla:

```
RENAME TABLE tbla TO tblb, tblc TO tbla;
```

In this case, the statement acquires metadata locks, in order, on tbla, tblb, and tblc (because tblb precedes tblc in name order):

Both statements acquire locks on tbla and tblc, in that order, but differ in whether the lock on the remaining table name is acquired before or after tblc.

Metadata lock acquisition order can make a difference in operation outcome when multiple transactions execute concurrently, as the following example illustrates.

Begin with two tables x and  $x_{new}$  that have identical structure. Three clients issue statements that involve these tables:

#### Client 1:

```
LOCK TABLE x WRITE, x_new WRITE;
```

The statement requests and acquires write locks in name order on x and x\_new.

#### Client 2:

```
INSERT INTO x VALUES(1);
```

The statement requests and blocks waiting for a write lock on x.

#### Client 3:

```
RENAME TABLE x TO x_old, x_new TO x;
```

The statement requests exclusive locks in name order on x,  $x_new$ , and  $x_old$ , but blocks waiting for the lock on x.

#### Client 1:

```
UNLOCK TABLES;
```

The statement releases the write locks on x and  $x_new$ . The exclusive lock request for x by Client 3 has higher priority than the write lock request by Client 2, so Client 3 acquires its lock on x, then also on  $x_new$  and  $x_old$ , performs the renaming, and releases its locks. Client 2 then acquires its lock on x, performs the insert, and releases its lock.

Lock acquisition order results in the RENAME TABLE executing before the INSERT. The x into which the insert occurs is the table that was named  $x_{new}$  when Client 2 issued the insert and was renamed to x by Client 3:

Now begin instead with tables named x and  $new_x$  that have identical structure. Again, three clients issue statements that involve these tables:

#### Client 1:

```
LOCK TABLE x WRITE, new_x WRITE;
```

The statement requests and acquires write locks in name order on new\_x and x.

#### Client 2:

```
INSERT INTO x VALUES(1);
```

The statement requests and blocks waiting for a write lock on x.

#### Client 3:

```
RENAME TABLE x TO old_x, new_x TO x;
```

The statement requests exclusive locks in name order on new\_x, old\_x, and x, but blocks waiting for the lock on new\_x.

#### Client 1:

```
UNLOCK TABLES;
```

The statement releases the write locks on x and  $new_x$ . For x, the only pending request is by Client 2, so Client 2 acquires its lock, performs the insert, and releases the lock. For  $new_x$ , the only pending request is by Client 3, which is permitted to acquire that lock (and also the lock on  $old_x$ ). The rename operation still blocks for the lock on x until the Client 2 insert finishes and releases its lock. Then Client 3 acquires the lock on x, performs the rename, and releases its lock.

In this case, lock acquisition order results in the INSERT executing before the RENAME TABLE. The x into which the insert occurs is the original x, now renamed to old\_x by the rename operation:

```
mysql> SELECT * FROM x;
Empty set (0.01 sec)
mysql> SELECT * FROM old_x;
```

If order of lock acquisition in concurrent statements makes a difference to an application in operation outcome, as in the preceding example, you may be able to adjust the table names to affect the order of lock acquisition.

Metadata locks are extended, as necessary, to tables related by a foreign key constraint to prevent conflicting DML and DDL operations from executing concurrently on the related tables. When updating a parent table, a metadata lock is taken on the child table while updating foreign key metadata. Foreign key metadata is owned by the child table.

## **Metadata Lock Release**

To ensure transaction serializability, the server must not permit one session to perform a data definition language (DDL) statement on a table that is used in an uncompleted explicitly or implicitly started transaction in another session. The server achieves this by acquiring metadata locks on tables used within a transaction and deferring release of those locks until the transaction ends. A metadata lock on a table prevents changes to the table's structure. This locking approach has the implication that a table that is being used by a transaction within one session cannot be used in DDL statements by other sessions until the transaction ends.

This principle applies not only to transactional tables, but also to nontransactional tables. Suppose that a session begins a transaction that uses transactional table t and nontransactional table nt as follows:

```
START TRANSACTION;
SELECT * FROM t;
SELECT * FROM nt;
```

The server holds metadata locks on both t and nt until the transaction ends. If another session attempts a DDL or write lock operation on either table, it blocks until metadata lock release at transaction end. For example, a second session blocks if it attempts any of these operations:

```
DROP TABLE t;
ALTER TABLE t ...;
DROP TABLE nt;
ALTER TABLE nt ...;
LOCK TABLE t ... WRITE;
```

The same behavior applies for The LOCK TABLES ... READ. That is, explicitly or implicitly started transactions that update any table (transactional or nontransactional) block and are blocked by LOCK TABLES ... READ for that table.

If the server acquires metadata locks for a statement that is syntactically valid but fails during execution, it does not release the locks early. Lock release is still deferred to the end of the transaction because the failed statement is written to the binary log and the locks protect log consistency.

In autocommit mode, each statement is in effect a complete transaction, so metadata locks acquired for the statement are held only to the end of the statement.

Metadata locks acquired during a PREPARE statement are released once the statement has been prepared, even if preparation occurs within a multiple-statement transaction.

For XA transactions in PREPARED state, metadata locks are maintained across client disconnects and server restarts, until an XA COMMIT or XA ROLLBACK is executed.

# 10.11.5 External Locking

External locking is the use of file system locking to manage contention for MyISAM database tables by multiple processes. External locking is used in situations where a single process such as the MySQL

server cannot be assumed to be the only process that requires access to tables. Here are some examples:

- If you run multiple servers that use the same database directory (not recommended), each server must have external locking enabled.
- If you use myisamchk to perform table maintenance operations on MyISAM tables, you must either ensure that the server is not running, or that the server has external locking enabled so that it locks table files as necessary to coordinate with myisamchk for access to the tables. The same is true for use of myisampack to pack MyISAM tables.

If the server is run with external locking enabled, you can use myisamchk at any time for read operations such a checking tables. In this case, if the server tries to update a table that myisamchk is using, the server waits for myisamchk to finish before it continues.

If you use myisamchk for write operations such as repairing or optimizing tables, or if you use myisampack to pack tables, you *must* always ensure that the mysqld server is not using the table. If you do not stop mysqld, at least do a mysqladmin flush-tables before you run myisamchk. Your tables may become corrupted if the server and myisamchk access the tables simultaneously.

With external locking in effect, each process that requires access to a table acquires a file system lock for the table files before proceeding to access the table. If all necessary locks cannot be acquired, the process is blocked from accessing the table until the locks can be obtained (after the process that currently holds the locks releases them).

External locking affects server performance because the server must sometimes wait for other processes before it can access tables.

External locking is unnecessary if you run a single server to access a given data directory (which is the usual case) and if no other programs such as myisamchk need to modify tables while the server is running. If you only *read* tables with other programs, external locking is not required, although myisamchk might report warnings if the server changes tables while myisamchk is reading them.

With external locking disabled, to use myisamchk, you must either stop the server while myisamchk executes or else lock and flush the tables before running myisamchk. To avoid this requirement, use the CHECK TABLE and REPAIR TABLE statements to check and repair MyISAM tables.

For mysqld, external locking is controlled by the value of the skip\_external\_locking system variable. When this variable is enabled, external locking is disabled, and vice versa. External locking is disabled by default.

Use of external locking can be controlled at server startup by using the --external-locking or --skip-external-locking option.

If you do use external locking option to enable updates to  ${\tt MyISAM}$  tables from many MySQL processes, do not start the server with the  ${\tt delay\_key\_write}$  system variable set to ALL or use the DELAY\_KEY\_WRITE=1 table option for any shared tables. Otherwise, index corruption can occur.

The easiest way to satisfy this condition is to always use --external-locking together with -- delay-key-write=OFF. (This is not done by default because in many setups it is useful to have a mixture of the preceding options.)

# 10.12 Optimizing the MySQL Server

This section discusses optimization techniques for the database server, primarily dealing with system configuration rather than tuning SQL statements. The information in this section is appropriate for DBAs who want to ensure performance and scalability across the servers they manage; for developers constructing installation scripts that include setting up the database; and people running MySQL themselves for development, testing, and so on who want to maximize their own productivity.

# 10.12.1 Optimizing Disk I/O

This section describes ways to configure storage devices when you can devote more and faster storage hardware to the database server. For information about optimizing an InnoDB configuration to improve I/O performance, see Section 10.5.8, "Optimizing InnoDB Disk I/O".

- Disk seeks are a huge performance bottleneck. This problem becomes more apparent when
  the amount of data starts to grow so large that effective caching becomes impossible. For large
  databases where you access data more or less randomly, you can be sure that you need at least
  one disk seek to read and a couple of disk seeks to write things. To minimize this problem, use disks
  with low seek times.
- Increase the number of available disk spindles (and thereby reduce the seek overhead) by either symlinking files to different disks or striping the disks:
  - · Using symbolic links

This means that, for MyISAM tables, you symlink the index file and data files from their usual location in the data directory to another disk (that may also be striped). This makes both the seek and read times better, assuming that the disk is not used for other purposes as well. See Section 10.12.2, "Using Symbolic Links".

Symbolic links are not supported for use with InnoDB tables. However, it is possible to place InnoDB data and log files on different physical disks. For more information, see Section 10.5.8, "Optimizing InnoDB Disk I/O".

#### Striping

Striping means that you have many disks and put the first block on the first disk, the second block on the second disk, and the *N*-th block on the (*N MOD number\_of\_disks*) disk, and so on. This means if your normal data size is less than the stripe size (or perfectly aligned), you get much better performance. Striping is very dependent on the operating system and the stripe size, so benchmark your application with different stripe sizes. See Section 10.13.2, "Using Your Own Benchmarks".

The speed difference for striping is *very* dependent on the parameters. Depending on how you set the striping parameters and number of disks, you may get differences measured in orders of magnitude. You have to choose to optimize for random or sequential access.

- For reliability, you may want to use RAID 0+1 (striping plus mirroring), but in this case, you need
  2 × N drives to hold N drives of data. This is probably the best option if you have the money for it.
  However, you may also have to invest in some volume-management software to handle it efficiently.
- A good option is to vary the RAID level according to how critical a type of data is. For example, store semi-important data that can be regenerated on a RAID 0 disk, but store really important data such as host information and logs on a RAID 0+1 or RAID N disk. RAID N can be a problem if you have many writes, due to the time required to update the parity bits.
- You can also set the parameters for the file system that the database uses:

If you do not need to know when files were last accessed (which is not really useful on a database server), you can mount your file systems with the -o noatime option. That skips updates to the last access time in inodes on the file system, which avoids some disk seeks.

On many operating systems, you can set a file system to be updated asynchronously by mounting it with the -o async option. If your computer is reasonably stable, this should give you better performance without sacrificing too much reliability. (This flag is on by default on Linux.)

### Using NFS with MySQL

You should be cautious when considering whether to use NFS with MySQL. Potential issues, which vary by operating system and NFS version, include the following:

- MySQL data and log files placed on NFS volumes becoming locked and unavailable for use. Locking
  issues may occur in cases where multiple instances of MySQL access the same data directory
  or where MySQL is shut down improperly, due to a power outage, for example. NFS version 4
  addresses underlying locking issues with the introduction of advisory and lease-based locking.
  However, sharing a data directory among MySQL instances is not recommended.
- Data inconsistencies introduced due to messages received out of order or lost network traffic. To avoid this issue, use TCP with hard and intr mount options.
- Maximum file size limitations. NFS Version 2 clients can only access the lowest 2GB of a file (signed 32 bit offset). NFS Version 3 clients support larger files (up to 64 bit offsets). The maximum supported file size also depends on the local file system of the NFS server.

Using NFS within a professional SAN environment or other storage system tends to offer greater reliability than using NFS outside of such an environment. However, NFS within a SAN environment may be slower than directly attached or bus-attached non-rotational storage.

If you choose to use NFS, NFS Version 4 or later is recommended, as is testing your NFS setup thoroughly before deploying into a production environment.

## 10.12.2 Using Symbolic Links

You can move databases or tables from the database directory to other locations and replace them with symbolic links to the new locations. You might want to do this, for example, to move a database to a file system with more free space or increase the speed of your system by spreading your tables to different disks.

For InnoDB tables, use the DATA DIRECTORY clause of the CREATE TABLE statement instead of symbolic links, as explained in Section 17.6.1.2, "Creating Tables Externally". This new feature is a supported, cross-platform technique.

The recommended way to do this is to symlink entire database directories to a different disk. Symlink MyISAM tables only as a last resort.

To determine the location of your data directory, use this statement:

SHOW VARIABLES LIKE 'datadir';

### 10.12.2.1 Using Symbolic Links for Databases on Unix

On Unix, symlink a database using this procedure:

1. Create the database using CREATE DATABASE:

```
mysql> CREATE DATABASE mydb1;
```

Using CREATE DATABASE creates the database in the MySQL data directory and permits the server to update the data dictionary with information about the database directory.

- 2. Stop the server to ensure that no activity occurs in the new database while it is being moved.
- 3. Move the database directory to some disk where you have free space. For example, use tar or mv. If you use a method that copies rather than moves the database directory, remove the original database directory after copying it.
- 4. Create a soft link in the data directory to the moved database directory:

```
$> ln -s /path/to/mydb1 /path/to/datadir
```

The command creates a symlink named mydb1 in the data directory.

5. Restart the server.

### 10.12.2.2 Using Symbolic Links for MylSAM Tables on Unix



#### Note

Symbolic link support as described here, along with the --symbolic-links option that controls it, and is deprecated; expect these to be removed in a future version of MySQL. In addition, the option is disabled by default.

Symlinks are fully supported only for MyISAM tables. For files used by tables for other storage engines, you may get strange problems if you try to use symbolic links. For InnoDB tables, use the alternative technique explained in Section 17.6.1.2, "Creating Tables Externally" instead.

Do not symlink tables on systems that do not have a fully operational <code>realpath()</code> call. (Linux and Solaris support <code>realpath()</code>). To determine whether your system supports symbolic links, check the value of the <code>have symlink</code> system variable using this statement:

SHOW VARIABLES LIKE 'have\_symlink';

The handling of symbolic links for MyISAM tables works as follows:

- In the data directory, you always have the data (.MYD) file and the index (.MYI) file. The data file and index file can be moved elsewhere and replaced in the data directory by symlinks.
- You can symlink the data file and the index file independently to different directories.
- To instruct a running MySQL server to perform the symlinking, use the DATA DIRECTORY and INDEX DIRECTORY options to CREATE TABLE. See Section 15.1.20, "CREATE TABLE Statement". Alternatively, if mysqld is not running, symlinking can be accomplished manually using ln -s from the command line.



#### Note

The path used with either or both of the DATA DIRECTORY and INDEX DIRECTORY options may not include the MySQL data directory. (Bug #32167)

• myisamchk does not replace a symlink with the data file or index file. It works directly on the file to which the symlink points. Any temporary files are created in the directory where the data file or index file is located. The same is true for the ALTER TABLE, OPTIMIZE TABLE, and REPAIR TABLE statements.



### Note

When you drop a table that is using symlinks, both the symlink and the file to which the symlink points are dropped. This is an extremely good reason not to run mysqld as the root operating system user or permit operating system users to have write access to MySQL database directories.

- If you rename a table with ALTER TABLE ... RENAME or RENAME TABLE and you do not move the table to another database, the symlinks in the database directory are renamed to the new names and the data file and index file are renamed accordingly.
- If you use ALTER TABLE ... RENAME or RENAME TABLE to move a table to another database, the table is moved to the other database directory. If the table name changed, the symlinks in the new database directory are renamed to the new names and the data file and index file are renamed accordingly.
- If you are not using symlinks, start mysqld with the --skip-symbolic-links option to ensure that no one can use mysqld to drop or rename a file outside of the data directory.

These table symlink operations are not supported:

• ALTER TABLE ignores the DATA DIRECTORY and INDEX DIRECTORY table options.

### 10.12.2.3 Using Symbolic Links for Databases on Windows

On Windows, symbolic links can be used for database directories. This enables you to put a database directory at a different location (for example, on a different disk) by setting up a symbolic link to it. Use of database symlinks on Windows is similar to their use on Unix, although the procedure for setting up the link differs.

Suppose that you want to place the database directory for a database named mydb at D:\data\mydb. To do this, create a symbolic link in the MySQL data directory that points to D:\data\mydb. However, before creating the symbolic link, make sure that the D:\data\mydb directory exists by creating it if necessary. If you already have a database directory named mydb in the data directory, move it to D:\data. Otherwise, the symbolic link has no effect. To avoid problems, make sure that the server is not running when you move the database directory.

On Windows, you can create a symlink using the mklink command. This command requires administrative privileges.

- 1. Make sure that the desired path to the database exists. For this example, we use D:\data\mydb, and a database named mydb.
- 2. If the database does not already exist, issue CREATE DATABASE mydb in the mysql client to create it.
- 3. Stop the MySQL service.
- 4. Using Windows Explorer or the command line, move the directory mydb from the data directory to D:\data, replacing the directory of the same name.
- 5. If you are not already using the command prompt, open it, and change location to the data directory, like this:

```
C:\> cd \path\to\datadir
```

If your MySQL installation is in the default location, you can use this:

```
C:\> cd C:\ProgramData\MySQL\MySQL Server 8.4\Data
```

6. In the data directory, create a symlink named mydb that points to the location of the database directory:

```
C:\> mklink /d mydb D:\data\mydb
```

7. Start the MySQL service.

After this, all tables created in the database mydb are created in D:\data\mydb.

Alternatively, on any version of Windows supported by MySQL, you can create a symbolic link to a MySQL database by creating a .sym file in the data directory that contains the path to the destination directory. The file should be named  $db\_name$ .sym, where  $db\_name$  is the database name.

Support for database symbolic links on Windows using .sym files is enabled by default. If you do not need .sym file symbolic links, you can disable support for them by starting mysqld with the --skip-symbolic-links option. To determine whether your system supports .sym file symbolic links, check the value of the have\_symlink system variable using this statement:

```
SHOW VARIABLES LIKE 'have_symlink';
```

To create a .sym file symlink, use this procedure:

1. Change location into the data directory:

C:\> cd \path\to\datadir

2. In the data directory, create a text file named mydb.sym that contains this path name: D:\data\mydb\



#### Note

The path name to the new database and tables should be absolute. If you specify a relative path, the location is relative to the mydb.sym file.

After this, all tables created in the database mydb are created in D: \data\mydb.

## 10.12.3 Optimizing Memory Use

## 10.12.3.1 How MySQL Uses Memory

MySQL allocates buffers and caches to improve performance of database operations. The default configuration is designed to permit a MySQL server to start on a virtual machine that has approximately 512MB of RAM. You can improve MySQL performance by increasing the values of certain cache and buffer-related system variables. You can also modify the default configuration to run MySQL on systems with limited memory.

The following list describes some of the ways that MySQL uses memory. Where applicable, relevant system variables are referenced. Some items are storage engine or feature specific.

The InnoDB buffer pool is a memory area that holds cached InnoDB data for tables, indexes, and
other auxiliary buffers. For efficiency of high-volume read operations, the buffer pool is divided into
pages that can potentially hold multiple rows. For efficiency of cache management, the buffer pool
is implemented as a linked list of pages; data that is rarely used is aged out of the cache, using a
variation of the LRU algorithm. For more information, see Section 17.5.1, "Buffer Pool".

The size of the buffer pool is important for system performance:

- InnoDB allocates memory for the entire buffer pool at server startup, using malloc() operations. The innodb\_buffer\_pool\_size system variable defines the buffer pool size. Typically, a recommended innodb\_buffer\_pool\_size value is 50 to 75 percent of system memory. innodb\_buffer\_pool\_size can be configured dynamically, while the server is running. For more information, see Section 17.8.3.1, "Configuring InnoDB Buffer Pool Size".
- On systems with a large amount of memory, you can improve concurrency by dividing the buffer pool into multiple buffer pool instances. The innodb\_buffer\_pool\_instances system variable defines the number of buffer pool instances.
- A buffer pool that is too small may cause excessive churning as pages are flushed from the buffer pool only to be required again a short time later.
- A buffer pool that is too large may cause swapping due to competition for memory.
- The storage engine interface enables the optimizer to provide information about the size of the
  record buffer to be used for scans that the optimizer estimates are likely to read multiple rows. The
  buffer size can vary based on the size of the estimate. InnobB uses this variable-size buffering
  capability to take advantage of row prefetching, and to reduce the overhead of latching and B-tree
  navigation.
- All threads share the MyISAM key buffer. The key\_buffer\_size system variable determines its size.

For each MyISAM table the server opens, the index file is opened once; the data file is opened once for each concurrently running thread that accesses the table. For each concurrent thread, a table structure, column structures for each column, and a buffer of size 3 \* N are allocated (where N is

the maximum row length, not counting BLOB columns). A BLOB column requires five to eight bytes plus the length of the BLOB data. The MyISAM storage engine maintains one extra row buffer for internal use.

- The myisam\_use\_mmap system variable can be set to 1 to enable memory-mapping for all MyISAM tables.
- If an internal in-memory temporary table becomes too large (as determined by tmp\_table\_size and max\_heap\_table\_size), MySQL automatically converts the table from in-memory to on-disk format, which uses the InnoDB storage engine. You can increase the permissible temporary table size as described in Section 10.4.4, "Internal Temporary Table Use in MySQL".

For MEMORY tables explicitly created with CREATE TABLE, only the max\_heap\_table\_size system variable determines how large a table can grow, and there is no conversion to on-disk format.

- The MySQL Performance Schema is a feature for monitoring MySQL server execution at a low level.
  The Performance Schema dynamically allocates memory incrementally, scaling its memory use to
  actual server load, instead of allocating required memory during server startup. Once memory is
  allocated, it is not freed until the server is restarted. For more information, see Section 29.17, "The
  Performance Schema Memory-Allocation Model".
- Each thread that the server uses to manage client connections requires some thread-specific space. The following list indicates these and which system variables control their size:
  - A stack (thread stack)
  - A connection buffer (net buffer length)
  - A result buffer (net\_buffer\_length)

The connection buffer and result buffer each begin with a size equal to net\_buffer\_length bytes, but are dynamically enlarged up to max\_allowed\_packet bytes as needed. The result buffer shrinks to net\_buffer\_length bytes after each SQL statement. While a statement is running, a copy of the current statement string is also allocated.

Each connection thread uses memory for computing statement digests. The server allocates <code>max\_digest\_length</code> bytes per session. See Section 29.10, "Performance Schema Statement Digests and Sampling".

- · All threads share the same base memory.
- When a thread is no longer needed, the memory allocated to it is released and returned to the system unless the thread goes back into the thread cache. In that case, the memory remains allocated.
- Each request that performs a sequential scan of a table allocates a *read buffer*. The read buffer size system variable determines the buffer size.
- When reading rows in an arbitrary sequence (for example, following a sort), a random-read buffer
  may be allocated to avoid disk seeks. The read\_rnd\_buffer\_size system variable determines
  the buffer size.
- All joins are executed in a single pass, and most joins can be done without even using a temporary table. Most temporary tables are memory-based hash tables. Temporary tables with a large row length (calculated as the sum of all column lengths) or that contain BLOB columns are stored on disk.
- Most requests that perform a sort allocate a sort buffer and zero to two temporary files depending on the result set size. See Section B.3.3.5, "Where MySQL Stores Temporary Files".
- Almost all parsing and calculating is done in thread-local and reusable memory pools. No memory
  overhead is needed for small items, thus avoiding the normal slow memory allocation and freeing.
  Memory is allocated only for unexpectedly large strings.

- For each table having BLOB columns, a buffer is enlarged dynamically to read in larger BLOB values. If you scan a table, the buffer grows as large as the largest BLOB value.
- MySQL requires memory and descriptors for the table cache. Handler structures for all in-use tables
  are saved in the table cache and managed as "First In, First Out" (FIFO). The table\_open\_cache
  system variable defines the initial table cache size; see Section 10.4.3.1, "How MySQL Opens and
  Closes Tables".

MySQL also requires memory for the table definition cache. The table\_definition\_cache system variable defines the number of table definitions that can be stored in the table definition cache. If you use a large number of tables, you can create a large table definition cache to speed up the opening of tables. The table definition cache takes less space and does not use file descriptors, unlike the table cache.

- A FLUSH TABLES statement or mysqladmin flush-tables command closes all tables that are not in use at once and marks all in-use tables to be closed when the currently executing thread finishes. This effectively frees most in-use memory. FLUSH TABLES does not return until all tables have been closed.
- The server caches information in memory 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.
- In a replication topology, the following settings affect memory usage, and can be adjusted as required:
  - The max\_allowed\_packet system variable on a replication source limits the maximum message size that the source sends to its replicas for processing. This setting defaults to 64M.
  - The system variable replica\_pending\_jobs\_size\_max on a multithreaded replica sets the
    maximum amount of memory that is made available for holding messages awaiting processing.
    This setting defaults to 128M. The memory is only allocated when needed, but it might be used
    if your replication topology handles large transactions sometimes. It is a soft limit, and larger
    transactions can be processed.
  - The rpl\_read\_size system variable on a replication source or replica controls the minimum amount of data in bytes that is read from the binary log files and relay log files. The default is 8192 bytes. A buffer the size of this value is allocated for each thread that reads from the binary log and relay log files, including dump threads on sources and coordinator threads on replicas.
  - The binlog\_transaction\_dependency\_history\_size system variable limits the number of row hashes held as an in-memory history.
  - The max\_binlog\_cache\_size system variable specifies the upper limit of memory usage by an individual transaction.
  - The max\_binlog\_stmt\_cache\_size system variable specifies the upper limit of memory usage by the statement cache.

ps and other system status programs may report that mysqld uses a lot of memory. This may be caused by thread stacks on different memory addresses. For example, the Solaris version of ps counts the unused memory between stacks as used memory. To verify this, check available swap with swap -s. We test mysqld with several memory-leakage detectors (both commercial and Open Source), so there should be no memory leaks.

### 10.12.3.2 Monitoring MySQL Memory Usage

The following example demonstrates how to use Performance Schema and sys schema to monitor MySQL memory usage.

Most Performance Schema memory instrumentation is disabled by default. Instruments can be enabled by updating the ENABLED column of the Performance Schema setup\_instruments table. Memory instruments have names in the form of memory/code\_area/instrument\_name, where code\_area is a value such as sql or innodb, and instrument\_name is the instrument detail.

1. To view available MySQL memory instruments, query the Performance Schema setup\_instruments table. The following query returns hundreds of memory instruments for all code areas.

```
mysql> SELECT * FROM performance_schema.setup_instruments
    WHERE NAME LIKE '%memory%';
```

You can narrow results by specifying a code area. For example, you can limit results to InnoDB memory instruments by specifying innodb as the code area.

```
mysql> SELECT * FROM performance_schema.setup_instruments
      WHERE NAME LIKE '%memory/innodb%';
NAME
                                         | ENABLED | TIMED
 memory/innodb/adaptive hash index
                                                    NO
 memory/innodb/buf_buf_pool
                                          NO
                                                    | NO
 memory/innodb/dict_stats_bg_recalc_pool_t | NO
                                                     NO
 memory/innodb/dict_stats_index_map_t
                                                     NO
                                           NO
 memory/innodb/dict_stats_n_diff_on_level NO
                                                    | NO
                                          NO
 memory/innodb/other
                                                     NO
 memory/innodb/row_log_buf
                                           NO
                                                     NO
 memory/innodb/row_merge_sort
                                          l NO
                                                     NO
 memory/innodb/std
                                                     NO
                                           NO
 memory/innodb/trx_sys_t::rw_trx_ids
```

Depending on your MySQL installation, code areas may include performance\_schema, sql, client, innodb, myisam, csv, memory, blackhole, archive, partition, and others.

 To enable memory instruments, add a performance-schema-instrument rule to your MySQL configuration file. For example, to enable all memory instruments, add this rule to your configuration file and restart the server:

performance-schema-instrument='memory/%=COUNTED'



#### Note

Enabling memory instruments at startup ensures that memory allocations that occur at startup are counted.

After restarting the server, the ENABLED column of the Performance Schema setup\_instruments table should report YES for memory instruments that you enabled. The TIMED column in the setup\_instruments table is ignored for memory instruments because memory operations are not timed.

```
mysql> SELECT * FROM performance_schema.setup_instruments
     WHERE NAME LIKE '%memory/innodb%';
                                         | ENABLED | TIMED
 memory/innodb/adaptive hash index
                                                     NO
 memory/innodb/buf_buf_pool
                                            NO
                                                      NO
 memory/innodb/dict_stats_bg_recalc_pool_t | NO
                                                      NO
                                           NO
 memory/innodb/dict_stats_index_map_t
                                                      NO
 memory/innodb/dict_stats_n_diff_on_level NO
                                                      NO
 memory/innodb/other
                                            NO
                                                      NO
 memory/innodb/row_log_buf
                                            NO
                                                      NO
 memory/innodb/row_merge_sort
                                            NO
                                                      NO
 memory/innodb/std
                                            NO
                                                      NO
 memory/innodb/trx_sys_t::rw_trx_ids
                                           NO
                                                      NO
```

3. Query memory instrument data. In this example, memory instrument data is queried in the Performance Schema memory\_summary\_global\_by\_event\_name table, which summarizes data by EVENT\_NAME. The EVENT\_NAME is the name of the instrument.

The following query returns memory data for the InnoDB buffer pool. For column descriptions, see Section 29.12.20.10, "Memory Summary Tables".

The same underlying data can be queried using the sys schema memory\_global\_by\_current\_bytes table, which shows current memory usage within the server globally, broken down by allocation type.

This sys schema query aggregates currently allocated memory (current\_alloc) by code area:

For more information about sys schema, see Chapter 30, MySQL sys Schema.

### 10.12.3.3 Enabling Large Page Support

Some hardware and operating system architectures support memory pages greater than the default (usually 4KB). The actual implementation of this support depends on the underlying hardware and operating system. Applications that perform a lot of memory accesses may obtain performance improvements by using large pages due to reduced Translation Lookaside Buffer (TLB) misses.

In MySQL, large pages can be used by InnoDB, to allocate memory for its buffer pool and additional memory pool.

Standard use of large pages in MySQL attempts to use the largest size supported, up to 4MB. Under Solaris, a "super large pages" feature enables uses of pages up to 256MB. This feature is available for recent SPARC platforms. It can be enabled or disabled by using the --super-large-pages or --skip-super-large-pages option.

MySQL also supports the Linux implementation of large page support (which is called HugeTLB in Linux).

Before large pages can be used on Linux, the kernel must be enabled to support them and it is necessary to configure the HugeTLB memory pool. For reference, the HugeTBL API is documented in the <code>Documentation/vm/hugetlbpage.txt</code> file of your Linux sources.

The kernels for some recent systems such as Red Hat Enterprise Linux may have the large pages feature enabled by default. To check whether this is true for your kernel, use the following command and look for output lines containing "huge":

```
$> grep -i huge /proc/meminfo
AnonHugePages: 2658304 kB
ShmemHugePages: 0 kB
HugePages_Total: 0
HugePages_Free: 0
HugePages_Rsvd: 0
HugePages_Surp: 0
HugePages_Surp: 0
Hugepagesize: 2048 kB
Hugetlb: 0 kB
```

The nonempty command output indicates that large page support is present, but the zero values indicate that no pages are configured for use.

If your kernel needs to be reconfigured to support large pages, consult the hugetlbpage.txt file for instructions.

Assuming that your Linux kernel has large page support enabled, configure it for use by MySQL using the following steps:

- Determine the number of large pages needed. This is the size of the InnoDB buffer pool divided by
  the large page size, which we can calculate as innodb\_buffer\_pool\_size/Hugepagesize.
  Assuming the default value for the innodb\_buffer\_pool\_size (128MB) and using the
  Hugepagesize value obtained from /proc/meminfo (2MB), this is 128MB / 2MB, or 64 Huge
  Pages. We call this value P.
- 2. As system root, open the file /etc/sysctl.conf in a text editor, and add the line shown here, where *P* is the number of large pages obtained in the previous step:

```
vm.nr_hugepages=P
```

Using the actual value obtained previously, the additional line should look like this:

```
vm.nr_hugepages=66
```

Save the updated file.

3. As system root, run the following command:

```
$> sudo sysctl -p
```



### Note

On some systems the large pages file may be named slightly differently; for example, some distributions call it nr\_hugepages. In the event sysctl returns an error relating to the file name, check the name of the corresponding file in /proc/sys/vm and use that instead.

To verify the large page configuration, check /proc/meminfo again as described previously. Now you should see some additional nonzero values in the output, similar to this:

```
$> grep -i huge /proc/meminfo
AnonHugePages: 2686976 kB
ShmemHugePages: 0 kB
HugePages_Total: 233
HugePages_Free: 233
HugePages_Rsvd: 0
HugePages_Surp: 0
Hugepagesize: 2048 kB
Hugetlb: 477184 kB
```

4. Optionally, you may wish to compact the Linux VM. You can do this using a sequence of commands, possibly in a script file, similar to what is shown here:

```
sync
sync
sync
echo 3 > /proc/sys/vm/drop_caches
echo 1 > /proc/sys/vm/compact_memory
```

See your operating platform documentation for more information about how to do this.

- 5. Check any configuration files such as my.cnf used by the server, and make sure that innodb\_buffer\_pool\_chunk\_size is set larger than the huge page size. The default for this variable is 128M.
- 6. Large page support in the MySQL server is disabled by default. To enable it, start the server with --large-pages. You can also do so by adding the following line to the [mysqld] section of the server my.cnf file:

```
large-pages=ON
```

With this option enabled, InnoDB uses large pages automatically for its buffer pool and additional memory pool. If InnoDB cannot do this, it falls back to use of traditional memory and writes a warning to the error log: Warning: Using conventional memory pool.

You can verify that MySQL is now using large pages by checking /proc/meminfo again after restarting mysqld, like this:

```
$> grep -i huge /proc/meminfo
AnonHugePages: 2516992 kB
ShmemHugePages: 0 kB
HugePages_Total: 233
HugePages_Free: 222
HugePages_Rsvd: 55
HugePages_Surp: 0
Hugepagesize: 2048 kB
Hugetlb: 477184 kB
```

# 10.13 Measuring Performance (Benchmarking)

To measure performance, consider the following factors:

- Whether you are measuring the speed of a single operation on a quiet system, or how a set of
  operations (a "workload") works over a period of time. With simple tests, you usually test how
  changing one aspect (a configuration setting, the set of indexes on a table, the SQL clauses in a
  query) affects performance. Benchmarks are typically long-running and elaborate performance tests,
  where the results could dictate high-level choices such as hardware and storage configuration, or
  how soon to upgrade to a new MySQL version.
- For benchmarking, sometimes you must simulate a heavy database workload to get an accurate picture.

- Performance can vary depending on so many different factors that a difference of a few percentage
  points might not be a decisive victory. The results might shift the opposite way when you test in a
  different environment.
- Certain MySQL features help or do not help performance depending on the workload. For completeness, always test performance with those features turned on and turned off. The most important feature to try with each workload is the adaptive hash index for InnoDB tables.

This section progresses from simple and direct measurement techniques that a single developer can do, to more complicated ones that require additional expertise to perform and interpret the results.

## 10.13.1 Measuring the Speed of Expressions and Functions

To measure the speed of a specific MySQL expression or function, invoke the BENCHMARK() function using the mysql client program. Its syntax is  $BENCHMARK(loop\_count,expr)$ . The return value is always zero, but mysql prints a line displaying approximately how long the statement took to execute. For example:

This result was obtained on a Pentium II 400MHz system. It shows that MySQL can execute 1,000,000 simple addition expressions in 0.32 seconds on that system.

The built-in MySQL functions are typically highly optimized, but there may be some exceptions. BENCHMARK() is an excellent tool for finding out if some function is a problem for your queries.

# 10.13.2 Using Your Own Benchmarks

Benchmark your application and database to find out where the bottlenecks are. After fixing one bottleneck (or by replacing it with a "dummy" module), you can proceed to identify the next bottleneck. Even if the overall performance for your application currently is acceptable, you should at least make a plan for each bottleneck and decide how to solve it if someday you really need the extra performance.

A free benchmark suite is the Open Source Database Benchmark, available at http://osdb.sourceforge.net/.

It is very common for a problem to occur only when the system is very heavily loaded. We have had many customers who contact us when they have a (tested) system in production and have encountered load problems. In most cases, performance problems turn out to be due to issues of basic database design (for example, table scans are not good under high load) or problems with the operating system or libraries. Most of the time, these problems would be much easier to fix if the systems were not already in production.

To avoid problems like this, benchmark your whole application under the worst possible load:

- The mysqlslap program can be helpful for simulating a high load produced by multiple clients issuing queries simultaneously. See Section 6.5.7, "mysqlslap A Load Emulation Client".
- You can also try benchmarking packages such as SysBench and DBT2, available at https://launchpad.net/sysbench, and http://osdldbt.sourceforge.net/#dbt2.

These programs or packages can bring a system to its knees, so be sure to use them only on your development systems.

# 10.13.3 Measuring Performance with performance\_schema

You can query the tables in the performance\_schema database to see real-time information about the performance characteristics of your server and the applications it is running. See Chapter 29, *MySQL Performance Schema* for details.

# 10.14 Examining Server Thread (Process) Information

To ascertain what your MySQL server is doing, it can be helpful to examine the process list, which indicates the operations currently being performed by the set of threads executing within the server. For example:

```
mysql> SHOW PROCESSLIST\G
      ****************** 1. row *****************
    Id: 5
  User: event_scheduler
  Host: localhost
    db: NULL
Command: Daemon
  Time: 2756681
 State: Waiting on empty queue
  Info: NULL
              ******* 2. row ****************
    Id: 20
  User: me
  Host: localhost:52943
   db: test
Command: Ouery
  Time: 0
 State: starting
  Info: SHOW PROCESSLIST
```

Threads can be killed with the KILL statement. See Section 15.7.8.4, "KILL Statement".

## 10.14.1 Accessing the Process List

The following discussion enumerates the sources of process information, the privileges required to see process information, and describes the content of process list entries.

- Sources of Process Information
- · Privileges Required to Access the Process List
- · Content of Process List Entries

### **Sources of Process Information**

Process information is available from these sources:

- The SHOW PROCESSLIST statement: Section 15.7.7.31, "SHOW PROCESSLIST Statement"
- The mysqladmin processlist command: Section 6.5.2, "mysqladmin A MySQL Server Administration Program"
- The INFORMATION\_SCHEMA PROCESSLIST table: Section 28.3.23, "The INFORMATION\_SCHEMA PROCESSLIST Table"
- The Performance Schema processlist table: Section 29.12.22.7, "The processlist Table"
- The Performance Schema threads table columns with names having a prefix of PROCESSLIST\_: Section 29.12.22.8, "The threads Table"
- The sys schema processlist and session views: Section 30.4.3.22, "The processlist and x \$processlist Views", and Section 30.4.3.33, "The session and x\$session Views"

The threads table compares to SHOW PROCESSLIST, INFORMATION\_SCHEMA PROCESSLIST, and mysgladmin processlist as follows:

 Access to the threads table does not require a mutex and has minimal impact on server performance. The other sources have negative performance consequences because they require a mutex.



#### Note

An alternative implementation for SHOW PROCESSLIST is available based on the Performance Schema processlist table, which, like the threads table, does not require a mutex and has better performance characteristics. For details, see Section 29.12.22.7, "The processlist Table".

- The threads table displays background threads, which the other sources do not. It also provides
  additional information for each thread that the other sources do not, such as whether the thread is a
  foreground or background thread, and the location within the server associated with the thread. This
  means that the threads table can be used to monitor thread activity the other sources cannot.
- You can enable or disable Performance Schema thread monitoring, as described in Section 29.12.22.8, "The threads Table".

For these reasons, DBAs who perform server monitoring using one of the other thread information sources may wish to monitor using the threads table instead.

The sys schema processlist view presents information from the Performance Schema threads table in a more accessible format. The sys schema session view presents information about user sessions like the sys schema processlist view, but with background processes filtered out.

### **Privileges Required to Access the Process List**

For most sources of process information, 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.

The Performance Schema threads table also provides thread information, but table access uses a different privilege model. See Section 29.12.22.8, "The threads Table".

### **Content of Process List Entries**

Each process list entry contains several pieces of information. The following list describes them using the labels from SHOW PROCESSLIST output. Other process information sources use similar labels.

- Id is the connection identifier for the client associated with the thread.
- User and Host indicate the account associated with the thread.
- db is the default database for the thread, or NULL if none has been selected.
- Command and State indicate what the thread is doing.

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.

The following sections list the possible Command values, and State values grouped by category. The meaning for some of these values is self-evident. For others, additional description is provided.



#### **Note**

Applications that examine process list information should be aware that the commands and states are subject to change.

• Time indicates how long the thread has been in its current state. The thread's notion of the current time may be altered in some cases: The thread can change the time with SET TIMESTAMP =

*value*. 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 19.2.3, "Replication Threads".

 Info indicates the statement the thread is executing, or NULL if it is executing no statement. For SHOW PROCESSLIST, this value contains only the first 100 characters of the statement. To see complete statements, use SHOW FULL PROCESSLIST (or query a different process information source).

### 10.14.2 Thread Command Values

A thread can have any of the following Command values:

• Binlog Dump

This is a thread on a replication source for sending binary log contents to a replica.

Change user

The thread is executing a change user operation.

• Close stmt

The thread is closing a prepared statement.

• Connect

Used by replication receiver threads connected to the source, and by replication worker threads.

• Connect Out

A replica is connecting to its source.

• Create DB

The thread is executing a create database operation.

Daemon

This thread is internal to the server, not a thread that services a client connection.

Debug

The thread is generating debugging information.

• Delayed insert

The thread is a delayed insert handler.

• Drop DB

The thread is executing a drop database operation.

- Error
- Execute

The thread is executing a prepared statement.

Fetch

The thread is fetching the results from executing a prepared statement.

• Field List

The thread is retrieving information for table columns.

• Init DB

The thread is selecting a default database.

Kill

The thread is killing another thread.

Long Data

The thread is retrieving long data in the result of executing a prepared statement.

Ping

The thread is handling a server ping request.

• Prepare

The thread is preparing a prepared statement.

Processlist

The thread is producing information about server threads.

Query

Employed for user clients while executing queries by single-threaded replication applier threads, as well as by the replication coordinator thread.

• Ouit

The thread is terminating.

• Refresh

The thread is flushing table, logs, or caches, or resetting status variable or replication server information.

• Register Slave

The thread is registering a replica server.

• Reset stmt

The thread is resetting a prepared statement.

• Set option

The thread is setting or resetting a client statement execution option.

• Shutdown

The thread is shutting down the server.

• Sleep

The thread is waiting for the client to send a new statement to it.

Statistics

The thread is producing server status information.

Time

Unused.

### 10.14.3 General Thread States

The following list describes thread State values that are associated with general query processing and not more specialized activities such as replication. Many of these are useful only for finding bugs in the server.

• After create

This occurs when the thread creates a table (including internal temporary tables), at the end of the function that creates the table. This state is used even if the table could not be created due to some error.

altering table

The server is in the process of executing an in-place ALTER TABLE.

Analyzing

The thread is calculating a MyISAM table key distributions (for example, for ANALYZE TABLE).

checking permissions

The thread is checking whether the server has the required privileges to execute the statement.

Checking table

The thread is performing a table check operation.

• cleaning up

The thread has processed one command and is preparing to free memory and reset certain state variables.

closing tables

The thread is flushing the changed table data to disk and closing the used tables. This should be a fast operation. If not, verify that you do not have a full disk and that the disk is not in very heavy use.

committing alter table to storage engine

The server has finished an in-place ALTER TABLE and is committing the result.

converting HEAP to ondisk

The thread is converting an internal temporary table from a MEMORY table to an on-disk table.

copy to tmp table

The thread is processing an ALTER TABLE statement. This state occurs after the table with the new structure has been created but before rows are copied into it.

For a thread in this state, the Performance Schema can be used to obtain about the progress of the copy operation. See Section 29.12.5, "Performance Schema Stage Event Tables".

• Copying to group table

If a statement has different ORDER BY and GROUP BY criteria, the rows are sorted by group and copied to a temporary table.

Copying to tmp table

The server is copying to a temporary table in memory.

· Copying to tmp table on disk

The server is copying to a temporary table on disk. The temporary result set has become too large (see Section 10.4.4, "Internal Temporary Table Use in MySQL"). Consequently, the thread is changing the temporary table from in-memory to disk-based format to save memory.

• Creating index

The thread is processing ALTER TABLE ... ENABLE KEYS for a MyISAM table.

• Creating sort index

The thread is processing a SELECT that is resolved using an internal temporary table.

• creating table

The thread is creating a table. This includes creation of temporary tables.

• Creating tmp table

The thread is creating a temporary table in memory or on disk. If the table is created in memory but later is converted to an on-disk table, the state during that operation is Copying to tmp table on disk.

• deleting from main table

The server is executing the first part of a multiple-table delete. It is deleting only from the first table, and saving columns and offsets to be used for deleting from the other (reference) tables.

• deleting from reference tables

The server is executing the second part of a multiple-table delete and deleting the matched rows from the other tables.

• discard\_or\_import\_tablespace

The thread is processing an ALTER TABLE ... DISCARD TABLESPACE or ALTER TABLE ... IMPORT TABLESPACE statement.

end

This occurs at the end but before the cleanup of ALTER TABLE, CREATE VIEW, DELETE, INSERT, SELECT, or UPDATE statements.

For the end state, the following operations could be happening:

- · Writing an event to the binary log
- · Freeing memory buffers, including for blobs
- executing

The thread has begun executing a statement.

• Execution of init\_command

The thread is executing statements in the value of the init\_command system variable.

freeing items

The thread has executed a command. This state is usually followed by cleaning up.

FULLTEXT initialization

The server is preparing to perform a natural-language full-text search.

init

This occurs before the initialization of ALTER TABLE, DELETE, INSERT, SELECT, or UPDATE statements. Actions taken by the server in this state include flushing the binary log and the InnoDB log.

• Killed

Someone has sent a KILL statement to the thread and it should abort next time it checks the kill flag. The flag is checked in each major loop in MySQL, but in some cases it might still take a short time for the thread to die. If the thread is locked by some other thread, the kill takes effect as soon as the other thread releases its lock.

Locking system tables

The thread is trying to lock a system table (for example, a time zone or log table).

• logging slow query

The thread is writing a statement to the slow-query log.

• login

The initial state for a connection thread until the client has been authenticated successfully.

manage keys

The server is enabling or disabling a table index.

• Opening system tables

The thread is trying to open a system table (for example, a time zone or log table).

• Opening tables

The thread is trying to open a table. This is should be very fast procedure, unless something prevents opening. For example, an ALTER TABLE or a LOCK TABLE statement can prevent opening a table until the statement is finished. It is also worth checking that your table\_open\_cache value is large enough.

For system tables, the Opening system tables state is used instead.

• optimizing

The server is performing initial optimizations for a query.

preparing

This state occurs during query optimization.

preparing for alter table

The server is preparing to execute an in-place ALTER TABLE.

Purging old relay logs

The thread is removing unneeded relay log files.

query end

This state occurs after processing a query but before the freeing items state.

• Receiving from client

The server is reading a packet from the client.

Removing duplicates

The query was using SELECT DISTINCT in such a way that MySQL could not optimize away the distinct operation at an early stage. Because of this, MySQL requires an extra stage to remove all duplicated rows before sending the result to the client.

removing tmp table

The thread is removing an internal temporary table after processing a SELECT statement. This state is not used if no temporary table was created.

• rename

The thread is renaming a table.

• rename result table

The thread is processing an ALTER TABLE statement, has created the new table, and is renaming it to replace the original table.

• Reopen tables

The thread got a lock for the table, but noticed after getting the lock that the underlying table structure changed. It has freed the lock, closed the table, and is trying to reopen it.

Repair by sorting

The repair code is using a sort to create indexes.

• Repair done

The thread has completed a multithreaded repair for a MyISAM table.

• Repair with keycache

The repair code is using creating keys one by one through the key cache. This is much slower than Repair by sorting.

• Rolling back

The thread is rolling back a transaction.

Saving state

For MyISAM table operations such as repair or analysis, the thread is saving the new table state to the .MYI file header. State includes information such as number of rows, the AUTO\_INCREMENT counter, and key distributions.

• Searching rows for update

The thread is doing a first phase to find all matching rows before updating them. This has to be done if the UPDATE is changing the index that is used to find the involved rows.

• Sending data

This state is now included in the Executing state.

Sending to client

The server is writing a packet to the client.

• setup

The thread is beginning an ALTER TABLE operation.

Sorting for group

The thread is doing a sort to satisfy a GROUP BY.

Sorting for order

The thread is doing a sort to satisfy an ORDER BY.

Sorting index

The thread is sorting index pages for more efficient access during a MyISAM table optimization operation.

• Sorting result

For a SELECT statement, this is similar to Creating sort index, but for nontemporary tables.

• starting

The first stage at the beginning of statement execution.

• statistics

The server is calculating statistics to develop a query execution plan. If a thread is in this state for a long time, the server is probably disk-bound performing other work.

• System lock

The thread has called mysql\_lock\_tables() and the thread state has not been updated since. This is a very general state that can occur for many reasons.

For example, the thread is going to request or is waiting for an internal or external system lock for the table. This can occur when InnoDB waits for a table-level lock during execution of LOCK TABLES. If this state is being caused by requests for external locks and you are not using multiple mysqld servers that are accessing the same MyISAM tables, you can disable external system locks with the --skip-external-locking option. However, external locking is disabled by default, so it is likely that this option has no effect. For SHOW PROFILE, this state means the thread is requesting the lock (not waiting for it).

For system tables, the Locking system tables state is used instead.

update

The thread is getting ready to start updating the table.

Updating

The thread is searching for rows to update and is updating them.

updating main table

The server is executing the first part of a multiple-table update. It is updating only the first table, and saving columns and offsets to be used for updating the other (reference) tables.

• updating reference tables

The server is executing the second part of a multiple-table update and updating the matched rows from the other tables.

User lock

The thread is going to request or is waiting for an advisory lock requested with a GET\_LOCK() call. For SHOW PROFILE, this state means the thread is requesting the lock (not waiting for it).

• User sleep

The thread has invoked a SLEEP() call.

• Waiting for commit lock

FLUSH TABLES WITH READ LOCK is waiting for a commit lock.

• waiting for handler commit

The thread is waiting for a transaction to commit versus other parts of query processing.

Waiting for tables

The thread got a notification that the underlying structure for a table has changed and it needs to reopen the table to get the new structure. However, to reopen the table, it must wait until all other threads have closed the table in question.

This notification takes place if another thread has used FLUSH TABLES or one of the following statements on the table in question: FLUSH TABLES tbl\_name, ALTER TABLE, RENAME TABLE, REPAIR TABLE, ANALYZE TABLE, OR OPTIMIZE TABLE.

• Waiting for table flush

The thread is executing FLUSH TABLES and is waiting for all threads to close their tables, or the thread got a notification that the underlying structure for a table has changed and it needs to reopen the table to get the new structure. However, to reopen the table, it must wait until all other threads have closed the table in question.

This notification takes place if another thread has used FLUSH TABLES or one of the following statements on the table in question: FLUSH TABLES tbl\_name, ALTER TABLE, RENAME TABLE, REPAIR TABLE, ANALYZE TABLE, Or OPTIMIZE TABLE.

• Waiting for *lock\_type* lock

The server is waiting to acquire a THR\_LOCK lock or a lock from the metadata locking subsystem, where <code>lock\_type</code> indicates the type of lock.

This state indicates a wait for a THR\_LOCK:

• Waiting for table level lock

These states indicate a wait for a metadata lock:

- Waiting for event metadata lock
- Waiting for global read lock
- Waiting for schema metadata lock
- Waiting for stored function metadata lock
- Waiting for stored procedure metadata lock
- Waiting for table metadata lock
- Waiting for trigger metadata lock

For information about table lock indicators, see Section 10.11.1, "Internal Locking Methods". For information about metadata locking, see Section 10.11.4, "Metadata Locking". To see which locks are blocking lock requests, use the Performance Schema lock tables described at Section 29.12.13, "Performance Schema Lock Tables".

Waiting on cond

A generic state in which the thread is waiting for a condition to become true. No specific state information is available.

Writing to net

The server is writing a packet to the network.

# 10.14.4 Replication Source Thread States

The following list shows the most common states you may see in the State column for the Binlog Dump thread of the replication source. If you see no Binlog Dump threads on a source, this means that replication is not running; that is, that no replicas are currently connected.

In MySQL 8.0, incompatible changes were made to instrumentation names. Monitoring tools that work with these instrumentation names might be impacted. If the incompatible changes have an impact for you, set the terminology\_use\_previous system variable to BEFORE\_8\_0\_26 to make MySQL Server use the old versions of the names for the objects specified in the previous list. This enables monitoring tools that rely on the old names to continue working until they can be updated to use the new names.

Set the terminology\_use\_previous system variable with session scope to support individual functions, or global scope to be a default for all new sessions. When global scope is used, the slow query log contains the old versions of the names.

Finished reading one binlog; switching to next binlog

The thread has finished reading a binary log file and is opening the next one to send to the replica.

Master has sent all binlog to slave; waiting for more updates
 Source has sent all binlog to replica; waiting for more updates

The thread has read all remaining updates from the binary logs and sent them to the replica. The thread is now idle, waiting for new events to appear in the binary log resulting from new updates occurring on the source.

Sending binlog event to slave
 Sending binlog event to replica

Binary logs consist of *events*, where an event is usually an update plus some other information. The thread has read an event from the binary log and is now sending it to the replica.

Waiting to finalize termination

A very brief state that occurs as the thread is stopping.

## 10.14.5 Replication I/O (Receiver) Thread States

The following list shows the most common states you see in the State column for a replication I/O (receiver) thread on a replica server. This state also appears in the Replica\_IO\_State column displayed by SHOW REPLICA STATUS, so you can get a good view of what is happening by using that statement.

In MySQL 8.0, incompatible changes were made to instrumentation names. Monitoring tools that work with these instrumentation names might be impacted. If the incompatible changes have an impact for you, set the terminology\_use\_previous system variable to BEFORE\_8\_0\_26 to make MySQL Server use the old versions of the names for the objects specified in the previous list. This enables monitoring tools that rely on the old names to continue working until they can be updated to use the new names.

Set the terminology\_use\_previous system variable with session scope to support individual functions, or global scope to be a default for all new sessions. When global scope is used, the slow query log contains the old versions of the names.

Checking master version
 Checking source version

A state that occurs very briefly, after the connection to the source is established.

Connecting to master
 Connecting to source

The thread is attempting to connect to the source.

Queueing master event to the relay log
 Queueing source event to the relay log

The thread has read an event and is copying it to the relay log so that the SQL thread can process it.

Reconnecting after a failed binlog dump request

The thread is trying to reconnect to the source.

Reconnecting after a failed master event read
 Reconnecting after a failed source event read

The thread is trying to reconnect to the source. When connection is established again, the state becomes Waiting for master to send event.

• Registering slave on master

```
Registering replica on source
```

A state that occurs very briefly after the connection to the source is established.

• Requesting binlog dump

A state that occurs very briefly, after the connection to the source is established. The thread sends to the source a request for the contents of its binary logs, starting from the requested binary log file name and position.

· Waiting for its turn to commit

A state that occurs when the replica thread is waiting for older worker threads to commit if replica preserve commit order is enabled.

· Waiting for master to send event

```
Waiting for source to send event
```

The thread has connected to the source and is waiting for binary log events to arrive. This can last for a long time if the source is idle. If the wait lasts for replica\_net\_timeout seconds, a timeout occurs. At that point, the thread considers the connection to be broken and makes an attempt to reconnect.

Waiting for master update

```
Waiting for source update
```

The initial state before Connecting to master or Connecting to source.

Waiting for slave mutex on exit

```
Waiting for replica mutex on exit
```

A state that occurs briefly as the thread is stopping.

• Waiting for the slave SQL thread to free enough relay log space

```
Waiting for the replica SQL thread to free enough relay log space
```

You are using a nonzero relay\_log\_space\_limit value, and the relay logs have grown large enough that their combined size exceeds this value. The I/O (receiver) thread is waiting until the SQL (applier) thread frees enough space by processing relay log contents so that it can delete some relay log files.

• Waiting to reconnect after a failed binlog dump request

If the binary log dump request failed (due to disconnection), the thread goes into this state while it sleeps, then tries to reconnect periodically. The interval between retries can be specified using the CHANGE REPLICATION SOURCE TO.

Waiting to reconnect after a failed master event read

```
Waiting to reconnect after a failed source event read
```

An error occurred while reading (due to disconnection). The thread is sleeping for the number of seconds set by the CHANGE REPLICATION SOURCE TO statement before attempting to reconnect.

## 10.14.6 Replication SQL Thread States

The following list shows the most common states you may see in the State column for a replication SQL thread on a replica server.

In MySQL 8.0, incompatible changes were made to instrumentation names. Monitoring tools that work with these instrumentation names might be impacted. If the incompatible changes have an impact for you, set the terminology\_use\_previous system variable to BEFORE\_8\_0\_26 to make MySQL Server use the old versions of the names for the objects specified in the previous list. This enables monitoring tools that rely on the old names to continue working until they can be updated to use the new names.

Set the terminology\_use\_previous system variable with session scope to support individual functions, or global scope to be a default for all new sessions. When global scope is used, the slow query log contains the old versions of the names.

• Making temporary file (append) before replaying LOAD DATA INFILE

The thread is executing a LOAD DATA statement and is appending the data to a temporary file containing the data from which the replica reads rows.

• Making temporary file (create) before replaying LOAD DATA INFILE

The thread is executing a LOAD DATA statement and is creating a temporary file containing the data from which the replica reads rows. This state can only be encountered if the original LOAD DATA statement was logged by a source running a version of MySQL lower than MySQL 5.0.3.

• Reading event from the relay log

The thread has read an event from the relay log so that the event can be processed.

Slave has read all relay log; waiting for more updates

```
Replica has read all relay log; waiting for more updates
```

The thread has processed all events in the relay log files, and is now waiting for the I/O (receiver) thread to write new events to the relay log.

Waiting for an event from Coordinator

Using the multithreaded replica (replica\_parallel\_workers is greater than 1), one of the replica worker threads is waiting for an event from the coordinator thread.

Waiting for slave mutex on exit

```
Waiting for replica mutex on exit
```

A very brief state that occurs as the thread is stopping.

Waiting for Slave Workers to free pending events

```
Waiting for Replica Workers to free pending events
```

This waiting action occurs when the total size of events being processed by Workers exceeds the size of the replica\_pending\_jobs\_size\_max system variable. The Coordinator resumes scheduling when the size drops below this limit. This state occurs only when replica\_parallel\_workers is set greater than 0.

Waiting for the next event in relay log

The initial state before Reading event from the relay log.

Waiting until SOURCE\_DELAY seconds after source executed event

The SQL thread has read an event but is waiting for the replica delay to lapse. This delay is set with the SOURCE\_DELAY option of the CHANGE REPLICATION SOURCE TO.

The Info column for the SQL thread may also show the text of a statement. This indicates that the thread has read an event from the relay log, extracted the statement from it, and may be executing it.

## 10.14.7 Replication Connection Thread States

These thread states occur on a replica server but are associated with connection threads, not with the I/O or SQL threads.

Changing master

Changing replication source

The thread is processing a CHANGE REPLICATION SOURCE TO statement.

• Killing slave

The thread is processing a STOP REPLICA statement.

• Opening master dump table

This state occurs after Creating table from master dump.

• Reading master dump table data

This state occurs after Opening master dump table.

Rebuilding the index on master dump table

This state occurs after Reading master dump table data.

### 10.14.8 NDB Cluster Thread States

- Committing events to binlog
- Opening mysql.ndb\_apply\_status
- Processing events

The thread is processing events for binary logging.

Processing events from schema table

The thread is doing the work of schema replication.

- Shutting down
- Syncing ndb table schema operation and binlog

This is used to have a correct binary log of schema operations for NDB.

Waiting for allowed to take ndbcluster global schema lock

The thread is waiting for permission to take a global schema lock.

• Waiting for event from ndbcluster

The server is acting as an SQL node in an NDB Cluster, and is connected to a cluster management node.

• Waiting for first event from ndbcluster

- Waiting for ndbcluster binlog update to reach current position
- Waiting for ndbcluster global schema lock

The thread is waiting for a global schema lock held by another thread to be released.

- Waiting for ndbcluster to start
- Waiting for schema epoch

The thread is waiting for a schema epoch (that is, a global checkpoint).

### 10.14.9 Event Scheduler Thread States

These states occur for the Event Scheduler thread, threads that are created to execute scheduled events, or threads that terminate the scheduler.

• Clearing

The scheduler thread or a thread that was executing an event is terminating and is about to end.

• Initialized

The scheduler thread or a thread that executes an event has been initialized.

• Waiting for next activation

The scheduler has a nonempty event queue but the next activation is in the future.

• Waiting for scheduler to stop

The thread issued SET GLOBAL event\_scheduler=OFF and is waiting for the scheduler to stop.

• Waiting on empty queue

The scheduler's event queue is empty and it is sleeping.

# 10.15 Tracing the Optimizer

The MySQL optimizer includes the capability to perform tracing; the interface is provided by a set of optimizer\_trace\_xxx system variables and the INFORMATION\_SCHEMA.OPTIMIZER\_TRACE table.

# 10.15.1 Typical Usage

To perform optimizer tracing entails the following steps:

- 1. Enable tracing by executing SET optimizer\_trace="enabled=ON".
- 2. Execute the statement to be traced. See Section 10.15.3, "Traceable Statements", for a listing of statements which can be traced.
- 3. Examine the contents of the  ${\tt INFORMATION\_SCHEMA.OPTIMIZER\_TRACE}$  table.
- 4. To examine traces for multiple queries, repeat the previous two steps as needed.
- 5. To disable tracing after you have finished, execute SET optimizer trace="enabled=OFF".

You can trace only statements which are executed within the current session; you cannot see traces from other sessions.

# 10.15.2 System Variables Controlling Tracing

The following system variables affect optimizer tracing:

- optimizer\_trace: Enables or disables optimizer tracing. See Section 10.15.8, "The optimizer\_trace System Variable".
- optimizer\_trace\_features: Enables or disables selected features of the MySQL Optimizer, using the syntax shown here:

```
SET optimizer_trace_features=option=value[,option=value][,...]
option:
    {greedy_search | range_optimizer | dynamic_range | repeated_subselect}}
value:
    {on | off | default}
```

See Section 10.15.10, "Selecting Optimizer Features to Trace", for more information on the effects of these.

- optimizer\_trace\_max\_mem\_size: Maximum amount of memory that can be used for storing all traces.
- optimizer\_trace\_limit: The maximum number of optimizer traces to be shown. See Section 10.15.4, "Tuning Trace Purging", for more information.
- optimizer\_trace\_offset: Offset of the first trace shown. See Section 10.15.4, "Tuning Trace Purging".
- end\_markers\_in\_json: If set to 1, causes the trace to repeat the key (if present) near the closing bracket. This also affects the output of EXPLAIN FORMAT=JSON in those versions of MySQL which support this statement. See Section 10.15.9, "The end\_markers\_in\_json System Variable".

## 10.15.3 Traceable Statements

Statements which are traceable are listed here:

- SELECT
- INSERT
- REPLACE
- UPDATE
- DELETE
- EXPLAIN with any of the preceding statements
- SET
- DO
- DECLARE, CASE, IF, and RETURN as used in stored routines
- CALL

Tracing is supported for both INSERT and REPLACE statements using VALUES, VALUES ROW, or SELECT.

Traces of multi-table UPDATE and DELETE statements are supported.

Tracing of SET optimizer\_trace is not supported.

For statements which are prepared and executed in separate steps, preparation and execution are traced separately.

## 10.15.4 Tuning Trace Purging

By default, each new trace overwrites the previous trace. Thus, if a statement contains substatements (such as invoking stored procedures, stored functions, or triggers), the topmost statement and substatements each generate one trace, but at the end of execution, the trace for only the last substatement is visible.

A user who wants to see the trace of a different substatement can enable or disable tracing for the desired substatement, but this requires editing the routine code, which may not always be possible. Another solution is to tune trace purging. This is done by setting the <code>optimizer\_trace\_offset</code> and <code>optimizer\_trace\_limit</code> system variables, like this:

```
SET optimizer_trace_offset=offset, optimizer_trace_limit=limit;
```

offset is a signed integer (default -1); limit is a positive integer (default 1). Such a SET statement has the following effects:

- All traces previously stored are cleared from memory.
- A subsequent SELECT from the OPTIMIZER\_TRACE table returns the first limit traces of the offset oldest stored traces (if offset >= 0), or the first limit traces of the -offset newest stored traces (if offset < 0).

#### Examples:

- SET optimizer\_trace\_offset=-1, optimizer\_trace\_limit=1: The most recent trace is shown (the default).
- SET optimizer\_trace\_offset=-2, optimizer\_trace\_limit=1: The next-to-last trace is shown.
- SET optimizer\_trace\_offset=-5, optimizer\_trace\_limit=5: The last five traces are shown.

Negative values for *offset* can thus prove useful when the substatements of interest are the last few in a stored routine. For example:

```
SET optimizer_trace_offset=-5, optimizer_trace_limit=5;

CALL stored_routine(); # more than 5 substatements in this routine

SELECT * FROM information_schema.OPTIMIZER_TRACE; # see only the last 5 traces
```

A positive *offset* can be useful when one knows that the interesting substatements are the first few in a stored routine.

The more accurately these two variables are set, the less memory is used. For example, SET optimizer\_trace\_offset=0, optimizer\_trace\_limit=5 requires sufficient memory to store five traces, so if only the three first are needed, is is better to use SET optimizer\_trace\_offset=0, optimizer\_trace\_limit=3, since tracing stops after limit traces. A stored routine may have a loop which executes many substatements and thus generates many traces, which can use a lot of memory; in such cases, choosing appropriate values for offset and limit can restrict tracing to, for example, a single iteration of the loop. This also decreases the impact of tracing on execution speed.

If offset is greater than or equal to 0, only <code>limit</code> traces are kept in memory. If offset is less than 0, that is not true: instead, <code>-offset</code> traces are kept in memory. Even if <code>limit</code> is smaller than <code>-offset</code>, excluding the last statement, the last statement must still be traced because it will be within the limit after executing one more statement. Since an offset less than 0 is counted from the end, the "window" moves as more statements execute.

Using optimizer\_trace\_offset and optimizer\_trace\_limit, which are restrictions at the trace producer level, provide better (greater) speed and (less) memory usage than setting offsets or

limits at the trace consumer (SQL) level with SELECT \* FROM OPTIMIZER\_TRACE LIMIT limit OFFSET offset, which saves almost nothing.

## 10.15.5 Tracing Memory Usage

Each stored trace is a string, which is extended (using realloc()) as optimization progresses by appending more data to it. The optimizer\_trace\_max\_mem\_size server system variable sets a limit on the total amount of memory used by all traces currently being stored. If this limit is reached, the current trace is not extended, which means the trace is incomplete; in this case the MISSING\_BYTES\_BEYOND\_MAX\_MEM\_SIZE column shows the number of bytes missing from the trace.

## 10.15.6 Privilege Checking

In complex scenarios where the query uses SQL SECURITY DEFINER views or stored routines, it may be that a user is denied from seeing the trace of its query because it lacks some extra privileges on those objects. In that case, the trace will be shown as empty and the INSUFFICIENT\_PRIVILEGES column will show "1".

## 10.15.7 Interaction with the --debug Option

Anything written to the trace is automatically written to the debug file.

## 10.15.8 The optimizer\_trace System Variable

The optimizer\_trace system variable has these on/off switches:

- enabled: Enables (ON) or disables (OFF) tracing
- one\_line: If set to ON, the trace contains no whitespace, thus conserving space. This renders the trace difficult to read for humans, still usable by JSON parsers, since they ignore whitespace.

# 10.15.9 The end\_markers\_in\_json System Variable

When reading a very large JSON document, it can be difficult to pair its closing bracket and opening brackets; setting end\_markers\_in\_json=ON repeats the structure's key, if it has one, near the closing bracket. This variable affects both optimizer traces and the output of EXPLAIN FORMAT=JSON.



#### Note

If end\_markers\_in\_json is enabled, the repetition of the key means the result is not a valid JSON document, and causes JSON parsers to throw an error.

# 10.15.10 Selecting Optimizer Features to Trace

Some features in the optimizer can be invoked many times during statement optimization and execution, and thus can make the trace grow beyond reason. They are:

- *Greedy search*: With an *N*-table join, this could explore factorial(*N*) plans.
- Range optimizer
- Dynamic range optimization: Shown as range checked for each record in EXPLAIN output; each outer row causes a re-run of the range optimizer.
- Subqueries: A subquery in which the WHERE clause may be executed once per row.

Those features can be excluded from tracing by setting one or more switches of the optimizer\_trace\_features system variable to OFF. These switches are listed here:

• greedy\_search: Greedy search is not traced.

- range\_optimizer: The range optimizer is not traced.
- dynamic\_range: Only the first call to the range optimizer on this JOIN\_TAB::SQL\_SELECT is traced.
- repeated\_subselect: Only the first execution of this Item\_subselect is traced.

### 10.15.11 Trace General Structure

A trace follows the actual execution path very closely; for each join, there is a join preparation object, a join optimization object, and a join execution object. Query transformations (IN to EXISTS, outer join to inner join, and so on), simplifications (elimination of clauses), and equality propagation are shown in subobjects. Calls to the range optimizer, cost evaluations, reasons why an access path is chosen over another one, or why a sorting method is chosen over another one, are shown as well.

## 10.15.12 Example

Here we take an example from the test suite.

```
# Tracing of ORDER BY & GROUP BY simplification.
SET optimizer_trace="enabled=on",end_markers_in_json=on; # make readable
SET optimizer_trace_max_mem_size=1000000; # avoid small default
CREATE TABLE t1 (
  pk INT, col_int_key INT,
  col_varchar_key VARCHAR(1);
  col_varchar_nokey VARCHAR(1)
INSERT INTO t1 VALUES
  (10,7,'v','v'),(11,0,'s','s'),(12,9,'l','l'),(13,3,'y','y'),(14,4,'c','c'),
  (15,2,'i','i'),(16,5,'h','h'),(17,3,'q','q'),(18,1,'a','a'),(19,3,'v','v'),
  (20,6,'u','u'),(21,7,'s','s'),(22,5,'y','y'),(23,1,'z','z'),(24,204,'h','h'),(25,224,'p','p'),(26,9,'e','e'),(27,5,'i','i'),(28,0,'y','y'),(29,3,'w','w');
CREATE TABLE t2 (
  pk INT, col_int_key INT,
  col_varchar_key VARCHAR(1),
  col_varchar_nokey VARCHAR(1),
  PRIMARY KEY (pk)
INSERT INTO t.2 VALUES
  (1,4,'b','b'),(2,8,'y','y'),(3,0,'p','p'),(4,0,'f','f'),(5,0,'p','p'),
  (6,7,'d','d'),(7,7,'f','f'),(8,5,'j','j'),(9,3,'e','e'),(10,188,'u','u'), (11,4,'v','v'),(12,9,'u','u'),(13,6,'i','i'),(14,1,'x','x'),(15,5,'l','l')
  (16,6,'q','q'),(17,2,'n','n'),(18,4,'r','r'),(19,231,'c','c'),(20,4,'h','h'),
  (21,3,'k','k'),(22,3,'t','t'),(23,7,'t','t'),(24,6,'k','k'),(25,7,'g','g'),
  (26,9,'z','z'),(27,4,'n','n'),(28,4,'j','j'),(29,2,'l','l'),(30,1,'d','d'),
   (31,2,'t','t'), (32,194,'y','y'), (33,2,'i','i'), (34,3,'j','j'), (35,8,'r','r'), \\
  (36,4,'b','b'),(37,9,'o','o'),(38,4,'k','k'),(39,5,'a','a'),(40,5,'f','f'),
  (41,9,'t','t'),(42,3,'c','c'),(43,8,'c','c'),(44,0,'r','r'),(45,98,'k','k'),
  (46,3,'1','1'),(47,1,'0','0'),(48,0,'t','t'),(49,189,'v','v'),(50,8,'x','x'),
  (51,3,'j','j'),(52,3,'x','x'),(53,9,'k','k'),(54,6,'o','o'),(55,8,'z','z'),
  (56,3,'n','n'),(57,9,'c','c'),(58,5,'d','d'),(59,9,'s','s'),(60,2,'j','j'),
  (61,2,'w','w'),(62,5,'f','f'),(63,8,'p','p'),(64,6,'o','o'),(65,9,'f','f'),
  (66,0,'x','x'),(67,3,'q','q'),(68,6,'g','g'),(69,5,'x','x'),(70,8,'p','p'),
  (71,2,'q','q'),(72,120,'q','q'),(73,25,'v','v'),(74,1,'g','g'),(75,3,'l','l'),
  (76,1,'w','w'),(77,3,'h','h'),(78,153,'c','c'),(79,5,'o','o'),(80,9,'o','o'),
   (81,1,'v','v'), (82,8,'y','y'), (83,7,'d','d'), (84,6,'p','p'), (85,2,'z','z'), \\
   (86,4,'t','t'), (87,7,'b','b'), (88,3,'y','y'), (89,8,'k','k'), (90,4,'c','c'), \\
  (91,6,'z','z'),(92,1,'t','t'),(93,7,'o','o'),(94,1,'u','u'),(95,0,'t','t'),
  (96,2,'k','k'),(97,7,'u','u'),(98,2,'b','b'),(99,1,'m','m'),(100,5,'o','o');
SELECT SUM(alias2.col_varchar_nokey) AS c1, alias2.pk AS c2
  FROM t1 AS alias1
  STRAIGHT_JOIN t2 AS alias2
```

```
ON alias2.pk = alias1.col_int_key
WHERE alias1.pk
GROUP BY c2
ORDER BY alias1.col_int_key, alias2.pk;
c1
   | c2 |
   0 1
       1
   0
        2
   0 |
       3
   0
        4
   0
        5
   0
        6
   0 | 9 |
```



#### Note

For reference, the complete trace is shown uninterrupted at the end of this section.

Now we can examine the trace, whose first column (QUERY), containing the original statement to be traced, is shown here:

This can be useful mark when several traces are stored.

The TRACE column begins by showing that execution of the statement is made up of discrete steps, like this:

```
"steps": [ {
```

This is followed by the preparation of the join for the first (and only) SELECT in the statement being traced, as shown here:

The output just shown displays the query as it is used for preparing the join; all columns (fields) have been resolved to their databases and tables, and each SELECT is annotated with a sequence number, which can be useful when studying subqueries.

The next portion of the trace shows how the join is optimized, starting with condition processing:

```
{
    "join_optimization": {
```

```
"select#": 1,
  "steps": [
     {
       "condition_processing": {
          "condition": "WHERE",
          "original_condition": "(`test`.`aliasl`.`pk` and \
(`test`.`alias2`.`pk` = `test`.`aliasl`.`col_int_key`))",
          "steps": [
               "transformation": "equality_propagation",
               "resulting_condition": "(`test`.`alias1`.`pk` and \
              multiple equal(`test`.`alias2`.`pk`, \
               `test`.`alias1`.`col_int_key`))"
               "transformation": "constant_propagation",
              "resulting_condition": "(`test`.`alias1`.`pk` and \ multiple equal(`test`.`alias2`.`pk`, \
               `test`.`alias1`.`col_int_key`))"
               "transformation": "trivial_condition_removal",
              "resulting_condition": "(`test`.`alias1`.`pk` and \
multiple equal(`test`.`alias2`.`pk`, \
               `test`.`alias1`.`col_int_key`))"
  ] /* steps */
} /* condition_processing */
```

Next, the optimizer checks for possible ref accesses, and identifies one:

A ref access which rejects NULL has been identified: no NULL in test.alias1.col\_int\_key can have a match. (Observe that it could have a match, were the operator a null-safe equals <=>).

Next, for every table in the query, we estimate the cost of, and number of records returned by, a table scan or a range access.

We need to find an optimal order for the tables. Normally, greedy search is used, but since the statement uses a straight join, only the requested order is explored, and one or more access methods are selected. As shown in this portion of the trace, we can choose a table scan:

```
"database": "test",
    "table": "alias2",
    "const_keys_added": {
      "keys": [
        "PRIMARY"
      ] /* keys */,
      "cause": "group_by"
    } /* const_keys_added */,
    range_analysis": {
      "table_scan": {
        "records": 100,
        "cost": 24.588
      } /* table_scan */,
      "potential_range_indices": [
          "index": "PRIMARY",
          "usable": true,
          "key_parts": [
            "pk"
          ] /* key_parts */
      ] /* potential_range_indices */,
      "setup_range_conditions": [
      ] /* setup_range_conditions */,
      "group_index_range": {
        "chosen": false,
        "cause": "not_single_table"
         } /* group_index_range */
    /* range_analysis */
] /* records_estimation */
```

As just shown in the second portion of the range analysis, it is not possible to use <code>GROUP\_MIN\_MAX</code> because it accepts only one table, and we have two in the join. This means that no range access is possible.

The optimizer estimates that reading the first table, and applying any required conditions to it, yields 20 rows:

For alias2, we choose ref access on the primary key rather than a table scan, because the number of records returned by the latter (75) is far greater than that returned by ref access (1), as shown here:

```
"records": 1,
    "cost": 20.2,
    "chosen": true

},

{
    "access_type": "scan",
    "using_join_cache": true,
    "records": 75,
    "cost": 7.4917,
    "chosen": false
    }

    ] /* considered_access_paths */
    */ best_access_path */,
    "cost_for_plan": 30.098,
    "records_for_plan": 20,
    "chosen": true
    }

] /* rest_of_plan */
}

] /* considered_execution_plans */
},
```

Now that the order of tables is fixed, we can split the WHERE condition into chunks which can be tested early (pushdown of conditions down the join tree):

This condition can be tested on rows of alias1 without reading rows from alias2.

```
{
    "database": "test",
    "table": "alias2",
    "attached": null
}

] /* attached_conditions_summary */
} /* attaching_conditions_to_tables */
},
{
```

Now we try to simplify the ORDER BY:

Because the WHERE clause contains alias2.pk=alias1.col\_int\_key, ordering by both columns is unnecessary; we can order by the first column alone, since the second column is always equal to it.

```
"resulting_clause_is_simple": true,
    "resulting_clause": "`test`.`alias1`.`col_int_key`"
} /* clause_processing */
```

},

The shorter ORDER BY clause (which is not visible in in the output of EXPLAIN) can be implemented as an index scan, since it uses only a single column of one table.

```
"clause_processing": {
        "clause": "GROUP BY",
        "original_clause": "`test`.`alias2`.`pk`",
        "items": [
            "item": "`test`.`alias2`.`pk`"
        ] /* items */,
        "resulting_clause_is_simple": false,
        "resulting_clause": "`test`.`alias2`.`pk`"
      } /* clause_processing */
      "refine_plan": [
        {
          "database": "test",
          "table": "alias1",
          "scan_type": "table"
          "database": "test",
          "table": "alias2"
      ] /* refine_plan */
  ] /* steps */
} /* join_optimization */
```

Now the join is executed:

```
"join_execution": {
    "select#": 1,
    "steps": [
    ] /* steps */
    } /* join_execution */
    }
] /* steps */
} 0 0
```

All traces have the same basic structure. If a statement uses subqueries, there can be multiple preparations, optimizations, and executions, as well as subquery-specific transformations.

The complete trace is shown here:

```
mysql> SELECT * FROM INFORMATION_SCHEMA.OPTIMIZER_TRACE\G
    ******************* 1. row *************
                           QUERY: SELECT SUM(alias2.col_varchar_nokey) AS c1, alias2.pk AS c2
  FROM t1 AS alias1
  STRAIGHT_JOIN t2 AS alias2
  ON alias2.pk = alias1.col_int_key
  WHERE alias1.pk
  GROUP BY c2
  ORDER BY alias1.col_int_key, alias2.pk
                           TRACE: {
  "steps": [
    {
      "join_preparation": {
        "select#": 1,
        "steps": [
         {
            "expanded_query": "/* select#1 */ select sum(`alias2`.`col_varchar_nokey`) AS `c1`,`alias2`.`pl
           "transformations_to_nested_joins": {
```

```
"transformations": [
          "JOIN_condition_to_WHERE",
          "parenthesis_removal"
        ] /* transformations */,
        "expanded_query": "/* select#1 */ select sum(`alias2`.`col_varchar_nokey`) AS `c1`,`alias
      } /* transformations_to_nested_joins */
      "functional_dependencies_of_GROUP_columns": {
        "all_columns_of_table_map_bits": [
        ] /* all_columns_of_table_map_bits */,
        "columns": [
          "test.alias2.pk",
          "test.alias1.col_int_key"
        ] /* columns */
      } /* functional_dependencies_of_GROUP_columns */
 ] /* steps */
} /* join_preparation */
"join_optimization": {
  "select#": 1,
  "steps": [
      "condition_processing": {
        "condition": "WHERE",
        "original_condition": "((0 <> `alias1`.`pk`) and (`alias2`.`pk` = `alias1`.`col_int_key`)
        "steps": [
          {
            "transformation": "equality_propagation",
            "resulting_condition": "((0 <> `alias1`.`pk`) and multiple equal(`alias2`.`pk`, `alia
            "transformation": "constant_propagation",
            "resulting_condition": "((0 <> `alias1`.`pk`) and multiple equal(`alias2`.`pk`, `alia
            "transformation": "trivial_condition_removal",
            "resulting_condition": "((0 <> `alias1`.`pk`) and multiple equal(`alias2`.`pk`, `alia
        ] /* steps */
      } /* condition_processing */
      "substitute_generated_columns": {
      } /* substitute_generated_columns */
      "table_dependencies": [
        {
          "table": "`tl` `alias1`",
          "row_may_be_null": false,
          "map_bit": 0,
          "depends_on_map_bits": [
          ] /* depends_on_map_bits */
          "table": "`t2` `alias2`",
          "row_may_be_null": false,
          "map_bit": 1,
          "depends_on_map_bits": [
          ] /* depends_on_map_bits */
      ] /* table_dependencies */
      "ref_optimizer_key_uses": [
          "table": "`t2` `alias2`",
```

```
"field": "pk",
    "equals": "`alias1`.`col_int_key`",
    "null_rejecting": true
] /* ref_optimizer_key_uses */
"rows_estimation": [
    "table": "`t1` `alias1`",
    "table_scan": {
      "rows": 20,
      "cost": 0.25
    } /* table_scan */
    "table": "`t2` `alias2`",
    "const_keys_added": {
     "keys": [
       "PRIMARY"
      ] /* keys */,
      "cause": "group_by"
    } /* const_keys_added */,
    "range_analysis": {
      "table_scan": {
        "rows": 100,
        "cost": 12.35
      } /* table_scan */,
      "potential_range_indexes": [
          "index": "PRIMARY",
          "usable": true,
          "key_parts": [
            "pk"
          ] /* key_parts */
      ] /* potential_range_indexes */,
      "setup_range_conditions": [
      ] /* setup_range_conditions */,
      "group_index_skip_scan": {
        "chosen": false,
"cause": "not_single_table"
      } /* group_index_skip_scan */,
      "skip_scan_range": {
        "chosen": false,
        "cause": "not_single_table"
      } /* skip_scan_range */
    } /* range_analysis */
] /* rows_estimation */
"considered_execution_plans": [
 {
    "plan_prefix": [
    ] /* plan_prefix */,
    "table": "`t1` `alias1`",
    "best_access_path": {
      "considered_access_paths": [
          "rows_to_scan": 20,
          "filtering_effect": [
          ] /* filtering_effect */,
          "final_filtering_effect": 0.9,
          "access_type": "scan",
          "resulting_rows": 18,
          "cost": 2.25,
          "chosen": true
      ] /* considered_access_paths */
    } /* best_access_path */,
    "condition_filtering_pct": 100,
```

```
"rows_for_plan": 18,
    "cost_for_plan": 2.25,
    "rest_of_plan": [
        "plan_prefix": [
          "`t1` `alias1`"
        ] /* plan_prefix */,
        "table": "`t2` `alias2`",
        "best_access_path": {
          "considered_access_paths": [
              "access_type": "eq_ref",
              "index": "PRIMARY",
              "rows": 1,
              "cost": 6.3,
              "chosen": true,
              "cause": "clustered_pk_chosen_by_heuristics"
              "rows_to_scan": 100,
              "filtering_effect": [
              ] /* filtering_effect */,
              "final_filtering_effect": 1,
              "access_type": "scan",
              "using_join_cache": true,
              "buffers_needed": 1,
              "resulting_rows": 100,
              "cost": 180.25,
              "chosen": false
          ] /* considered_access_paths */
        } /* best_access_path */,
        "condition_filtering_pct": 100,
        "rows_for_plan": 18,
        "cost_for_plan": 8.55,
        "chosen": true
    ] /* rest_of_plan */
] /* considered_execution_plans */
"attaching_conditions_to_tables": {
  "original_condition": "((`alias2`.`pk` = `alias1`.`col_int_key`) and (0 <> `alias1`.`pk`)
  "attached_conditions_computation": [
 ] /* attached_conditions_computation */,
  "attached_conditions_summary": [
      "table": "`t1` `alias1`",
      "attached": "((0 <> `alias1`.`pk`) and (`alias1`.`col_int_key` is not null))"
      "table": "`t2` `alias2`",
      "attached": "(`alias2`.`pk` = `alias1`.`col_int_key`)"
 ] /* attached_conditions_summary */
} /* attaching_conditions_to_tables */
"optimizing_distinct_group_by_order_by": {
  "simplifying_order_by": {
    "original_clause": "`alias1`.`col_int_key`,`alias2`.`pk`",
    "items": [
        "item": "`alias1`.`col_int_key`"
        "item": "`alias2`.`pk`",
        "eq_ref_to_preceding_items": true
    ] /* items */,
    "resulting_clause_is_simple": true,
```

```
"resulting_clause": "`alias1`.`col_int_key`"
        } /* simplifying_order_by */,
        "simplifying_group_by": {
          "original_clause": "`c2`",
          "items": [
              "item": "`alias2`.`pk`"
          ] /* items */,
          "resulting_clause_is_simple": false,
          "resulting_clause": "`c2`'
        } /* simplifying_group_by */
      } /* optimizing_distinct_group_by_order_by */
      "finalizing_table_conditions": [
          "table": "`t1` `alias1`",
          "original_table_condition": "((0 <> `alias1`.`pk`) and (`alias1`.`col_int_key` is not null
          "final_table_condition ": "((0 <> `aliasl`.`pk`) and (`aliasl`.`col_int_key` is not null
          "table": "`t2` `alias2`",
          "original_table_condition": "(`alias2`.`pk` = `alias1`.`col_int_key`)",
          "final_table_condition ": null
      ] /* finalizing_table_conditions */
      "refine_plan": [
          "table": "`t1` `alias1`"
          "table": "`t2` `alias2`"
      ] /* refine_plan */
      "considering_tmp_tables": [
          "adding_tmp_table_in_plan_at_position": 2,
          "write_method": "continuously_update_group_row"
         "adding_sort_to_table": ""
        } /* filesort */
     ] /* considering_tmp_tables */
  ] /* steps */
} /* join_optimization */
"join_execution": {
  "select#": 1,
  "steps": [
      "temp_table_aggregate": {
        "select#": 1,
        "steps": [
            "creating_tmp_table": {
              "tmp_table_info": {
                "table": "<temporary>",
                "in_plan_at_position": 2,
                "columns": 3,
                "row_length": 18,
                "key_length": 4,
                "unique_constraint": false,
                "makes_grouped_rows": true,
                "cannot_insert_duplicates": false,
                "location": "TempTable"
```

```
} /* tmp_table_info */
                  } /* creating_tmp_table */
              ] /* steps */
             /* temp_table_aggregate */
            "sorting_table": "<temporary>",
            "filesort_information": [
                "direction": "asc",
                "expression": "`alias1`.`col_int_key`"
            ] /* filesort_information */,
            "filesort_priority_queue_optimization": {
              "usable": false,
              "cause": "not applicable (no LIMIT)"
            } /* filesort_priority_queue_optimization */,
            "filesort_execution": [
            ] /* filesort_execution */,
            "filesort_summary": {
              "memory_available": 262144,
              "key_size": 9,
              "row_size": 26,
              "max_rows_per_buffer": 7710,
              "num_rows_estimate": 18446744073709551615,
              "num_rows_found": 8,
              "num_initial_chunks_spilled_to_disk": 0,
              "peak_memory_used": 32832,
              "sort_algorithm": "std::sort",
              "unpacked_addon_fields": "skip_heuristic",
              "sort_mode": "<fixed_sort_key, additional_fields>"
            } /* filesort_summary */
        ] /* steps */
      } /* join_execution */
  ] /* steps */
MISSING_BYTES_BEYOND_MAX_MEM_SIZE: 0
          INSUFFICIENT_PRIVILEGES: 0
```

# 10.15.13 Displaying Traces in Other Applications

Examining a trace in the mysql command-line client can be made less difficult using the pager less command (or your operating platform's equivalent). An alternative can be to send the trace to a file, similarly to what is shown here:

```
SELECT TRACE INTO DUMPFILE file
FROM INFORMATION_SCHEMA.OPTIMIZER_TRACE;
```

You can then pass this file to a JSON-aware text editor or other viewer, such as the JsonView add-on for Firefox and Chrome, which shows objects in color and allows objects to be expanded or collapsed.

INTO DUMPFILE is preferable to INTO OUTFILE for this purpose, since the latter escapes newlines. As noted previously, you should ensure that end\_markers\_in\_json is OFFwhen executing the SELECT INTO statement, so that the output is valid JSON.

# 10.15.14 Preventing the Use of Optimizer Trace

If, for some reason, you wish to prevent users from seeing traces of their queries, start the server with the options shown here:

```
--maximum-optimizer-trace-max-mem-size=0 --optimizer-trace-max-mem-size=0
```

This sets the maximum size to 0 and prevents users from changing this limit, thus truncating all traces to 0 bytes.

# 10.15.15 Testing Optimizer Trace

This feature is tested in mysql-test/suite/opt\_trace and unittest/gunit/opt\_trace-t.

# **10.15.16 Optimizer Trace Implementation**

See the files  $sql/opt\_trace*$ , starting with  $sql/opt\_trace.h$ . A trace is started by creating an instance of  $Opt\_trace\_start$ ; information is added to this trace by creating instances of  $Opt\_trace\_object$  and  $Opt\_trace\_array$ , and by using the add() methods of these classes.