# MySQL Internationalization and Localization

# MySQL Internationalization and Localization

## Abstract

This is the MySQL Internationalization and Localization extract from the MySQL 6.0 Reference Manual.

Document generated on: 2009-06-02 (revision: 15165)

For more information on the terms of this license, for details on how the MySQL documentation is built and produced, or if you are interested in doing a translation, please contact the Documentation Team.

For additional licensing information, including licenses for libraries used by MySQL, see Preface, Notes, Licenses.

If you want help with using MySQL, please visit either the MySQL Forums or MySQL Mailing Lists where you can discuss your issues with other MySQL users.

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

# Internationalization and Localization

This chapter covers issues of internationalization (MySQL's capabilities for adapting to local use) and localization (selecting particular local conventions):

- MySQL support for character sets in SQL statements.

- How to configure the server to support different character sets.

- Selecting the language for error messages.

- How to set the server's time zone and enable per-connection time zone support.

- Selecting the locale for day and month names.

# Chapter 1. Character Set Support

MySQL includes character set support that enables you to store data using a variety of character sets and perform comparisons according to a variety of collations. You can specify character sets at the server, database, table, and column level. MySQL supports the use of character sets for the `MyISAM`, `MEMORY`, and `InnoDB` storage engines.

This chapter discusses the following topics:

- What are character sets and collations?

- The multiple-level default system for character set assignment

- Syntax for specifying character sets and collations

- Affected functions and operations

- Unicode support

- The character sets and collations that are available, with notes

Character set issues affect not only data storage, but also communication between client programs and the MySQL server. If you want the client program to communicate with the server using a character set different from the default, you'll need to indicate which one. For example, to use the `utf8` Unicode character set, issue this statement after connecting to the server:

```
SET NAMES 'utf8';
```

For more information about configuring character sets for application use and character set-related issues in client/server communication, see Section 1.5, "Configuring the Character Set and Collation for Applications", and Section 1.4, "Connection Character Sets and Collations".

## 1.1. Character Sets and Collations in General

A *character set* is a set of symbols and encodings. A *collation* is a set of rules for comparing characters in a character set. Let's make the distinction clear with an example of an imaginary character set.

Suppose that we have an alphabet with four letters: "A", "B", "a", "b". We give each letter a number: "A" = 0, "B" = 1, "a" = 2, "b" = 3. The letter "A" is a symbol, the number 0 is the **encoding** for "A", and the combination of all four letters and their encodings is a **character set**.

Suppose that we want to compare two string values, "A" and "B". The simplest way to do this is to look at the encodings: 0 for "A" and 1 for "B". Because 0 is less than 1, we say "A" is less than "B". What we've just done is apply a collation to our character set. The collation is a set of rules (only one rule in this case): "compare the encodings." We call this simplest of all possible collations a *binary* collation.

But what if we want to say that the lowercase and uppercase letters are equivalent? Then we would have at least two rules: (1) treat the lowercase letters "a" and "b" as equivalent to "A" and "B"; (2) then compare the encodings. We call this a *case-insensitive* collation. It is a little more complex than a binary collation.

In real life, most character sets have many characters: not just "A" and "B" but whole alphabets, sometimes multiple alphabets or eastern writing systems with thousands of characters, along with many special symbols and punctuation marks. Also in real life, most collations have many rules, not just for whether to distinguish lettercase, but also for whether to distinguish accents (an "accent" is a mark attached to a character as in German "Ã#"), and for multiple-character mappings (such as the rule that "Ã#" = "OE" in one of the two German collations).

MySQL can do these things for you:

- Store strings using a variety of character sets

- Compare strings using a variety of collations

- Mix strings with different character sets or collations in the same server, the same database, or even the same table

- Allow specification of character set and collation at any level

In these respects, MySQL is far ahead of most other database management systems. However, to use these features effectively, you

need to know what character sets and collations are available, how to change the defaults, and how they affect the behavior of string operators and functions.

# 1.2. Character Sets and Collations in MySQL

The MySQL server can support multiple character sets. To list the available character sets, use the SHOW CHARACTER SET statement. A partial listing follows. For more complete information, see Section 1.13, "Character Sets and Collations That MySQL Supports".

```
mysql> SHOW CHARACTER SET;
+----------+-----------------------------+---------------------+--------+
| Charset  | Description                 | Default collation   | Maxlen |
+----------+-----------------------------+---------------------+--------+
| big5     | Big5 Traditional Chinese    | big5_chinese_ci     |      2 |
| dec8     | DEC West European           | dec8_swedish_ci     |      1 |
| cp850    | DOS West European           | cp850_general_ci    |      1 |
| hp8      | HP West European            | hp8_english_ci      |      1 |
| koi8r    | KOI8-R Relcom Russian       | koi8r_general_ci    |      1 |
| latin1   | cp1252 West European        | latin1_swedish_ci   |      1 |
| latin2   | ISO 8859-2 Central European | latin2_general_ci   |      1 |
| swe7     | 7bit Swedish                | swe7_swedish_ci     |      1 |
| ascii    | US ASCII                    | ascii_general_ci    |      1 |
| ujis     | EUC-JP Japanese             | ujis_japanese_ci    |      3 |
| sjis     | Shift-JIS Japanese          | sjis_japanese_ci    |      2 |
| hebrew   | ISO 8859-8 Hebrew           | hebrew_general_ci   |      1 |
| tis620   | TIS620 Thai                 | tis620_thai_ci      |      1 |
| euckr    | EUC-KR Korean               | euckr_korean_ci     |      2 |
| koi8u    | KOI8-U Ukrainian            | koi8u_general_ci    |      1 |
| gb2312   | GB2312 Simplified Chinese   | gb2312_chinese_ci   |      2 |
| greek    | ISO 8859-7 Greek            | greek_general_ci    |      1 |
| cp1250   | Windows Central European    | cp1250_general_ci   |      1 |
| gbk      | GBK Simplified Chinese      | gbk_chinese_ci      |      2 |
| latin5   | ISO 8859-9 Turkish          | latin5_turkish_ci   |      1 |
...
```

Any given character set always has at least one collation. It may have several collations. To list the collations for a character set, use the SHOW COLLATION statement. For example, to see the collations for the latin1 (cp1252 West European) character set, use this statement to find those collation names that begin with latin1:

```
mysql> SHOW COLLATION LIKE 'latin1%';
+-------------------+---------+----+---------+----------+---------+
| Collation         | Charset | Id | Default | Compiled | Sortlen |
+-------------------+---------+----+---------+----------+---------+
| latin1_german1_ci | latin1  |  5 |         |          |       0 |
| latin1_swedish_ci | latin1  |  8 | Yes     | Yes      |       1 |
| latin1_danish_ci  | latin1  | 15 |         |          |       0 |
| latin1_german2_ci | latin1  | 31 |         | Yes      |       2 |
| latin1_bin        | latin1  | 47 |         | Yes      |       1 |
| latin1_general_ci | latin1  | 48 |         |          |       0 |
| latin1_general_cs | latin1  | 49 |         |          |       0 |
| latin1_spanish_ci | latin1  | 94 |         |          |       0 |
+-------------------+---------+----+---------+----------+---------+
```

The latin1 collations have the following meanings.

| Collation | Meaning |
| --- | --- |
| latin1_german1_ci | German DIN-1 |
| latin1_swedish_ci | Swedish/Finnish |
| latin1_danish_ci | Danish/Norwegian |
| latin1_german2_ci | German DIN-2 |
| latin1_bin | Binary according to latin1 encoding |
| latin1_general_ci | Multilingual (Western European) |
| latin1_general_cs | Multilingual (ISO Western European), case sensitive |
| latin1_spanish_ci | Modern Spanish |

Collations have these general characteristics:

• Two different character sets cannot have the same collation.

• Each character set has one collation that is the *default collation*. For example, the default collation for latin1 is latin1_swedish_ci. The output for SHOW CHARACTER SET indicates which collation is the default for each displayed character set.

- There is a convention for collation names: They start with the name of the character set with which they are associated, they usually include a language name, and they end with `_ci` (case insensitive), `_cs` (case sensitive), or `_bin` (binary).

In cases where a character set has multiple collations, it might not be clear which collation is most suitable for a given application. To avoid choosing the wrong collation, it can be helpful to perform some comparisons with representative data values to make sure that a given collation sorts values the way you expect.

Collation-Charts.Org is a useful site for information that shows how one collation compares to another.

# 1.3. Specifying Character Sets and Collations

There are default settings for character sets and collations at four levels: server, database, table, and column. The description in the following sections may appear complex, but it has been found in practice that multiple-level defaulting leads to natural and obvious results.

`CHARACTER SET` is used in clauses that specify a character set. `CHARSET` can be used as a synonym for `CHARACTER SET`.

Character set issues affect not only data storage, but also communication between client programs and the MySQL server. If you want the client program to communicate with the server using a character set different from the default, you'll need to indicate which one. For example, to use the `utf8` Unicode character set, issue this statement after connecting to the server:

```
SET NAMES 'utf8';
```

For more information about character set-related issues in client/server communication, see Section 1.4, "Connection Character Sets and Collations".

## 1.3.1. Server Character Set and Collation

MySQL Server has a server character set and a server collation. These can be set at server startup on the command line or in an option file and changed at runtime.

Initially, the server character set and collation depend on the options that you use when you start `mysqld`. You can use `--character-set-server` for the character set. Along with it, you can add `--collation-server` for the collation. If you don't specify a character set, that is the same as saying `--character-set-server=latin1`. If you specify only a character set (for example, `latin1`) but not a collation, that is the same as saying `--character-set-server=latin1 --collation-server=latin1_swedish_ci` because `latin1_swedish_ci` is the default collation for `latin1`. Therefore, the following three commands all have the same effect:

```
shell> mysqld
shell> mysqld --character-set-server=latin1
shell> mysqld --character-set-server=latin1 \
          --collation-server=latin1_swedish_ci
```

One way to change the settings is by recompiling. If you want to change the default server character set and collation when building from sources, use: `--with-charset` and `--with-collation` as arguments for `configure`. For example:

```
shell> ./configure --with-charset=latin1
```

Or:

```
shell> ./configure --with-charset=latin1 \
          --with-collation=latin1_german1_ci
```

Both `mysqld` and `configure` verify that the character set/collation combination is valid. If not, each program displays an error message and terminates.

The server character set and collation are used as default values if the database character set and collation are not specified in `CREATE DATABASE` statements. They have no other purpose.

The current server character set and collation can be determined from the values of the `character_set_server` and `collation_server` system variables. These variables can be changed at runtime.

## 1.3.2. Database Character Set and Collation

Every database has a database character set and a database collation. The `CREATE DATABASE` and `ALTER DATABASE` statements have optional clauses for specifying the database character set and collation:

```
CREATE DATABASE db_name
```

```
    [[DEFAULT] CHARACTER SET charset_name]
    [[DEFAULT] COLLATE collation_name]
ALTER DATABASE db_name
    [[DEFAULT] CHARACTER SET charset_name]
    [[DEFAULT] COLLATE collation_name]
```

The keyword `SCHEMA` can be used instead of `DATABASE`.

All database options are stored in a text file named `db.opt` that can be found in the database directory.

The `CHARACTER SET` and `COLLATE` clauses make it possible to create databases with different character sets and collations on the same MySQL server.

Example:

```
CREATE DATABASE db_name CHARACTER SET latin1 COLLATE latin1_swedish_ci;
```

MySQL chooses the database character set and database collation in the following manner:

- If both `CHARACTER SET X` and `COLLATE Y` are specified, character set `X` and collation `Y` are used.

- If `CHARACTER SET X` is specified without `COLLATE`, character set `X` and its default collation are used. To see the default collation for each character set, use the `SHOW COLLATION` statement.

- If `COLLATE Y` is specified without `CHARACTER SET`, the character set associated with `Y` and collation `Y` are used.

- Otherwise, the server character set and server collation are used.

The database character set and collation are used as default values for table definitions if the table character set and collation are not specified in `CREATE TABLE` statements. The database character set also is used by `LOAD DATA INFILE`. The character set and collation have no other purposes.

The character set and collation for the default database can be determined from the values of the `character_set_database` and `collation_database` system variables. The server sets these variables whenever the default database changes. If there is no default database, the variables have the same value as the corresponding server-level system variables, `character_set_server` and `collation_server`.

## 1.3.3. Table Character Set and Collation

Every table has a table character set and a table collation. The `CREATE TABLE` and `ALTER TABLE` statements have optional clauses for specifying the table character set and collation:

```
CREATE TABLE tbl_name (column_list)
    [[DEFAULT] CHARACTER SET charset_name]
    [COLLATE collation_name]]
ALTER TABLE tbl_name
    [[DEFAULT] CHARACTER SET charset_name]
    [COLLATE collation_name]
```

Example:

```
CREATE TABLE t1 ( ... )
CHARACTER SET latin1 COLLATE latin1_danish_ci;
```

MySQL chooses the table character set and collation in the following manner:

- If both `CHARACTER SET X` and `COLLATE Y` are specified, character set `X` and collation `Y` are used.

- If `CHARACTER SET X` is specified without `COLLATE`, character set `X` and its default collation are used. To see the default collation for each character set, use the `SHOW COLLATION` statement.

- If `COLLATE Y` is specified without `CHARACTER SET`, the character set associated with `Y` and collation `Y` are used.

- Otherwise, the database character set and collation are used.

The table character set and collation are used as default values for column definitions if the column character set and collation are not specified in individual column definitions. The table character set and collation are MySQL extensions; there are no such things in standard SQL.

## 1.3.4. Column Character Set and Collation

Every "character" column (that is, a column of type CHAR, VARCHAR, or TEXT) has a column character set and a column collation. Column definition syntax for CREATE TABLE and ALTER TABLE has optional clauses for specifying the column character set and collation:

```
col_name {CHAR | VARCHAR | TEXT} (col_length)
    [CHARACTER SET charset_name]
    [COLLATE collation_name]
```

These clauses can also be used for ENUM and SET columns:

```
col_name {ENUM | SET} (val_list)
    [CHARACTER SET charset_name]
    [COLLATE collation_name]
```

Examples:

```
CREATE TABLE t1
(
    col1 VARCHAR(5)
      CHARACTER SET latin1
      COLLATE latin1_german1_ci
);
ALTER TABLE t1 MODIFY
    col1 VARCHAR(5)
      CHARACTER SET latin1
      COLLATE latin1_swedish_ci;
```

MySQL chooses the column character set and collation in the following manner:

- If both CHARACTER SET X and COLLATE Y are specified, character set X and collation Y are used.

```
CREATE TABLE t1
(
    col1 CHAR(10) CHARACTER SET utf8 COLLATE utf8_unicode_ci
) CHARACTER SET latin1 COLLATE latin1_bin;
```

  The character set and collation are specified for the column, so they are used. The column has character set utf8 and collation utf8_unicode_ci.

- If CHARACTER SET X is specified without COLLATE, character set X and its default collation are used.

```
CREATE TABLE t1
(
    col1 CHAR(10) CHARACTER SET utf8
) CHARACTER SET latin1 COLLATE latin1_bin;
```

  The character set is specified for the column, but the collation is not. The column has character set utf8 and the default collation for utf8, which is utf8_general_ci. To see the default collation for each character set, use the SHOW COLLATION statement.

- If COLLATE Y is specified without CHARACTER SET, the character set associated with Y and collation Y are used.

```
CREATE TABLE t1
(
    col1 CHAR(10) COLLATE utf8_polish_ci
) CHARACTER SET latin1 COLLATE latin1_bin;
```

  The collation is specified for the column, but the character set is not. The column has collation utf8_polish_ci and the character set is the one associated with the collation, which is utf8.

- Otherwise, the table character set and collation are used.

```
CREATE TABLE t1
(
    col1 CHAR(10)
) CHARACTER SET latin1 COLLATE latin1_bin;
```

  Neither the character set nor collation are specified for the column, so the table defaults are used. The column has character set latin1 and collation latin1_bin.

The CHARACTER SET and COLLATE clauses are standard SQL.

If you use `ALTER TABLE` to convert a column from one character set to another, MySQL attempts to map the data values, but if the character sets are incompatible, there may be data loss.

# 1.3.5. Character String Literal Character Set and Collation

Every character string literal has a character set and a collation.

A character string literal may have an optional character set introducer and `COLLATE` clause:

```
[_charset_name]'string' [COLLATE collation_name]
```

Examples:

```
SELECT 'string';
SELECT _latin1'string';
SELECT _latin1'string' COLLATE latin1_danish_ci;
```

For the simple statement `SELECT 'string'`, the string has the character set and collation defined by the `character_set_connection` and `collation_connection` system variables.

The `_charset_name` expression is formally called an *introducer*. It tells the parser, "the string that is about to follow uses character set *X*." Because this has confused people in the past, we emphasize that an introducer does not change the string to the introducer character set like `CONVERT()` would do. It does not change the string's value, although padding may occur. The introducer is just a signal. An introducer is also legal before standard hex literal and numeric hex literal notation (`x'literal'` and `0xnnnn`), or before bit-field literal notation (`b'literal'` and `0bnnnn`).

Examples:

```
SELECT _latin1 x'AABBCC';
SELECT _latin1 0xAABBCC;
SELECT _latin1 b'1100011';
SELECT _latin1 0b1100011;
```

MySQL determines a literal's character set and collation in the following manner:

- If both `_X` and `COLLATE Y` are specified, character set *X* and collation *Y* are used.

- If `_X` is specified but `COLLATE` is not specified, character set *X* and its default collation are used. To see the default collation for each character set, use the `SHOW COLLATION` statement.

- Otherwise, the character set and collation given by the `character_set_connection` and `collation_connection` system variables are used.

Examples:

- A string with `latin1` character set and `latin1_german1_ci` collation:

  ```
  SELECT _latin1'Müller' COLLATE latin1_german1_ci;
  ```

- A string with `latin1` character set and its default collation (that is, `latin1_swedish_ci`):

  ```
  SELECT _latin1'Müller';
  ```

- A string with the connection default character set and collation:

  ```
  SELECT 'Müller';
  ```

Character set introducers and the `COLLATE` clause are implemented according to standard SQL specifications.

An introducer indicates the character set for the following string, but does not change now how the parser performs escape processing within the string. Escapes are always interpreted by the parser according to the character set given by `character_set_connection`.

The following examples show that escape processing occurs using `character_set_connection` even in the presence of an introducer. The examples use `SET NAMES` (which changes `character_set_connection`, as discussed in Section 1.4, "Connection Character Sets and Collations"), and display the resulting strings using the `HEX()` function so that the exact string

contents can be seen.

Example 1:

```
mysql> SET NAMES latin1;
Query OK, 0 rows affected (0.01 sec)
mysql> SELECT HEX('à\n'), HEX(_sjis'à\n');
+-----------+----------------+
| HEX('à\n') | HEX(_sjis'à\n') |
+-----------+----------------+
| E00A       | E00A           |
+-----------+----------------+
1 row in set (0.00 sec)
```

Here, "à" (hex value E0) is followed by "\n", the escape sequence for newline. The escape sequence is interpreted using the character_set_connection value of latin1 to produce a literal newline (hex value 0A). This happens even for the second string. That is, the introducer of _sjis does not affect the parser's escape processing.

Example 2:

```
mysql> SET NAMES sjis;
Query OK, 0 rows affected (0.00 sec)
mysql> SELECT HEX('à\n'), HEX(_latin1'à\n');
+-----------+------------------+
| HEX('à\n') | HEX(_latin1'à\n') |
+-----------+------------------+
| E05C6E     | E05C6E           |
+-----------+------------------+
1 row in set (0.04 sec)
```

Here, character_set_connection is sjis, a character set in which the sequence of "à" followed by "\" (hex values 05 and 5C) is a valid multi-byte character. Hence, the first two bytes of the string are interpreted as a single sjis character, and the "\" is not interpreted as an escape character. The following "n" (hex value 6E) is not interpreted as part of an escape sequence. This is true even for the second string; the introducer of _latin1 does not affect escape processing.

## 1.3.6. National Character Set

Standard SQL defines NCHAR or NATIONAL CHAR as a way to indicate that a CHAR column should use some predefined character set. MySQL 6.0 uses utf8 as this predefined character set. For example, these data type declarations are equivalent:

```
CHAR(10) CHARACTER SET utf8
NATIONAL CHARACTER(10)
NCHAR(10)
```

As are these:

```
VARCHAR(10) CHARACTER SET utf8
NATIONAL VARCHAR(10)
NCHAR VARCHAR(10)
NATIONAL CHARACTER VARYING(10)
NATIONAL CHAR VARYING(10)
```

You can use N'literal' (or n'literal') to create a string in the national character set. These statements are equivalent:

```
SELECT N'some text';
SELECT n'some text';
SELECT _utf8'some text';
```

For information on upgrading character sets to MySQL 6.0 from versions prior to 4.1, see the *MySQL 3.23, 4.0, 4.1 Reference Manual*.

## 1.3.7. Examples of Character Set and Collation Assignment

The following examples show how MySQL determines default character set and collation values.

**Example 1: Table and Column Definition**

```
CREATE TABLE t1
(
    c1 CHAR(10) CHARACTER SET latin1 COLLATE latin1_german1_ci
) DEFAULT CHARACTER SET latin2 COLLATE latin2_bin;
```

Here we have a column with a latin1 character set and a latin1_german1_ci collation. The definition is explicit, so that is straightforward. Notice that there is no problem with storing a latin1 column in a latin2 table.

**Example 2: Table and Column Definition**

```
CREATE TABLE t1
(
    c1 CHAR(10) CHARACTER SET latin1
) DEFAULT CHARACTER SET latin1 COLLATE latin1_danish_ci;
```

This time we have a column with a latin1 character set and a default collation. Although it might seem natural, the default collation is not taken from the table level. Instead, because the default collation for latin1 is always latin1_swedish_ci, column c1 has a collation of latin1_swedish_ci (not latin1_danish_ci).

**Example 3: Table and Column Definition**

```
CREATE TABLE t1
(
    c1 CHAR(10)
) DEFAULT CHARACTER SET latin1 COLLATE latin1_danish_ci;
```

We have a column with a default character set and a default collation. In this circumstance, MySQL checks the table level to determine the column character set and collation. Consequently, the character set for column c1 is latin1 and its collation is latin1_danish_ci.

**Example 4: Database, Table, and Column Definition**

```
CREATE DATABASE d1
    DEFAULT CHARACTER SET latin2 COLLATE latin2_czech_ci;
USE d1;
CREATE TABLE t1
(
    c1 CHAR(10)
);
```

We create a column without specifying its character set and collation. We're also not specifying a character set and a collation at the table level. In this circumstance, MySQL checks the database level to determine the table settings, which thereafter become the column settings.) Consequently, the character set for column c1 is latin2 and its collation is latin2_czech_ci.

## 1.3.8. Compatibility with Other DBMSs

For MaxDB compatibility these two statements are the same:

```
CREATE TABLE t1 (f1 CHAR(N) UNICODE);
CREATE TABLE t1 (f1 CHAR(N) CHARACTER SET ucs2);
```

# 1.4. Connection Character Sets and Collations

Several character set and collation system variables relate to a client's interaction with the server. Some of these have been mentioned in earlier sections:

- The server character set and collation can be determined from the values of the character_set_server and collation_server system variables.

- The character set and collation of the default database can be determined from the values of the character_set_database and collation_database system variables.

Additional character set and collation system variables are involved in handling traffic for the connection between a client and the server. Every client has connection-related character set and collation system variables.

Consider what a "connection" is: It is what you make when you connect to the server. The client sends SQL statements, such as queries, over the connection to the server. The server sends responses, such as result sets, over the connection back to the client. This leads to several questions about character set and collation handling for client connections, each of which can be answered in terms of system variables:

- What character set is the statement in when it leaves the client?

  The server takes the character_set_client system variable to be the character set in which statements are sent by the client.

- What character set should the server translate a statement to after receiving it?

For this, the server uses the `character_set_connection` and `collation_connection` system variables. It converts statements sent by the client from `character_set_client` to `character_set_connection` (except for string literals that have an introducer such as `_latin1` or `_utf8`). `collation_connection` is important for comparisons of literal strings. For comparisons of strings with column values, `collation_connection` does not matter because columns have their own collation, which has a higher collation precedence.

- What character set should the server translate to before shipping result sets or error messages back to the client?

    The `character_set_results` system variable indicates the character set in which the server returns query results to the client. This includes result data such as column values, and result metadata such as column names.

You can fine-tune the settings for these variables, or you can depend on the defaults (in which case, you can skip the rest of this section). If you do not use the defaults, you must change the character settings *for each connection to the server.*

There are two statements that affect the connection character sets:

```
SET NAMES 'charset_name'
SET CHARACTER SET charset_name
```

`SET NAMES` indicates what character set the client will use to send SQL statements to the server. Thus, `SET NAMES 'cp1251'` tells the server "future incoming messages from this client are in character set `cp1251`." It also specifies the character set that the server should use for sending results back to the client. (For example, it indicates what character set to use for column values if you use a `SELECT` statement.)

A `SET NAMES 'x'` statement is equivalent to these three statements:

```
SET character_set_client = x;
SET character_set_results = x;
SET character_set_connection = x;
```

Setting `character_set_connection` to *x* also sets `collation_connection` to the default collation for *x*. It is not necessary to set that collation explicitly. To specify a particular collation for the character sets, use the optional `COLLATE` clause:

```
SET NAMES 'charset_name' COLLATE 'collation_name'
```

`SET CHARACTER SET` is similar to `SET NAMES` but sets `character_set_connection` and `collation_connection` to `character_set_database` and `collation_database`. A `SET CHARACTER SET x` statement is equivalent to these three statements:

```
SET character_set_client = x;
SET character_set_results = x;
SET collation_connection = @@collation_database;
```

Setting `collation_connection` also sets `character_set_connection` to the character set associated with the collation (equivalent to executing `SET character_set_connection = @@character_set_database`). It is not necessary to set `character_set_connection` explicitly.

When a client connects, it sends to the server the name of the character set that it wants to use. The server uses the name to set the `character_set_client`, `character_set_results`, and `character_set_connection` system variables. In effect, the server performs a `SET NAMES` operation using the character set name.

With the `mysql` client, it is not necessary to execute `SET NAMES` every time you start up if you want to use a character set different from the default. You can add the `--default-character-set` option setting to your `mysql` statement line, or in your option file. For example, the following option file setting changes the three character set variables set to `koi8r` each time you invoke `mysql`:

```
[mysql]
default-character-set=koi8r
```

If you are using the `mysql` client with auto-reconnect enabled (which is not recommended), it is preferable to use the `charset` command rather than `SET NAMES`. For example:

```
mysql> charset utf8
Charset changed
```

The `charset` command issues a `SET NAMES` statement, and also changes the default character set that is used if `mysql` reconnects after the connection has dropped.

Example: Suppose that `column1` is defined as `CHAR(5) CHARACTER SET latin2`. If you do not say `SET NAMES` or `SET CHARACTER SET`, then for `SELECT column1 FROM t`, the server sends back all the values for `column1` using the character set that the client specified when it connected. On the other hand, if you say `SET NAMES 'latin1'` or `SET CHARACTER SET latin1` before issuing the `SELECT` statement, the server converts the `latin2` values to `latin1` just before sending results back. Conversion may be lossy if there are characters that are not in both character sets.

If you do not want the server to perform any conversion of result sets, set `character_set_results` to `NULL` or `binary`:

```
SET character_set_results = NULL;
```

> **Note**
>
> `ucs2`, `utf16`, and `utf32` cannot be used as a client character set, which means that they do not work for `SET NAMES` or `SET CHARACTER SET`.

To see the values of the character set and collation system variables that apply to your connection, use these statements:

```
SHOW VARIABLES LIKE 'character_set%';
SHOW VARIABLES LIKE 'collation%';
```

You must also consider the environment within which your MySQL applications execute. See Section 1.5, "Configuring the Character Set and Collation for Applications".

# 1.5. Configuring the Character Set and Collation for Applications

For applications that store data using the default MySQL character set and collation (`latin1`, `latin1_swedish_ci`), no special configuration should be needed. If applications require data storage using a different character set or collation, you can configure character set information several ways:

- Specify character settings per database. For example, applications that use one database might require `utf8`, whereas applications that use another database might require `sjis`.

- Specify character settings at server startup. This causes the server to use the given settings for all applications that do not make other arrangements.

- Specify character settings at configuration time, if you build MySQL from source. This causes the server to use the given settings for all applications, without having to specify them at server startup.

When different applications require different character settings, the per-database technique provides a good deal of flexibility. If most or all applications use the same character set, specifying character settings at server startup or configuration time may be most convenient.

For the per-database or server-startup techniques, the settings control the character set for data storage. Applications must also tell the server which character set to use for client/server communications, as described in the following instructions.

The examples shown here assume use of the `utf8` character set and `utf8_general_ci` collation.

**Specify character settings per database.** To create a database such that its tables will use a given default character set and collation for data storage, use a `CREATE DATABASE` statement like this:

```
CREATE DATABASE mydb
  DEFAULT CHARACTER SET utf8
  DEFAULT COLLATE utf8_general_ci;
```

Tables created in the database will use `utf8` and `utf8_general_ci` by default for any character columns.

Applications that use the database should also configure their connection to the server each time they connect. This can be done by executing a `SET NAMES 'utf8'` statement after connecting. The statement can be used regardless of connection method: The `mysql` client, PHP scripts, and so forth.

In some cases, it may be possible to configure the connection to use the desired character set some other way. For example, for connections made using `mysql`, you can specify the `--default-character-set=utf8` command-line option to achieve the same effect as `SET NAMES 'utf8'`.

For more information about configuring client connections, see Section 1.4, "Connection Character Sets and Collations".

**Specify character settings at server startup.** To select a character set and collation at server startup, use the `--character-set-server` and `--collation-server` options. For example, to specify the options in an option file, in-

clude these lines:

```
[mysqld]
character-set-server=utf8
collation-server=utf8_general_ci
```

These settings apply server-wide and apply as the defaults for databases created by any application, and for tables created in those databases.

It is still necessary for applications to configure their connection using SET NAMES or equivalent after they connect, as described previously. You might be tempted to start the server with the --init_connect="SET NAMES 'utf8'" option to cause SET NAMES to be executed automatically for each client that connects. However, this will yield inconsistent results because the init_connect value is not executed for users who have the SUPER privilege.

**Specify character settings at MySQL configuration time.** To select a character set and collation when you configure and build MySQL from source, use the --with-charset and --with-collation options:

```
shell> ./configure --with-charset=utf8 --with-collation=utf8_general_ci
```

The resulting server uses utf8 and utf8_general_ci as the default for databases and tables and for client connections. It is unnecessary to use --character-set-server and --collation-server at server startup. It is also unnecessary for applications to configure their connection using SET NAMES or equivalent after they connect to the server.

Regardless of how you configure the MySQL character set for application use, you must also consider the environment within which those applications execute. If you will send statements using UTF-8 text taken from a file that you create in an editor, you should edit the file with the locale of your environment set to UTF-8 so that the file's encoding is correct and so that the operating system handles it correctly. If you use the mysql client from within a terminal window, the window must be configured to use UTF-8 or characters may not display properly. For a script that executes in a Web environment, the script must handle character encoding properly for its interaction with the MySQL server, and it must generate pages that correctly indicate the encoding so that browsers know how to display the content of the pages. For example, you can include this <meta> tag within your <head> element:

```
<meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
```

# 1.6. Collation Issues

The following sections discuss various aspects of character set collations.

## 1.6.1. Using COLLATE in SQL Statements

With the COLLATE clause, you can override whatever the default collation is for a comparison. COLLATE may be used in various parts of SQL statements. Here are some examples:

- With ORDER BY:

```
SELECT k
FROM t1
ORDER BY k COLLATE latin1_german2_ci;
```

- With AS:

```
SELECT k COLLATE latin1_german2_ci AS k1
FROM t1
ORDER BY k1;
```

- With GROUP BY:

```
SELECT k
FROM t1
GROUP BY k COLLATE latin1_german2_ci;
```

- With aggregate functions:

```
SELECT MAX(k COLLATE latin1_german2_ci)
FROM t1;
```

- With DISTINCT:

```
SELECT DISTINCT k COLLATE latin1_german2_ci
FROM t1;
```

- With `WHERE`:

```
SELECT *
FROM t1
WHERE _latin1 'Müller' COLLATE latin1_german2_ci = k;
```

```
SELECT *
FROM t1
WHERE k LIKE _latin1 'Müller' COLLATE latin1_german2_ci;
```

- With `HAVING`:

```
SELECT k
FROM t1
GROUP BY k
HAVING k = _latin1 'Müller' COLLATE latin1_german2_ci;
```

## 1.6.2. `COLLATE` Clause Precedence

The `COLLATE` clause has high precedence (higher than `||`), so the following two expressions are equivalent:

```
x || y COLLATE z
x || (y COLLATE z)
```

## 1.6.3. `BINARY` Operator

The `BINARY` operator casts the string following it to a binary string. This is an easy way to force a comparison to be done byte by byte rather than character by character. `BINARY` also causes trailing spaces to be significant.

```
mysql> SELECT 'a' = 'A';
        -> 1
mysql> SELECT BINARY 'a' = 'A';
        -> 0
mysql> SELECT 'a' = 'a ';
        -> 1
mysql> SELECT BINARY 'a' = 'a ';
        -> 0
```

`BINARY str` is shorthand for `CAST(str AS BINARY)`.

The `BINARY` attribute in character column definitions has a different effect. A character column defined with the `BINARY` attribute is assigned the binary collation of the column's character set. Every character set has a binary collation. For example, the binary collation for the `latin1` character set is `latin1_bin`, so if the table default character set is `latin1`, these two column definitions are equivalent:

```
CHAR(10) BINARY
CHAR(10) CHARACTER SET latin1 COLLATE latin1_bin
```

The effect of `BINARY` as a column attribute differs from its effect prior to MySQL 4.1. Formerly, `BINARY` resulted in a column that was treated as a binary string. A binary string is a string of bytes that has no character set or collation, which differs from a nonbinary character string that has a binary collation. For both types of strings, comparisons are based on the numeric values of the string unit, but for nonbinary strings the unit is the character and some character sets allow multi-byte characters. The `BINARY` and `VARBINARY` Types.

The use of `CHARACTER SET binary` in the definition of a `CHAR`, `VARCHAR`, or `TEXT` column causes the column to be treated as a binary data type. For example, the following pairs of definitions are equivalent:

```
CHAR(10) CHARACTER SET binary
BINARY(10)
VARCHAR(10) CHARACTER SET binary
VARBINARY(10)
TEXT CHARACTER SET binary
BLOB
```

## 1.6.4. The `_bin` and `binary` Collations

This section describes how `_bin` collations for nonbinary strings differ from the `binary` "collation" for binary strings.

Nonbinary strings (as stored in the `CHAR`, `VARCHAR`, and `TEXT` data types) have a character set and collation. A given character set can have several collations, each of which defines a particular sorting and comparison order for the characters in the set. One of these is the binary collation for the character set, indicated by a `_bin` suffix in the collation name. For example, `latin1` and `utf8` have binary collations named `latin1_bin` and `utf8_bin`.

Binary strings (as stored in the `BINARY`, `VARBINARY`, and `BLOB` data types) have no character set or collation in the sense that nonbinary strings do. (Applied to a binary string, the `CHARSET()` and `COLLATION()` functions both return a value of `binary`.) Binary strings are sequences of bytes and the numeric values of those bytes determine sort order.

The `_bin` collations differ from the `binary` collation in several respects.

**The unit for sorting and comparison.** Binary strings are sequences of bytes. Sorting and comparison is always based on numeric byte values. Nonbinary strings are sequences of characters, which might be multi-byte. Collations for nonbinary strings define an ordering of the character values for sorting and comparison. For the `_bin` collation, this ordering is based solely on numeric values of the characters (which is similar to ordering for binary strings except that a `_bin` collation must take into account that a character might contain multiple bytes). For other collations, character ordering might take additional factors such as lettercase into account.

**Character set conversion.** A nonbinary string has a character set and is converted to another character set in many cases, even when the string has a `_bin` collation:

- When assigning column values from another column that has a different character set:

```
UPDATE t1 SET utf8_bin_column=latin1_column;
INSERT INTO t1 (latin1_column) SELECT utf8_bin_column FROM t2;
```

- When assigning column values for `INSERT` or `UPDATE` using a string literal:

```
SET NAMES latin1;
INSERT INTO t1 (utf8_bin_column) VALUES ('string-in-latin1');
```

- When sending results from the server to a client:

```
SET NAMES latin1;
SELECT utf8_bin_column FROM t2;
```

For binary string columns, no conversion occurs. For the preceding cases, the string value is copied byte-wise.

**Lettercase conversion.** Collations provide information about lettercase of characters, so characters in a nonbinary string can be converted from one lettercase to another, even for `_bin` collations that ignore lettercase for ordering:

```
mysql> SET NAMES latin1 COLLATE latin1_bin;
Query OK, 0 rows affected (0.02 sec)
mysql> SELECT LOWER('aA'), UPPER('zZ');
+-------------+-------------+
| LOWER('aA') | UPPER('zZ') |
+-------------+-------------+
| aa          | ZZ          |
+-------------+-------------+
1 row in set (0.13 sec)
```

The concept of lettercase does not apply to bytes in a binary string. To perform lettercase conversion, the string must be converted to a nonbinary string:

```
mysql> SET NAMES binary;
Query OK, 0 rows affected (0.00 sec)
mysql> SELECT LOWER('aA'), LOWER(CONVERT('aA' USING latin1));
+-------------+-----------------------------------+
| LOWER('aA') | LOWER(CONVERT('aA' USING latin1)) |
+-------------+-----------------------------------+
| aA          | aa                                |
+-------------+-----------------------------------+
1 row in set (0.00 sec)
```

**Trailing space handling in comparisons.** Nonbinary strings have `PADSPACE` behavior for all collations, including `_bin` collations. Trailing spaces are insignificant in comparisons:

```
mysql> SET NAMES utf8 COLLATE utf8_bin;
Query OK, 0 rows affected (0.00 sec)
mysql> SELECT 'a ' = 'a';
+------------+
| 'a ' = 'a' |
+------------+
|          1 |
+------------+
```

```
1 row in set (0.00 sec)
```

For binary strings, all characters are significant in comparisons, including trailing spaces:

```
mysql> SET NAMES binary;
Query OK, 0 rows affected (0.00 sec)
mysql> SELECT 'a ' = 'a';
+------------+
| 'a ' = 'a' |
+------------+
|          0 |
+------------+
1 row in set (0.00 sec)
```

**Trailing space handling for inserts and retrievals.** CHAR(*N*) columns store nonbinary strings. Values shorter than *N* characters are extended with spaces on insertion. For retrieval, trailing spaces are removed.

BINARY(*N*) columns store binary strings. Values shorter than *N* bytes are extended with 0x00 bytes on insertion. For retrieval, nothing is removed; a value of the declared length is always returned.

```
mysql> CREATE TABLE t1 (
    ->    a CHAR(10) CHARACTER SET utf8 COLLATE utf8_bin,
    ->    b BINARY(10)
    -> );
Query OK, 0 rows affected (0.09 sec)
mysql> INSERT INTO t1 VALUES ('a','a');
Query OK, 1 row affected (0.01 sec)
mysql> SELECT HEX(a), HEX(b) FROM t1;
+--------+----------------------+
| HEX(a) | HEX(b)               |
+--------+----------------------+
| 61     | 61000000000000000000 |
+--------+----------------------+
1 row in set (0.04 sec)
```

# 1.6.5. Special Cases Where Collation Determination Is Tricky

In the great majority of statements, it is obvious what collation MySQL uses to resolve a comparison operation. For example, in the following cases, it should be clear that the collation is the collation of column x:

```
SELECT x FROM T ORDER BY x;
SELECT x FROM T WHERE x = x;
SELECT DISTINCT x FROM T;
```

However, when multiple operands are involved, there can be ambiguity. For example:

```
SELECT x FROM T WHERE x = 'Y';
```

Should this query use the collation of the column x, or of the string literal 'Y'?

Standard SQL resolves such questions using what used to be called "coercibility" rules. Basically, this means: Both x and 'Y' have collations, so which collation takes precedence? This can be difficult to resolve, but the following rules cover most situations:

- An explicit COLLATE clause has a coercibility of 0. (Not coercible at all.)

- The concatenation of two strings with different collations has a coercibility of 1.

- The collation of a column or a stored routine parameter or local variable has a coercibility of 2.

- A "system constant" (the string returned by functions such as USER() or VERSION()) has a coercibility of 3.

- A literal's collation has a coercibility of 4.

- NULL or an expression that is derived from NULL has a coercibility of 5.

The preceding coercibility values are current for MySQL 6.0.

Those rules resolve ambiguities in the following manner:

- Use the collation with the lowest coercibility value.

- If both sides have the same coercibility, then:

- If both sides are Unicode, or both sides are not Unicode, it is an error.

- If one of the sides has a Unicode character set, and another side has a non-Unicode character set, the side with Unicode character set wins, and automatic character set conversion is applied to the non-Unicode side. For example, the following statement will not return an error:

```
SELECT CONCAT(utf8_column, latin1_column) FROM t1;
```

  It will return a result, and the character set of the result will be `utf8`. The collation of the result will be the collation of `utf8_column`. Values of `latin1_column` will be automatically converted to `utf8` before concatenating.

- For an operation with operands from the same character set but that mix a `_bin` collation and a `_ci` or `_cs` collation, the `_bin` collation is used. This is similar to how operations that mix nonbinary and binary strings evaluate the operands as binary strings, except that it is for collations rather than data types.

Although automatic conversion is not in the SQL standard, the SQL standard document does say that every character set is (in terms of supported characters) a "subset" of Unicode. Because it is a well-known principle that "what applies to a superset can apply to a subset," we believe that a collation for Unicode can apply for comparisons with non-Unicode strings.

Examples:

| | |
|---|---|
| `column1 = 'A'` | Use collation of `column1` |
| `column1 = 'A' COLLATE x` | Use collation of `'A' COLLATE x` |
| `column1 COLLATE x = 'A' COLLATE y` | Error |

The `COERCIBILITY()` function can be used to determine the coercibility of a string expression:

```
mysql> SELECT COERCIBILITY('A' COLLATE latin1_swedish_ci);
        -> 0
mysql> SELECT COERCIBILITY(VERSION());
        -> 3
mysql> SELECT COERCIBILITY('A');
        -> 4
```

See Information Functions.

## 1.6.6. Collations Must Be for the Right Character Set

Each character set has one or more collations, but each collation is associated with one and only one character set. Therefore, the following statement causes an error message because the `latin2_bin` collation is not legal with the `latin1` character set:

```
mysql> SELECT _latin1 'x' COLLATE latin2_bin;
ERROR 1253 (42000): COLLATION 'latin2_bin' is not valid
for CHARACTER SET 'latin1'
```

## 1.6.7. Examples of the Effect of Collation

**Example 1: Sorting German Umlauts**

Suppose that column `X` in table `T` has these `latin1` column values:

```
Muffler
Müller
MX Systems
MySQL
```

Suppose also that the column values are retrieved using the following statement:

```
SELECT X FROM T ORDER BY X COLLATE collation_name;
```

The following table shows the resulting order of the values if we use `ORDER BY` with different collations.

| `latin1_swedish_ci` | `latin1_german1_ci` | `latin1_german2_ci` |
|---|---|---|
| Muffler | Muffler | Müller |
| MX Systems | Müller | Muffler |

| MÃ¼ller | MX Systems | MX Systems |
|---------|-----------|------------|
| MySQL | MySQL | MySQL |

The character that causes the different sort orders in this example is the U with two dots over it (ü), which the Germans call "U-umlaut."

- The first column shows the result of the SELECT using the Swedish/Finnish collating rule, which says that U-umlaut sorts with Y.

- The second column shows the result of the SELECT using the German DIN-1 rule, which says that U-umlaut sorts with U.

- The third column shows the result of the SELECT using the German DIN-2 rule, which says that U-umlaut sorts with UE.

**Example 2: Searching for German Umlauts**

Suppose that you have three tables that differ only by the character set and collation used:

```
mysql> CREATE TABLE german1 (
    ->    c CHAR(10)
    -> ) CHARACTER SET latin1 COLLATE latin1_german1_ci;
mysql> CREATE TABLE german2 (
    ->    c CHAR(10)
    -> ) CHARACTER SET latin1 COLLATE latin1_german2_ci;
mysql> CREATE TABLE germanutf8 (
    ->    c CHAR(10)
    -> ) CHARACTER SET utf8 COLLATE utf8_unicode_ci;
```

Each table contains two records:

```
mysql> INSERT INTO german1 VALUES ('Bar'), ('BÃ¤r');
mysql> INSERT INTO german2 VALUES ('Bar'), ('BÃ¤r');
mysql> INSERT INTO germanutf8 VALUES ('Bar'), ('BÃ¤r');
```

Two of the above collations have an A = Ã# equality, and one has no such equality (latin1_german2_ci). For that reason, you'll get these results in comparisons:

```
mysql> SELECT * FROM german1 WHERE c = 'BÃ¤r';
+------+
| c    |
+------+
| Bar  |
| BÃ¤r  |
+------+
mysql> SELECT * FROM german2 WHERE c = 'BÃ¤r';
+------+
| c    |
+------+
| BÃ¤r  |
+------+
mysql> SELECT * FROM germanutf8 WHERE c = 'BÃ¤r';
+------+
| c    |
+------+
| Bar  |
| BÃ¤r  |
+------+
```

This is not a bug but rather a consequence of the sorting that latin1_german1_ci or utf8_unicode_ci do (the sorting shown is done according to the German DIN 5007 standard).

# 1.7. String Repertoire

The *repertoire* of a character set is the collection of characters in the set.

String expressions have a repertoire attribute, which can have two values:

- ASCII: The expression can contain only characters in the Unicode range U+0000 to U+007F.

- UNICODE: The expression can contain characters in the Unicode range U+0000 to U+FFFF.

The ASCII range is a subset of UNICODE range, so a string with ASCII repertoire can be converted safely without loss of in-

formation to the character set of any string with `UNICODE` repertoire or to a character set that is a superset of `ASCII`. (All MySQL character sets are supersets of `ASCII` with the exception of `swe7`, which reuses some punctuation characters for Swedish accented characters.) The use of repertoire enables character set conversion in expressions for many cases where MySQL would otherwise return an "illegal mix of collations" error.

The following discussion provides examples of expressions and their repertoires, and describes how the use of repertoire changes string expression evaluation:

- The repertoire for string constants depends on string content:

```
SET NAMES utf8; SELECT 'abc';
SELECT _utf8'def';
SELECT N'MySQL';
```

  Although the character set is `utf8` in each of the preceding cases, the strings do not actually contain any characters outside the ASCII range, so their repertoire is `ASCII` rather than `UNICODE`.

- Columns having the `ascii` character set have `ASCII` repertoire because of their character set. In the following table, `c1` has `ASCII` repertoire:

```
CREATE TABLE t1 (c1 CHAR(1) CHARACTER SET ascii);
```

  The following example illustrates how repertoire enables a result to be determined in a case where an error occurs without repertoire:

```
CREATE TABLE t1 (
  c1 CHAR(1) CHARACTER SET latin1,
  c2 CHAR(1) CHARACTER SET ascii
);
INSERT INTO t1 VALUES ('a','b');
SELECT CONCAT(c1,c2) FROM t1;
```

  Without repertoire, this error occurs:

```
ERROR 1267 (HY000): Illegal mix of collations (latin1_swedish_ci,IMPLICIT)
and (ascii_general_ci,IMPLICIT) for operation 'concat'
```

  Using repertoire, subset to superset (`ascii` to `latin1`) conversion can occur and a result is returned:

```
+---------------+
| CONCAT(c1,c2) |
+---------------+
| ab            |
+---------------+
```

- Functions with one string argument inherit the repertoire of their argument. The result of `UPPER(_utf8'abc')` has `ASCII` repertoire, because its argument has `ASCII` repertoire.

- For functions that return a string but do not have string arguments and use `character_set_connection` as the result character set, the result repertoire is `ASCII` if `character_set_connection` is `ascii`, and `UNICODE` otherwise:

```
FORMAT(numeric_column, 4);
```

  Use of repertoire changes how MySQL evaluates the following example:

```
SET NAMES ascii;
CREATE TABLE t1 (a INT, b VARCHAR(10) CHARACTER SET latin1);
INSERT INTO t1 VALUES (1,'b');
SELECT CONCAT(FORMAT(a, 4), b) FROM t1;
```

  Without repertoire, this error occurs:

```
ERROR 1267 (HY000): Illegal mix of collations (ascii_general_ci,COERCIBLE)
and (latin1_swedish_ci,IMPLICIT) for operation 'concat'
```

  With repertoire, a result is returned:

```
+-----------------------+
| CONCAT(FORMAT(a, 4), b) |
+-----------------------+
| 1.0000b               |
+-----------------------+
```

- Functions with two or more string arguments use the "widest" argument repertoire for the result repertoire (`UNICODE` is wider than `ASCII`). Consider the following `CONCAT()` calls:

```
CONCAT(_ucs2 0x0041, _ucs2 0x0042)
CONCAT(_ucs2 0x0041, _ucs2 0x00C2)
```

  For the first call, the repertoire is `ASCII` because both arguments are within the range of the `ascii` character set. For the second call, the repertoire is `UNICODE` because the second argument is outside the `ascii` character set range.

- The repertoire for function return values is determined based only on the repertoire of the arguments that affect the result's character set and collation.

```
IF(column1 < column2, 'smaller', 'greater')
```

  The result repertoire is `ASCII` because the two string arguments (the second argument and the third argument) both have `ASCII` repertoire. The first argument does not matter for the result repertoire, even if the expression uses string values.

# 1.8. Operations Affected by Character Set Support

This section describes operations that take character set information into account.

## 1.8.1. Result Strings

MySQL has many operators and functions that return a string. This section answers the question: What is the character set and collation of such a string?

For simple functions that take string input and return a string result as output, the output's character set and collation are the same as those of the principal input value. For example, `UPPER(X)` returns a string whose character string and collation are the same as that of `X`. The same applies for `INSTR()`, `LCASE()`, `LOWER()`, `LTRIM()`, `MID()`, `REPEAT()`, `REPLACE()`, `REVERSE()`, `RIGHT()`, `RPAD()`, `RTRIM()`, `SOUNDEX()`, `SUBSTRING()`, `TRIM()`, `UCASE()`, and `UPPER()`.

Note: The `REPLACE()` function, unlike all other functions, always ignores the collation of the string input and performs a case-sensitive comparison.

If a string input or function result is a binary string, the string has no character set or collation. This can be checked by using the `CHARSET()` and `COLLATION()` functions, both of which return `binary` to indicate that their argument is a binary string:

```
mysql> SELECT CHARSET(BINARY 'a'), COLLATION(BINARY 'a');
+---------------------+-----------------------+
| CHARSET(BINARY 'a') | COLLATION(BINARY 'a') |
+---------------------+-----------------------+
| binary              | binary                |
+---------------------+-----------------------+
```

For operations that combine multiple string inputs and return a single string output, the "aggregation rules" of standard SQL apply for determining the collation of the result:

- If an explicit `COLLATE X` occurs, use `X`.

- If explicit `COLLATE X` and `COLLATE Y` occur, raise an error.

- Otherwise, if all collations are `X`, use `X`.

- Otherwise, the result has no collation.

For example, with `CASE ... WHEN a THEN b WHEN b THEN c COLLATE X END`, the resulting collation is `X`. The same applies for `UNION`, `||`, `CONCAT()`, `ELT()`, `GREATEST()`, `IF()`, and `LEAST()`.

For operations that convert to character data, the character set and collation of the strings that result from the operations are defined by the `character_set_connection` and `collation_connection` system variables. This applies only to `CAST()`, `CONV()`, `FORMAT()`, `HEX()`, and `SPACE()`.

If you are uncertain about the character set or collation of the result returned by a string function, you can use the `CHARSET()` or `COLLATION()` function to find out:

```
mysql> SELECT USER(), CHARSET(USER()), COLLATION(USER());
+----------------+-----------------+-------------------+
| USER()         | CHARSET(USER()) | COLLATION(USER()) |
```

```
+----------------+----------------+-------------------+
| test@localhost | utf8           | utf8_general_ci   |
+----------------+----------------+-------------------+
```

## 1.8.2. `CONVERT()` and `CAST()`

`CONVERT()` provides a way to convert data between different character sets. The syntax is:

```
CONVERT(expr USING transcoding_name)
```

In MySQL, transcoding names are the same as the corresponding character set names.

Examples:

```
SELECT CONVERT(_latin1'Müller' USING utf8);
INSERT INTO utf8table (utf8column)
    SELECT CONVERT(latin1field USING utf8) FROM latin1table;
```

`CONVERT(... USING ...)` is implemented according to the standard SQL specification.

You may also use `CAST()` to convert a string to a different character set. The syntax is:

```
CAST(character_string AS character_data_type CHARACTER SET charset_name)
```

Example:

```
SELECT CAST(_latin1'test' AS CHAR CHARACTER SET utf8);
```

If you use `CAST()` without specifying `CHARACTER SET`, the resulting character set and collation are defined by the `character_set_connection` and `collation_connection` system variables. If you use `CAST()` with `CHARACTER SET X`, the resulting character set and collation are `X` and the default collation of `X`.

You may not use a `COLLATE` clause inside a `CAST()`, but you may use it outside. That is, `CAST(... COLLATE ...)` is illegal, but `CAST(...) COLLATE ...` is legal.

Example:

```
SELECT CAST(_latin1'test' AS CHAR CHARACTER SET utf8) COLLATE utf8_bin;
```

## 1.8.3. `SHOW` Statements and `INFORMATION_SCHEMA`

Several `SHOW` statements provide additional character set information. These include `SHOW CHARACTER SET`, `SHOW COLLATION`, `SHOW CREATE DATABASE`, `SHOW CREATE TABLE` and `SHOW COLUMNS`. These statements are described here briefly. For more information, see `SHOW` Syntax.

`INFORMATION_SCHEMA` has several tables that contain information similar to that displayed by the `SHOW` statements. For example, the `CHARACTER_SETS` and `COLLATIONS` tables contain the information displayed by `SHOW CHARACTER SET` and `SHOW COLLATION`. See `INFORMATION_SCHEMA` Tables.

The `SHOW CHARACTER SET` command shows all available character sets. It takes an optional `LIKE` clause that indicates which character set names to match. For example:

```
mysql> SHOW CHARACTER SET LIKE 'latin%';
+---------+-----------------------------+-------------------+--------+
| Charset | Description                 | Default collation | Maxlen |
+---------+-----------------------------+-------------------+--------+
| latin1  | cp1252 West European        | latin1_swedish_ci |      1 |
| latin2  | ISO 8859-2 Central European | latin2_general_ci |      1 |
| latin5  | ISO 8859-9 Turkish          | latin5_turkish_ci |      1 |
| latin7  | ISO 8859-13 Baltic          | latin7_general_ci |      1 |
+---------+-----------------------------+-------------------+--------+
```

The output from `SHOW COLLATION` includes all available character sets. It takes an optional `LIKE` clause that indicates which collation names to match. For example:

```
mysql> SHOW COLLATION LIKE 'latin1%';
+-------------------+---------+----+---------+----------+---------+
| Collation         | Charset | Id | Default | Compiled | Sortlen |
+-------------------+---------+----+---------+----------+---------+
| latin1_german1_ci | latin1  |  5 |         |          |       0 |
| latin1_swedish_ci | latin1  |  8 | Yes     | Yes      |       0 |
| latin1_danish_ci  | latin1  | 15 |         |          |       0 |
| latin1_german2_ci | latin1  | 31 |         | Yes      |       2 |
```

```
| latin1_bin        | latin1  | 47 |         | Yes     |        0 |
| latin1_general_ci | latin1  | 48 |         |         |        0 |
| latin1_general_cs | latin1  | 49 |         |         |        0 |
| latin1_spanish_ci | latin1  | 94 |         |         |        0 |
+-------------------+---------+----+---------+---------+----------+
```

`SHOW CREATE DATABASE` displays the `CREATE DATABASE` statement that creates a given database:

```
mysql> SHOW CREATE DATABASE test;
+----------+----------------------------------------------------------------+
| Database | Create Database                                                |
+----------+----------------------------------------------------------------+
| test     | CREATE DATABASE `test` /*!40100 DEFAULT CHARACTER SET latin1 */ |
+----------+----------------------------------------------------------------+
```

If no `COLLATE` clause is shown, the default collation for the character set applies.

`SHOW CREATE TABLE` is similar, but displays the `CREATE TABLE` statement to create a given table. The column definitions indicate any character set specifications, and the table options include character set information.

The `SHOW COLUMNS` statement displays the collations of a table's columns when invoked as `SHOW FULL COLUMNS`. Columns with `CHAR`, `VARCHAR`, or `TEXT` data types have collations. Numeric and other non-character types have no collation (indicated by `NULL` as the `Collation` value). For example:

```
mysql> SHOW FULL COLUMNS FROM person\G
*************************** 1. row ***************************
     Field: id
      Type: smallint(5) unsigned
 Collation: NULL
      Null: NO
       Key: PRI
   Default: NULL
     Extra: auto_increment
Privileges: select,insert,update,references
   Comment:
*************************** 2. row ***************************
     Field: name
      Type: char(60)
 Collation: latin1_swedish_ci
      Null: NO
       Key:
   Default:
     Extra:
Privileges: select,insert,update,references
   Comment:
```

The character set is not part of the display but is implied by the collation name.

# 1.9. Unicode Support

The initial implementation of Unicode support (in MySQL 4.1) included two character sets for storing Unicode data:

- `ucs2`, the UCS-2 encoding of the Unicode character set using 16 bits per character

- `utf8`, a UTF-8 encoding of the Unicode character set using one to three bytes per character

These two character sets support the characters from the Basic Multilingual Plane (BMP) of Unicode Version 3.0. BMP characters have these characteristics:

- Their code values are between 0 and 65535 (or `U+0000` .. `U+FFFF`)

- They can be encoded with a fixed 16-bit word, as in `ucs2`

- They can be encoded with 8, 16, or 24 bits, as in `utf8`

- They are sufficient for almost all characters in major languages

Characters not supported by the aforementioned character sets include supplementary characters that lie outside the BMP. As of MySQL 6.0 (versions 6.0.4 and up), Unicode support is extended to include supplementary characters, which requires new character sets that have a broader range and therefore take more space. The following table shows a brief feature comparison of previous and current Unicode support.

| Before MySQL 6.0 | MySQL 6.0 |
|---|---|

| All Unicode 3.0 characters | All Unicode 5.0 characters |
|---|---|
| No supplementary characters | With supplementary characters |
| `ucs2` character set | No change |
| `utf8` character set for up to three bytes | Renamed to `utf8mb3` |
| | New `utf8` character set for up to four bytes |
| | New `utf16` character set |
| | New `utf32` character set |

Most of these changes are upward compatible. For example, the old `utf8` character set (the three-byte version) has been renamed to `utf8mb3`, so tables created before MySQL 6.0 that used `utf8` are treated as using `utf8mb3` after an in-place upgrade to MySQL 6.0. However, the table contents will not have changed in any way.

The new `utf8` character set, as well as `utf16` and `utf32`, support supplementary characters. The Unicode character sets in MySQL 6.0 are:

- `ucs2`, the UCS-2 encoding of the Unicode character set using 16 bits per character

- `utf16`, the UTF-16 encoding for the Unicode character set; like `ucs2` but with an extension for supplementary characters

- `utf32`, the UTF-32 encoding for the Unicode character set using 32 bits per character

- `utf8`, a UTF-8 encoding of the Unicode character set using one to four bytes per character

- `utf8mb3`, a UTF-8 encoding of the Unicode character set using one to three bytes per character (known as `utf8` before MySQL 6.0)

A similar set of collations is available for each Unicode character set. For example, each has a Danish collation, the names of which are `ucs2_danish_ci`, `utf16_danish_ci`, `ucs32_danish_ci`, `utf8_danish_ci`, and `utf8mb3_danish_ci`. All Unicode collations are listed at Section 1.13.1, "Unicode Character Sets", which also describes collation properties for supplementary characters.

Note that although many of the supplementary characters come from East Asian languages, what MySQL 6.0 adds is support for more Japanese and Chinese characters in Unicode character sets, not support for new Japanese and Chinese character sets.

The following discussion provides additional detail on the Unicode character sets in MySQL. For details on potential incompatibility issues for your applications, see Section 1.10, "Upgrading from Previous to Current Unicode Support". That section also describes how to convert tables from `utf8mb3` to the new (four-byte) `utf8` character set, and what constraints may apply in doing so.

The MySQL implementation of UCS-2, UTF-16, and UTF-32 stores characters in big-endian byte order and does not use a byte order mark (BOM) at the beginning of values. Other database systems might use little-endian byte order or a BOM, in which case, conversion of values will need to be performed when transferring data between those systems and MySQL.

MySQL uses no BOM for UTF-8 values.

Client applications that need to communicate with the server using Unicode should set the client character set accordingly; for example, by issuing a `SET NAMES 'utf8'` statement. `ucs2`, `utf16`, and `utf32` cannot be used as a client character set, which means that they do not work for `SET NAMES` or `SET CHARACTER SET`. (See Section 1.4, "Connection Character Sets and Collations".)

**The `ucs2` Character Set (UCS-2 Unicode Encoding)**

In UCS-2, every character is represented by a two-byte Unicode code with the most significant byte first. For example: `LATIN CAPITAL LETTER A` has the code `0x0041` and it is stored as a two-byte sequence: `0x00 0x41`. `CYRILLIC SMALL LETTER YERU` (Unicode `0x044B`) is stored as a two-byte sequence: `0x04 0x4B`. For Unicode characters and their codes, please refer to the Unicode Home Page.

In MySQL, the `ucs2` character set is a fixed 16-bit encoding for Unicode BMP characters.

**The `utf16` Character Set (UTF-16 Unicode Encoding)**

The `utf16` character set is the `ucs2` character set with an extension that enables encoding of supplementary characters:

- For a BMP character, `utf16` and `ucs2` have identical storage characteristics: The same code values, same encoding, and same

length.

- For a supplementary character, `utf16` has a special sequence for representing the character using 32 bits. This is called the "surrogate" mechanism: For a number greater than `0xffff`, take 10 bits and add them to `0xd800` and put them in the first 16-bit word, take 10 more bits and add them to `0xdc00` and put them in the next 16-bit word. Consequently, all supplementary characters require 32 bits, where the first 16 bits are a number between `0xd800` and `0xdbff`, and the last 16 bits are a number between `0xdc00` and `0xdfff`. Examples are in 15.5 Surrogates Area in the Unicode 4.0 document.

Because `utf16` supports surrogates and `ucs2` does not, there is a validity check that applies only in `utf16`: You cannot insert a top surrogate without a bottom surrogate, or vice versa. For example:

```
INSERT INTO t (ucs2_column) VALUES (0xd800); /* legal */
INSERT INTO t (utf16_column)VALUES (0xd800); /* illegal */
```

There is no validity check for characters that are technically valid but are not true Unicode (that is, characters that Unicode considers to be "unassigned code points" or "private use" characters or even "illegals" like `0xffff`). For example, since `U+F8FF` is the Apple Logo, this is legal:

```
INSERT INTO t (utf16_column)VALUES (0xf8ff); /* legal */
```

Such characters cannot be expected to mean the same thing to everyone.

Because MySQL must allow for the worst case (that one character requires four bytes) the maximum length of a `utf16` column or index is only half of the maximum length for a `ucs2` column or index. For example, at the time of writing the maximum length of a Falcon index key is 1100 bytes (assuming the default Falcon page size), so these statements create tables with the longest allowed indexes for `ucs2` and `utf16` columns:

```
CREATE TABLE tf (s1 VARCHAR(550) CHARACTER SET ucs2) ENGINE=FALCON;
CREATE INDEX i ON tf (s1);
CREATE TABLE tg (s1 VARCHAR(275) CHARACTER SET utf16) ENGINE=FALCON;
CREATE INDEX i ON tg (s1);
```

### The `utf32` Character Set (UTF-32 Unicode Encoding)

The `utf32` character set is fixed length (like `ucs2` and unlike `utf16`). `utf32` uses 32 bits for every character, unlike `ucs2` (which uses 16 bits for every character), and unlike `utf16` (which uses 16 bits for some characters and 32 bits for others).

`utf32` takes twice as much space as `ucs2` and more space than `utf16`, but `utf32` has the same advantage as `ucs2` that it is predictable for storage: The required number of bytes for `utf32` equals the number of characters times 4. Also, unlike `utf16`, there are no tricks for encoding in `utf32`, so the stored value equals the code value.

To demonstrate how the latter advantage is useful, here is an example that shows how to determine a `utf8` value given the `utf32` code value:

```
/* Assume code value = 100cc LINEAR B WHEELED CHARIOT */
CREATE TABLE tmp (utf32 CHAR(1) CHARACTER SET utf32,
                  utf8 CHAR(1) CHARACTER SET utf8);
INSERT INTO tmp VALUES (0x000100cc,NULL);
UPDATE tmp SET utf8 = utf32;
SELECT HEX(utf32),HEX(utf8) FROM tmp;
```

MySQL is very forgiving about additions of unassigned Unicode characters, private-use-area characters, and other code values not present in the official Unicode 5.0 Character Database. There is in fact only one validity check for `utf32`: No code value may be greater than `0x10ffff`. For example, this is illegal:

```
INSERT INTO t (utf32_column) VALUES (0x110000); /* illegal */
```

### The `utf8` Character Set (Four-Byte UTF-8 Unicode Encoding)

UTF-8 (Unicode Transformation Format with 8-bit units) is an alternative way to store Unicode data. It is implemented according to RFC 3629. RFC 3629 describes encoding sequences that take from one to four bytes. (An older standard for UTF-8 encoding is given by RFC 2279, which describes UTF-8 sequences that take from one to six bytes. RFC 3629 renders RFC 2279 obsolete; for this reason, sequences with five and six bytes are no longer used.)

The idea of UTF-8 is that various Unicode characters are encoded using byte sequences of different lengths:

- Basic Latin letters, digits, and punctuation signs use one byte.

- Most European and Middle East script letters fit into a two-byte sequence: extended Latin letters (with tilde, macron, acute,

grave and other accents), Cyrillic, Greek, Armenian, Hebrew, Arabic, Syriac, and others.

• Korean, Chinese, and Japanese ideographs use three-byte or four-byte sequences.

**Tip**: To save space with UTF-8, use `VARCHAR` instead of `CHAR`. Otherwise, MySQL must reserve four bytes for each character in a `CHAR CHARACTER SET utf8` column because that is the maximum possible length. For example, MySQL must reserve 40 bytes for a `CHAR(10) CHARACTER SET utf8` column.

Before MySQL 6.0, the character set named `utf8` contained only BMP characters and used a maximum of three bytes per character. As of MySQL 6.0, that character set uses a different name, `utf8mb3`, but its characteristics remain the same as before.

In terms of table content (for old tables created before MySQL 6.0), the old `utf8` (now `utf8mb3`)) and new `utf8` are compatible:

• For a BMP character, the 5.1 and 6.0 versions of `utf8` have identical storage characteristics: same code values, same encoding, same length.

• For a supplementary character, `utf8` in 5.1 cannot store the character at all, while `utf8` in 6.0 requires four bytes to store. Since `utf8` in 5.1 cannot store the character at all, you do not have any supplementary characters in `utf8` columns in your 5.1 database, so you need not worry about converting characters or losing data.

**The `utf8mb3` Character Set (Three-Byte UTF-8 Unicode Encoding)**

The `utf8mb3` character set in MySQL 6.0 is the same as the character set that was called `utf8` before 6.0. In other words, `utf8mb3` in MySQL 6.0 has exactly the same characteristics as `utf8` in MySQL 5.1:

• No support for supplementary characters (BMP characters only)

• Maximum of three bytes per multi-byte character

For the few cases where it is desirable to have complete compatibility in MySQL 6.0 with the old `utf8` character set, you can define with `CHARACTER SET utf8mb3`.

Exactly the same set of characters is available in `utf8mb3` as in `ucs2`. That is, they have the same repertoire.

# 1.10. Upgrading from Previous to Current Unicode Support

This section describes issues pertaining to Unicode support that you may face when upgrading from MySQL 5.1 to 6.0. It also provides guidelines for downgrading from 6.0 back to 5.1.

In most respects, upgrading from MySQL 5.1 to 6.0 should present few problems with regard to Unicode usage, although there are some potential areas of incompatibility. Some examples:

• For the variable-length character data types (`VARCHAR` and the `TEXT` types), the maximum length in characters for `utf8` columns is less in MySQL 6.0 than previously.

• For all character data types (`CHAR`, `VARCHAR`, and the `TEXT` types), the maximum number of characters for `utf8` columns that can be indexed is less in MySQL 6.0 than previously.

Consequently, if you want to upgrade tables from the old `utf8` (now `utf8mb3`) to the current `utf8`, it may be necessary to change some column or index definitions.

You can upgrade from MySQL 5.1 to 6.0 in two different ways:

• You can install 6.0 "in place" on top of an existing 5.1 installation. In this case, the names of `utf8` character sets and collations will change to use `utf8mb3`, but no changes to column or index lengths in table definitions will be needed. That is, there will be name differences but no structural changes.

• You can dump your 5.1 data using `mysqldump`, and then reload the dump into 6.0. Because definitions in the dump file will refer to `utf8`, the server will use `utf8` in the reloaded tables, but these tables will use the new (four-byte) `utf8`, not the old 5.1 (three-byte) `utf8`. In this case, the names of `utf8` character sets and collations will remain the same, but changes to column or index lengths might be needed for long columns. In other words, there will not be name differences, but you might need to make some structural changes.

If you upgrade MySQL in place by installing 6.0 on top of an existing 5.1 installation, the changes in the 6.0 Unicode implementation have these effects:

- For ucs2, there should be no issues

- Database and tables that used utf8 in MySQL 5.1 will be reported as using utf8mb3 in 6.0 (for example, if you examine object structure with SHOW statements or select definitions from INFORMATION_SCHEMA tables). However, no data will have been changed and the server can use utf8mb3 columns and indexes with no conversion required.

- When the server starts, it checks the structure of certain system tables in the mysql database. If it finds that columns have the utf8mb3 character set when it expects utf8, it writes warnings to the error log but continues to use the tables. To make the warnings go away, convert the system tables to utf8 by running mysql_upgrade.

If you use events, a known issue is that if you upgrade from MySQL 5.1 to 6.0.4 though 6.0.6, the event scheduler will not work, even after you run mysql_upgrade. (This is an issue only for an upgrade, not for a new installation of MySQL 6.0.x.) As of MySQL 6.0.7, mysql_upgrade handles upgrading the system tables properly. If you upgrade to 6.0.4 through 6.0.6, you can work around this upgrading problem by using these instructions:

1. In MySQL 5.1, before upgrading, create a dump file containing your mysql.event table:

   ```
   shell> mysqldump -uroot -p mysql event > event.sql
   ```

2. Stop the server, upgrade to MySQL 6.0, and start the server.

3. Recreate the mysql.event table using the dump file:

   ```
   shell> mysql -uroot -p mysql < event.sql
   ```

4. Run mysql_upgrade to upgrade the other system tables in the mysql database:

   ```
   shell> mysql_upgrade -uroot -p
   ```

5. Restart the server. The event scheduler should run normally.

To following example illustrates what you should see if you use MySQL 5.1 tables with an installation that has been upgraded in place to 6.0. Suppose that a table was originally defined as follows in MySQL 5.1:

```
CREATE TABLE t1 (
  col1 CHAR(10) CHARACTER SET utf8 COLLATE utf8_unicode_ci NOT NULL,
  col2 CHAR(10) CHARACTER SET utf8 COLLATE utf8_bin NOT NULL
) CHARACTER SET utf8;
```

In MySQL 5.1, SHOW CREATE TABLE produces this result:

```
mysql> SHOW CREATE TABLE t1\G
*************************** 1. row ***************************
       Table: t1
Create Table: CREATE TABLE `t1` (
  `col1` char(10) CHARACTER SET utf8 COLLATE utf8_unicode_ci NOT NULL,
  `col2` char(10) CHARACTER SET utf8 COLLATE utf8_bin NOT NULL
) ENGINE=MyISAM DEFAULT CHARSET=utf8
```

After upgrading to MySQL 6.0, SHOW CREATE TABLE produces this result instead:

```
mysql> SHOW CREATE TABLE t1\G
*************************** 1. row ***************************
       Table: t1
Create Table: CREATE TABLE `t1` (
  `col1` char(10) CHARACTER SET utf8mb3 COLLATE utf8mb3_unicode_ci NOT NULL,
  `col2` char(10) CHARACTER SET utf8mb3 COLLATE utf8mb3_bin NOT NULL
) ENGINE=MyISAM DEFAULT CHARSET=utf8mb3
```

Internally, the IDs for the character data have not changed. Only the names associated with the IDs have changed. For example, in MySQL 5.1, we have this:

```
mysql> SHOW COLLATION LIKE 'utf8%bin';
+-----------+---------+----+---------+----------+---------+
| Collation | Charset | Id | Default | Compiled | Sortlen |
+-----------+---------+----+---------+----------+---------+
| utf8_bin  | utf8    | 83 |         | Yes      |       1 |
```

```
+-----------+---------+----+---------+----------+---------+
```

In MySQL 6.0, we have this:

```
mysql> SHOW COLLATION LIKE 'utf8%bin';
+-------------+---------+----+---------+----------+---------+
| Collation   | Charset | Id | Default | Compiled | Sortlen |
+-------------+---------+----+---------+----------+---------+
| utf8mb3_bin | utf8mb3 | 83 |         | Yes      |       1 |
| utf8_bin    | utf8    | 46 |         | Yes      |       1 |
+-------------+---------+----+---------+----------+---------+
```

Thus, for table `t1`, the columnn `col2` still has a collation ID of 83. What has changed is the name associated with 83. In other words, the collation IDs for the old `utf8` character set now belong to `utf8mb3`. The collations IDs for the new `utf8` character set are new.

Tables can be converted from `utf8mb3` to `utf8` by using `ALTER TABLE`. The following statement does so for `t1`:

```
ALTER TABLE t1
  DEFAULT CHARACTER SET utf8,
  MODIFY col1 CHAR(10) CHARACTER SET utf8 COLLATE utf8_unicode_ci NOT NULL,
  MODIFY col2 CHAR(10) CHARACTER SET utf8 COLLATE utf8_bin NOT NULL;
```

In terms of table content, conversion from the old `utf8` to the new `utf8` presents no problems:

- For a BMP character, the 5.1 and 6.0 versions of `utf8` have identical storage characteristics: same code values, same encoding, same length.

- For a supplementary character, `utf8` in 5.1 cannot store the character at all, while `utf8` in 6.0 requires four bytes to store. Since `utf8` in 5.1 cannot store the character at all, you do not have any supplementary characters in `utf8` columns in your 5.1 database, so you need not worry about converting characters or losing data.

In terms of table structure, the catch when converting from the old `utf8` to the new `utf8` is that the maximum length of a column or index key is unchanged in terms of bytes. Therefore, it is smaller in terms of characters, because the maximum length of a character is four bytes instead of three. The things to watch for when converting your MySQL 5.1 tables are:

- Look at the definitions of `utf8` columns, and make sure they will not exceed the maximum length for the storage engine.

- Look at all indexes on `utf8` columns, and make sure they will not exceed the maximum length for the storage engine. Sometimes the maximum can change due to storage engine enhancements.

Check those things for the `CHAR`, `VARCHAR`, and `TEXT` data types.

If the preceding conditions apply for you, you will have to reduce the defined length of columns or indexes, or you will have to use `utf8mb3` rather than `utf8`.

Here are some examples where structural changes may be needed:

- A `TINYTEXT` column can hold up to 255 bytes, so it can hold up to 85 three-byte or 63 four-byte characters. Suppose that you have a `TINYTEXT` column that uses `utf8mb3` but must be able to contain more than 63 characters. You cannot convert it to `utf8` unless you also change the data type to a longer type such as `TEXT`.

  Similarly, a very long `VARCHAR` column may need to be changed to one of the longer `TEXT` types if you want to convert it from `utf8mb3` to `utf8`.

- `InnoDB` has a maximum index length of 767 bytes, so for `utf8mb3` or `utf8` columns, you can index a maximum of 255 or 191 characters, respectively. If you currently have `utf8mb3` columns with indexes longer than 191 characters, you will need to index a smaller number of characters. In an `InnoDB` table, this column and index definition is legal:

  ```
  col1 VARCHAR(500) CHARACTER SET utf8mb3, INDEX (col1(255))
  ```

  To use `utf8` instead, the index must be smaller:

  ```
  col1 VARCHAR(500) CHARACTER SET utf8, INDEX (col1(191))
  ```

The preceding types of changes are most likely to be required only if you have very long columns or indexes. Otherwise, you

should be able to convert your tables from `utf8mb3` to `utf8` without problems. You can do this either by dumping them with `mysqldump` in MySQL 5.1 and reloading the dump into 6.0, or by using `ALTER TABLE` after upgrading in place from 5.1 to 6.0.

The following items summarize other potential areas of incompatibility:

- Performance of four-byte UTF-8 (`utf8`) is slower than for three-byte UTF-8 (`utf8mb3`). There are tradeoffs: If you want the same character set and collation names, you can convert from `utf8mb3` to `utf8`. But that will introduce a small performance penalty. If you do not want to incur this penalty, continue to use `utf8mb3`.

- `SET NAMES 'utf8'` now causes use of four-byte character set. As long as no four-byte character are sent in either direction, there should be no problems. Otherwise, applications that expect to receive a maximum of three bytes per character may have problems. Conversely, applications that expect to send four-byte characters must ensure that the server understands them.

- Applications cannot send `utf16` or `utf32` character data to a pre-6.0 server.

- Applications that use the old collation IDs to detect `utf8` collations need to be updated unless it is okay that these IDs signify `utf8mb3` collations in MySQL 6.0.

- Applications that test for the string `'utf8'` in character set or collation names and expect them to refer to the three-byte `utf8` may need to be adjusted to recognize `'utf8mb3'` instead.

- Look at all application programs that process `utf8` data, and make sure that buffer lengths are not calculated with an algorithm such as "number of characters times 3."

- For replication, if four-byte characters are going to be used, all servers involved must understand them.

- Check for other uses of `utf8`, such as `_utf8` introducers, `utf8` stored routine variables, or `utf8` function return types.

Similar restrictions apply if you attempt to replicate from a MySQL 6.0 master to a MySQL 5.1 slave. `utf8mb3` data will be seen as `utf8` by the slave as long as it is represented in binary log events by ID, rather than by name (as in statements such as `CREATE TABLE`). But you cannot send `utf16` or `utf32` data, or data for the new `utf8`.

If you have upgraded to MySQL 6.0 and then decide to downgrade back to 5.1, these considerations apply:

- `ucs2` data should present no problems.

- For an in-place downgrade, `utf8mb3` data will be seen as `utf8` in 5.1 and should present no problems.

- Any definitions that refer to the `utf8mb3`, `utf16`, or `utf32` character set names will not be recognized by 5.1.

- For objects with definitions that refer to the `utf8` character set name, you can dump them with `mysqldump` in 6.0 and then reload the dump in 5.1, as long as there are no four-byte characters in the data. The 5.1 server will see `utf8` in the dump file object definitions and create new objects that use the 5.1 (three-byte) `utf8` character set.

- For tables in the `mysql` database, `mysql_upgrade` in 6.0 converts `utf8mb3` columns to `utf8`. Thus, if you run `mysql_upgrade` after upgrading to 6.0, these tables will not be legal in 5.1. You will need to convert `utf8` columns to `utf8mb3` with `ALTER TABLE` before downgrading, or you can do this:

  1. In MySQL 6.0, before downgrading, create a dump file containing the `mysql` database tables:

     ```
     shell> mysqldump -uroot -p mysql > mysql.sql
     ```

  2. Check the dump file to make sure that there are no instances of "utf8mb3". If there are, change them to "utf8".

  3. Stop the server, downgrade to MySQL 5.1, and start the server with the `--skip-grant-tables` option so that it does not try to use the `mysql` database tables.

  4. Recreate the dump file to recreate the `mysql` database tables:

     ```
     shell> mysql mysql < mysql.sql
     ```

  5. Restart the server. It should use the `mysql` database tables normally.

# 1.11. UTF-8 for Metadata

*Metadata* is "the data about the data." Anything that *describes* the database — as opposed to being the *contents* of the database —

is metadata. Thus column names, database names, user names, version names, and most of the string results from `SHOW` are metadata. This is also true of the contents of tables in `INFORMATION_SCHEMA`, because those tables by definition contain information about database objects.

Representation of metadata must satisfy these requirements:

- All metadata must be in the same character set. Otherwise, neither the `SHOW` commands nor `SELECT` statements for tables in `INFORMATION_SCHEMA` would work properly because different rows in the same column of the results of these operations would be in different character sets.

- Metadata must include all characters in all languages. Otherwise, users would not be able to name columns and tables using their own languages.

To satisfy both requirements, MySQL stores metadata in a Unicode character set, namely UTF-8. This does not cause any disruption if you never use accented or non-Latin characters. But if you do, you should be aware that metadata is in UTF-8.

The metadata requirements mean that the return values of the `USER()`, `CURRENT_USER()`, `SESSION_USER()`, `SYSTEM_USER()`, `DATABASE()`, and `VERSION()` functions have the UTF-8 character set by default.

The server sets the `character_set_system` system variable to the name of the metadata character set:

```
mysql> SHOW VARIABLES LIKE 'character_set_system';
+----------------------+-------+
| Variable_name        | Value |
+----------------------+-------+
| character_set_system | utf8  |
+----------------------+-------+
```

Storage of metadata using Unicode does *not* mean that the server returns headers of columns and the results of `DESCRIBE` functions in the `character_set_system` character set by default. When you use `SELECT column1 FROM t`, the name `column1` itself is returned from the server to the client in the character set determined by the value of the `character_set_results` system variable, which has a default value of `latin1`. If you want the server to pass metadata results back in a different character set, use the `SET NAMES` statement to force the server to perform character set conversion. `SET NAMES` sets the `character_set_results` and other related system variables. (See Section 1.4, "Connection Character Sets and Collations".) Alternatively, a client program can perform the conversion after receiving the result from the server. It is more efficient for the client perform the conversion, but this option is not always available for all clients.

If `character_set_results` is set to `NULL`, no conversion is performed and the server returns metadata using its original character set (the set indicated by `character_set_system`).

Error messages returned from the server to the client are converted to the client character set automatically, as with metadata.

If you are using (for example) the `USER()` function for comparison or assignment within a single statement, don't worry. MySQL performs some automatic conversion for you.

```
SELECT * FROM t1 WHERE USER() = latin1_column;
```

This works because the contents of `latin1_column` are automatically converted to UTF-8 before the comparison.

```
INSERT INTO t1 (latin1_column) SELECT USER();
```

This works because the contents of `USER()` are automatically converted to `latin1` before the assignment.

Although automatic conversion is not in the SQL standard, the SQL standard document does say that every character set is (in terms of supported characters) a "subset" of Unicode. Because it is a well-known principle that "what applies to a superset can apply to a subset," we believe that a collation for Unicode can apply for comparisons with non-Unicode strings. For more information about coercion of strings, see Section 1.6.5, "Special Cases Where Collation Determination Is Tricky".

# 1.12. Column Character Set Conversion

To convert a binary or nonbinary string column to use a particular character set, use `ALTER TABLE`. For successful conversion to occur, one of the following conditions must apply:

- If the column has a binary data type (`BINARY`, `VARBINARY`, `BLOB`), all the values that it contains must be encoded using a single character set (the character set you're converting the column to). If you use a binary column to store information in multiple character sets, MySQL has no way to know which values use which character set and cannot convert the data properly.

- If the column has a nonbinary data type (`CHAR`, `VARCHAR`, `TEXT`), its contents should be encoded in the column's character

set, not some other character set. If the contents are encoded in a different character set, you can convert the column to use a binary data type first, and then to a nonbinary column with the desired character set.

Suppose that a table `t` has a binary column named `col1` defined as `VARBINARY(50)`. Assuming that the information in the column is encoded using a single character set, you can convert it to a nonbinary column that has that character set. For example, if `col1` contains binary data representing characters in the `greek` character set, you can convert it as follows:

```
ALTER TABLE t MODIFY col1 VARCHAR(50) CHARACTER SET greek;
```

If your original column has a type of `BINARY(50)`, you could convert it to `CHAR(50)`, but the resulting values will be padded with `0x00` bytes at the end, which may be undesirable. To remove these bytes, use the `TRIM()` function:

```
UPDATE t SET col1 = TRIM(TRAILING 0x00 FROM col1);
```

Suppose that table `t` has a nonbinary column named `col1` defined as `CHAR(50) CHARACTER SET latin1` but you want to convert it to use `utf8` so that you can store values from many languages. The following statement accomplishes this:

```
ALTER TABLE t MODIFY col1 CHAR(50) CHARACTER SET utf8;
```

Conversion may be lossy if the column contains characters that are not in both character sets.

A special case occurs if you have old tables from MySQL 4.0 or earlier where a nonbinary column contains values that actually are encoded in a character set different from the server's default character set. For example, an application might have stored `sjis` values in a column, even though MySQL's default character set was `latin1`. It is possible to convert the column to use the proper character set but an additional step is required. Suppose that the server's default character set was `latin1` and `col1` is defined as `CHAR(50)` but its contents are `sjis` values. The first step is to convert the column to a binary data type, which removes the existing character set information without performing any character conversion:

```
ALTER TABLE t MODIFY col1 BLOB;
```

The next step is to convert the column to a nonbinary data type with the proper character set:

```
ALTER TABLE t MODIFY col1 CHAR(50) CHARACTER SET sjis;
```

This procedure requires that the table not have been modified already with statements such as `INSERT` or `UPDATE` after an upgrade to MySQL 4.1 or later. In that case, MySQL would store new values in the column using `latin1`, and the column will contain a mix of `sjis` and `latin1` values and cannot be converted properly.

If you specified attributes when creating a column initially, you should also specify them when altering the table with `ALTER TABLE`. For example, if you specified `NOT NULL` and an explicit `DEFAULT` value, you should also provide them in the `ALTER TABLE` statement. Otherwise, the resulting column definition will not include those attributes.

# 1.13. Character Sets and Collations That MySQL Supports

MySQL supports 70+ collations for 30+ character sets. This section indicates which character sets MySQL supports. There is one subsection for each group of related character sets. For each character set, the allowable collations are listed.

You can always list the available character sets and their default collations with the `SHOW CHARACTER SET` statement:

```
mysql> SHOW CHARACTER SET;
+----------+-----------------------------+---------------------+--------+
| Charset  | Description                 | Default collation   | Maxlen |
+----------+-----------------------------+---------------------+--------+
| big5     | Big5 Traditional Chinese    | big5_chinese_ci     |      2 |
| dec8     | DEC West European           | dec8_swedish_ci     |      1 |
| cp850    | DOS West European           | cp850_general_ci    |      1 |
| hp8      | HP West European            | hp8_english_ci      |      1 |
| koi8r    | KOI8-R Relcom Russian       | koi8r_general_ci    |      1 |
| latin1   | cp1252 West European        | latin1_swedish_ci   |      1 |
| latin2   | ISO 8859-2 Central European | latin2_general_ci   |      1 |
| swe7     | 7bit Swedish                | swe7_swedish_ci     |      1 |
| ascii    | US ASCII                    | ascii_general_ci    |      1 |
| ujis     | EUC-JP Japanese             | ujis_japanese_ci    |      3 |
| sjis     | Shift-JIS Japanese          | sjis_japanese_ci    |      2 |
| hebrew   | ISO 8859-8 Hebrew           | hebrew_general_ci   |      1 |
| tis620   | TIS620 Thai                 | tis620_thai_ci      |      1 |
| euckr    | EUC-KR Korean               | euckr_korean_ci     |      2 |
| koi8u    | KOI8-U Ukrainian            | koi8u_general_ci    |      1 |
| gb2312   | GB2312 Simplified Chinese    | gb2312_chinese_ci   |      2 |
| greek    | ISO 8859-7 Greek            | greek_general_ci    |      1 |
| cp1250   | Windows Central European    | cp1250_general_ci   |      1 |
| gbk      | GBK Simplified Chinese      | gbk_chinese_ci      |      2 |
| latin5   | ISO 8859-9 Turkish          | latin5_turkish_ci   |      1 |
| armscii8 | ARMSCII-8 Armenian          | armscii8_general_ci |      1 |
```

```
| utf8mb3  | UTF-8 Unicode              | utf8mb3_general_ci  |    3   |
| ucs2     | UCS-2 Unicode              | ucs2_general_ci     |    2   |
| cp866    | DOS Russian                | cp866_general_ci    |    1   |
| keybcs2  | DOS Kamenicky Czech-Slovak | keybcs2_general_ci  |    1   |
| macce    | Mac Central European       | macce_general_ci    |    1   |
| macroman | Mac West European          | macroman_general_ci |    1   |
| cp852    | DOS Central European       | cp852_general_ci    |    1   |
| latin7   | ISO 8859-13 Baltic         | latin7_general_ci   |    1   |
| utf8     | UTF-8 Unicode              | utf8_general_ci     |    4   |
| cp1251   | Windows Cyrillic           | cp1251_general_ci   |    1   |
| utf16    | UTF-16 Unicode             | utf16_general_ci    |    4   |
| cp1256   | Windows Arabic             | cp1256_general_ci   |    1   |
| cp1257   | Windows Baltic             | cp1257_general_ci   |    1   |
| utf32    | UTF-32 Unicode             | utf32_general_ci    |    4   |
| binary   | Binary pseudo charset      | binary              |    1   |
| geostd8  | GEOSTD8 Georgian           | geostd8_general_ci  |    1   |
| cp932    | SJIS for Windows Japanese  | cp932_japanese_ci   |    2   |
| eucjpms  | UJIS for Windows Japanese  | eucjpms_japanese_ci |    3   |
+----------+----------------------------+---------------------+--------+
```

The `utf8` character set prior to MySQL 6.0.4 encoded BMP characters using up to three bytes per character. In MySQL 6.0.4, this character set is renamed to `utf8mb3`, and a new `utf8` is added that encodes BMP and supplementary characters using up to four bytes per character. The `utf16` and `utf32` character sets also were added in MySQL 6.0.4.

In cases where a character set has multiple collations, it might not be clear which collation is most suitable for a given application. To avoid choosing the wrong collation, it can be helpful to perform some comparisons with representative data values to make sure that a given collation sorts values the way you expect.

Collation-Charts.Org is a useful site for information that shows how one collation compares to another.

# 1.13.1. Unicode Character Sets

MySQL 6.0 has several Unicode character sets:

- `ucs2`, the UCS-2 encoding of the Unicode character set using 16 bits per character

- `utf16`, the UTF-16 encoding for the Unicode character set; like `ucs2` but with an extension for supplementary characters

- `utf32`, the UTF-32 encoding for the Unicode character set using 32 bits per character

- `utf8`, a UTF-8 encoding of the Unicode character set using one to four bytes per character

- `utf8mb3`, a UTF-8 encoding of the Unicode character set using one to three bytes per character (known as `utf8` before MySQL 6.0)

You can store text in about 650 languages using these character sets. This section lists the collations available for each Unicode character set. For general information about the character sets, see Section 1.9, "Unicode Support".

A similar set of collations is available for each Unicode character set. These are shown in the following list, where *xxx* represents the character set name. For example, *xxx_danish_ci* represents the Danish collations, the specific names of which are `ucs2_danish_ci`, `utf16_danish_ci`, `ucs32_danish_ci`, `utf8_danish_ci`, and `utf8mb3_danish_ci`.

- *xxx_bin*

- *xxx_czech_ci*

- *xxx_danish_ci*

- *xxx_esperanto_ci*

- *xxx_estonian_ci*

- *xxx_general_ci* (default)

- *xxx_hungarian_ci*

- *xxx_icelandic_ci*

- *xxx_latvian_ci*

- *xxx_lithuanian_ci*

- *xxx_persian_ci*

- *xxx_polish_ci*

- *xxx_roman_ci*

- *xxx_romanian_ci*

- *xxx_sinhala_ci*

- *xxx_slovak_ci*

- *xxx_slovenian_ci*

- *xxx_spanish2_ci*

- *xxx_spanish_ci*

- *xxx_swedish_ci*

- *xxx_turkish_ci*

- *xxx_unicode_ci*

The collations for `utf8mb3`, `utf16`, and `utf32` were added in MySQL 6.0.4. In 6.0.4, the old (three-byte) `utf8` character set was renamed to `utf8mb3` and the new (four-byte) `utf8` character set was added. The collation IDs for the `utf8mb3` collations in MySQL 6.0 are the same as the IDs for the `utf8` collations in MySQL 5.1. In other words, the IDs are still associated with data for a particular encoding; it is the encoding name that has changed.

MySQL implements the *xxx_unicode_ci* collations according to the Unicode Collation Algorithm (UCA) described at [ht-tp://www.unicode.org/reports/tr10/](ht-tp://www.unicode.org/reports/tr10/). The collation uses the version-4.0.0 UCA weight keys: [ht-tp://www.unicode.org/Public/UCA/4.0.0/allkeys-4.0.0.txt](ht-tp://www.unicode.org/Public/UCA/4.0.0/allkeys-4.0.0.txt). Currently, the *xxx_unicode_ci* collations have only partial support for the Unicode Collation Algorithm. Some characters are not supported yet. Also, combining marks are not fully supported. This affects primarily Vietnamese, Yoruba, and some smaller languages such as Navajo. The following discussion uses `utf8_unicode_ci` for concreteness.

For any Unicode character set, operations performed using the *xxx_general_ci* collation are faster than those for the *xxx_unicode_ci* collation. For example, comparisons for the `utf8_general_ci` collation are faster, but slightly less correct, than comparisons for `utf8_unicode_ci`. The reason for this is that `utf8_unicode_ci` supports mappings such as expansions; that is, when one character compares as equal to combinations of other characters. For example, in German and some other languages "Ã#" is equal to "ss". `utf8_unicode_ci` also supports contractions and ignorable characters. `utf8_general_ci` is a legacy collation that does not support expansions, contractions, or ignorable characters. It can make only one-to-one comparisons between characters.

To further illustrate, the following equalities hold in both `utf8_general_ci` and `utf8_unicode_ci` (for the effect this has in comparisons or when doing searches, see Section 1.6.7, "Examples of the Effect of Collation"):

```
Ã# = A
Ã# = O
Ã# = U
```

A difference between the collations is that this is true for `utf8_general_ci`:

```
Ã# = s
```

Whereas this is true for `utf8_unicode_ci`:

```
Ã# = ss
```

MySQL implements language-specific collations for the `utf8` character set only if the ordering with `utf8_unicode_ci` does not work well for a language. For example, `utf8_unicode_ci` works fine for German and French, so there is no need to create special `utf8` collations for these two languages.

`utf8_general_ci` also is satisfactory for both German and French, except that "Ã#" is equal to "s", and not to "ss". If this is acceptable for your application, then you should use `utf8_general_ci` because it is faster. Otherwise, use `utf8_unicode_ci` because it is more accurate.

`utf8_swedish_ci`, like other `utf8` language-specific collations, is derived from `utf8_unicode_ci` with additional language rules. For example, in Swedish, the following relationship holds, which is not something expected by a German or French speaker:

```
Ã# = Y < Ã#
```

The `xxx_spanish_ci` and `xxx_spanish2_ci` collations correspond to modern Spanish and traditional Spanish, respectively. In both collations, "ñ" (n-tilde) is a separate letter between "n" and "o". In addition, for traditional Spanish, "ch" is a separate letter between "c" and "d", and "ll" is a separate letter between "l" and "m"

In the `xxx_roman_ci` collations, `I` and `J` compare as equal, and `U` and `V` compare as equal.

For all Unicode collations except the "binary" (`_bin`) collations, MySQL performs a table lookup to find a character's collating weight. This weight can be displayed using the `WEIGHT_STRING()` function. (See String Functions.) But if a character is not in the table (for example, because it is a "new" character), collating weight determination becomes more complex:

- For BMP characters in general collations (for example, `utf8_general_ci`), weight = code point.

- For BMP characters in UCA collations (for example, `utf8_unicode_ci`), the following algorithm applies:

```
if (code >= 0x3400 && code <= 0x4DB5)
  base= 0xFB80; /* CJK Ideograph Extension */
else if (code >= 0x4E00 && code <= 0x9FA5)
  base= 0xFB40; /* CJK Ideograph */
else
  base= 0xFBC0; /* All other characters */
aaaa= base +  (code >> 15);
bbbb= (code & 0x7FFF) | 0x8000;
```

The result is a sequence of two collating elements, `aaaa` followed by `bbbb`. For example:

```
mysql> SELECT HEX(WEIGHT_STRING(_ucs2 0x04CF COLLATE ucs2_unicode_ci));
+--------------------------------------------------------+
| HEX(WEIGHT_STRING(_ucs2 0x04CF COLLATE ucs2_unicode_ci)) |
+--------------------------------------------------------+
| FBC084CF                                               |
+--------------------------------------------------------+
```

Thus, `U+04cf CYRILLIC SMALL LETTER PALOCHKA` currently is, with all UCA collations, greater than `U+04c0 CYRILLIC LETTER PALOCHKA`. Eventually, after further collation tuning, all palochkas will sort together.

- For supplementary characters in general collations, the weight is the weight for `0xfffd REPLACEMENT CHARACTER`. For supplementary characters in UCA collations, their collating weight is `0xfffd`. That is, to MySQL, all supplementary characters are equal to each other, and greater than almost all BMP characters.

  An example with Deseret characters and `COUNT(DISTINCT)`:

```
CREATE TABLE t (s1 VARCHAR(5) CHARACTER SET utf32 COLLATE utf32_unicode_ci);
INSERT INTO t VALUES (0xfffd);    /* REPLACEMENT CHARACTER */
INSERT INTO t VALUES (0x010412); /* DESERET CAPITAL LETTER BEE */
INSERT INTO t VALUES (0x010413); /* DESERET CAPITAL LETTER TEE */
SELECT COUNT(DISTINCT s1) FROM t;
```

  The result is 1, because Deseret Bee = Deseret Tee = Replacement Character, in the MySQL Unicode collation.

  An example with cuneiform characters and `WEIGHT_STRING()`:

```
/*
The four characters in the INSERT string are
00000041  # LATIN CAPITAL LETTER A
0001218F # CUNEIFORM SIGN KAB
000121A7 # CUNEIFORM SIGN KISH
00000042  # LATIN CAPITAL LETTER B
*/
CREATE TABLE t (s1 CHAR(4) CHARACTER SET utf32 COLLATE utf32_unicode_ci);
INSERT INTO t VALUES (0x000000410001218f000121a700000042);
SELECT HEX(WEIGHT_STRING(s1)) FROM t;
```

  The result is:

```
0E33 FFFD FFFD 0E4A
```

  `0E33` and `0E4A` are primary weights as in UCA 4.0.0. `FFFD` is the weight for KAB and also for KISH.

The current rule that all supplementary characters are equal to each other is non-optimal. However, we don't expect the rule to cause trouble. These characters are very rare, so it will be very rare that a multi-character string consists entirely of supplementary characters. In Japan, since the supplementary characters are obscure Kanji ideographs, the typical user doesn't care what order they're in, anyway. If you really want to get rows sorted by MySQL's rule and secondarily by code point value, it is easy:

```
ORDER BY s1 COLLATE utf32_unicode_ci, s1 COLLATE utf32_bin
```

**The utf16_bin Collation**

There is a difference between "ordering by the character's code value" and "ordering by the character's binary representation," a difference that appears only with `utf16_bin`, because of surrogates.

Suppose that `utf16_bin` (the binary collation for `utf16`) was a binary comparison "byte by byte" rather than "character by character." If that were so, then the order of characters in `utf16_bin` would differ from the order in `utf8_bin`. For example, here is a chart showing two rare characters. The first character is in the range `E000-FFFF`, so it is greater than a surrogate but less than a supplementary. The second character is a supplementary.

```
Code point  Character                  utf8       utf16
----------  ---------                  ----       -----
0FF9D       HALFWIDTH KATAKANA LETTER N  EF BE 9D    FF 9D
10384       UGARITIC LETTER DELTA        F0 90 8E 84  D8 00 DF 84
```

The two characters in the chart are in order by code point value, because `0xff9d < 0x10384`. And they are in order by `utf8` value, because `0xef < 0xf0`. But they are not in order by `utf16` value, if we use byte-by-byte comparison, because `0xff > 0xd8`.

So MySQL's `utf16_bin` collation is not "byte by byte." It is "by code point." When MySQL sees a supplementary-character encoding in utf16, it converts to the character's code-point value, and then compares. Therefore `utf8_bin` and `utf16_bin` are the same ordering. This is consistent with the SQL:2008 standard requirement for a UCS_BASIC collation: "UCS_BASIC is a collation in which the ordering is determined entirely by the Unicode scalar values of the characters in the strings being sorted. It is applicable to the UCS character repertoire. Since every character repertoire is a subset of the UCS repertoire, the UCS_BASIC collation is potentially applicable to every character set. NOTE 11 — The Unicode scalar value of a character is its code point treated as an unsigned integer."

If the character set is `ucs2`, then comparison is byte-by-byte, but `ucs2` strings shouldn't contain surrogates, anyway.

For additional information about Unicode collations in MySQL, see Collation-Charts.Org (utf8).

# 1.13.2. West European Character Sets

Western European character sets cover most West European languages, such as French, Spanish, Catalan, Basque, Portuguese, Italian, Albanian, Dutch, German, Danish, Swedish, Norwegian, Finnish, Faroese, Icelandic, Irish, Scottish, and English.

- `ascii` (US ASCII) collations:
  - `ascii_bin`
  - `ascii_general_ci` (default)
- `cp850` (DOS West European) collations:
  - `cp850_bin`
  - `cp850_general_ci` (default)
- `dec8` (DEC Western European) collations:
  - `dec8_bin`
  - `dec8_swedish_ci` (default)
- `hp8` (HP Western European) collations:
  - `hp8_bin`
  - `hp8_english_ci` (default)
- `latin1` (cp1252 West European) collations:
  - `latin1_bin`
  - `latin1_danish_ci`
  - `latin1_general_ci`

- `latin1_general_cs`

- `latin1_german1_ci`

- `latin1_german2_ci`

- `latin1_spanish_ci`

- `latin1_swedish_ci` (default)

`latin1` is the default character set. MySQL's `latin1` is the same as the Windows `cp1252` character set. This means it is the same as the official `ISO 8859-1` or IANA (Internet Assigned Numbers Authority) `latin1`, except that IANA `latin1` treats the code points between `0x80` and `0x9f` as "undefined," whereas `cp1252`, and therefore MySQL's `latin1`, assign characters for those positions. For example, `0x80` is the Euro sign. For the "undefined" entries in `cp1252`, MySQL translates `0x81` to Unicode `0x0081`, `0x8d` to `0x008d`, `0x8f` to `0x008f`, `0x90` to `0x0090`, and `0x9d` to `0x009d`.

The `latin1_swedish_ci` collation is the default that probably is used by the majority of MySQL customers. Although it is frequently said that it is based on the Swedish/Finnish collation rules, there are Swedes and Finns who disagree with this statement.

The `latin1_german1_ci` and `latin1_german2_ci` collations are based on the DIN-1 and DIN-2 standards, where DIN stands for *Deutsches Institut fÃ¼r Normung* (the German equivalent of ANSI). DIN-1 is called the "dictionary collation" and DIN-2 is called the "phone book collation." For an example of the effect this has in comparisons or when doing searches, see Section 1.6.7, "Examples of the Effect of Collation".

- `latin1_german1_ci` (dictionary) rules:

```
Ã# = A
Ã# = O
Ã# = U
Ã# = s
```

- `latin1_german2_ci` (phone-book) rules:

```
Ã# = AE
Ã# = OE
Ã# = UE
Ã# = ss
```

For an example of the effect this has in comparisons or when doing searches, see Section 1.6.7, "Examples of the Effect of Collation".

In the `latin1_spanish_ci` collation, "Ã±" (n-tilde) is a separate letter between "n" and "o".

- `macroman` (Mac West European) collations:

  - `macroman_bin`

  - `macroman_general_ci` (default)

- `swe7` (7bit Swedish) collations:

  - `swe7_bin`

  - `swe7_swedish_ci` (default)

For additional information about Western European collations in MySQL, see Collation-Charts.Org (ascii, cp850, dec8, hp8, latin1, macroman, swe7).

## 1.13.3. Central European Character Sets

MySQL provides some support for character sets used in the Czech Republic, Slovakia, Hungary, Romania, Slovenia, Croatia, Poland, and Serbia (Latin).

- `cp1250` (Windows Central European) collations:

  - `cp1250_bin`

  - `cp1250_croatian_ci`

- `cp1250_czech_cs`

- `cp1250_general_ci` (default)

- `cp1250_polish_ci`

- `cp852` (DOS Central European) collations:

  - `cp852_bin`

  - `cp852_general_ci` (default)

- `keybcs2` (DOS Kamenicky Czech-Slovak) collations:

  - `keybcs2_bin`

  - `keybcs2_general_ci` (default)

- `latin2` (ISO 8859-2 Central European) collations:

  - `latin2_bin`

  - `latin2_croatian_ci`

  - `latin2_czech_cs`

  - `latin2_general_ci` (default)

  - `latin2_hungarian_ci`

- `macce` (Mac Central European) collations:

  - `macce_bin`

  - `macce_general_ci` (default)

For additional information about Central European collations in MySQL, see Collation-Charts.Org (cp1250, cp852, keybcs2, latin2, macce).

## 1.13.4. South European and Middle East Character Sets

South European and Middle Eastern character sets supported by MySQL include Armenian, Arabic, Georgian, Greek, Hebrew, and Turkish.

- `armscii8` (ARMSCII-8 Armenian) collations:

  - `armscii8_bin`

  - `armscii8_general_ci` (default)

- `cp1256` (Windows Arabic) collations:

  - `cp1256_bin`

  - `cp1256_general_ci` (default)

- `geostd8` (GEOSTD8 Georgian) collations:

  - `geostd8_bin`

  - `geostd8_general_ci` (default)

- `greek` (ISO 8859-7 Greek) collations:

  - `greek_bin`

  - `greek_general_ci` (default)

- `hebrew` (ISO 8859-8 Hebrew) collations:

- hebrew_bin

- hebrew_general_ci (default)

- latin5 (ISO 8859-9 Turkish) collations:

  - latin5_bin

  - latin5_turkish_ci (default)

For additional information about South European and Middle Eastern collations in MySQL, see Collation-Charts.Org (armscii8, cp1256, geostd8, greek, hebrew, latin5).

## 1.13.5. Baltic Character Sets

The Baltic character sets cover Estonian, Latvian, and Lithuanian languages.

- cp1257 (Windows Baltic) collations:

  - cp1257_bin

  - cp1257_general_ci (default)

  - cp1257_lithuanian_ci

- latin7 (ISO 8859-13 Baltic) collations:

  - latin7_bin

  - latin7_estonian_cs

  - latin7_general_ci (default)

  - latin7_general_cs

For additional information about Baltic collations in MySQL, see Collation-Charts.Org (cp1257, latin7).

## 1.13.6. Cyrillic Character Sets

The Cyrillic character sets and collations are for use with Belarusian, Bulgarian, Russian, Ukrainian, and Serbian (Cyrillic) languages.

- cp1251 (Windows Cyrillic) collations:

  - cp1251_bin

  - cp1251_bulgarian_ci

  - cp1251_general_ci (default)

  - cp1251_general_cs

  - cp1251_ukrainian_ci

- cp866 (DOS Russian) collations:

  - cp866_bin

  - cp866_general_ci (default)

- koi8r (KOI8-R Relcom Russian) collations:

  - koi8r_bin

  - koi8r_general_ci (default)

- koi8u (KOI8-U Ukrainian) collations:

    - koi8u_bin

    - koi8u_general_ci (default)

For additional information about Cyrillic collations in MySQL, see Collation-Charts.Org (cp1251, cp866, koi8r, koi8u). ).

# 1.13.7. Asian Character Sets

The Asian character sets that we support include Chinese, Japanese, Korean, and Thai. These can be complicated. For example, the Chinese sets must allow for thousands of different characters. See Section 1.13.7.1, "The cp932 Character Set", for additional information about the cp932 and sjis character sets.

For answers to some common questions and problems relating support for Asian character sets in MySQL, see Chapter 9, *MySQL 6.0 FAQ — MySQL Chinese, Japanese, and Korean Character Sets*.

- big5 (Big5 Traditional Chinese) collations:

    - big5_bin

    - big5_chinese_ci (default)

- cp932 (SJIS for Windows Japanese) collations:

    - cp932_bin

    - cp932_japanese_ci (default)

- eucjpms (UJIS for Windows Japanese) collations:

    - eucjpms_bin

    - eucjpms_japanese_ci (default)

- euckr (EUC-KR Korean) collations:

    - euckr_bin

    - euckr_korean_ci (default)

- gb2312 (GB2312 Simplified Chinese) collations:

    - gb2312_bin

    - gb2312_chinese_ci (default)

- gbk (GBK Simplified Chinese) collations:

    - gbk_bin

    - gbk_chinese_ci (default)

- sjis (Shift-JIS Japanese) collations:

    - sjis_bin

    - sjis_japanese_ci (default)

- tis620 (TIS620 Thai) collations:

    - tis620_bin

    - tis620_thai_ci (default)

- ujis (EUC-JP Japanese) collations:

    - ujis_bin

    - ujis_japanese_ci (default)

The `big5_chinese_ci` collation sorts on number of strokes.

For additional information about Asian collations in MySQL, see Collation-Charts.Org (big5, cp932, eucjpms, euckr, gb2312, gbk, sjis, tis620, ujis).

## 1.13.7.1. The `cp932` Character Set

**Why is `cp932` needed?**

In MySQL, the `sjis` character set corresponds to the `Shift_JIS` character set defined by IANA, which supports JIS X0201 and JIS X0208 characters. (See http://www.iana.org/assignments/character-sets.)

However, the meaning of "SHIFT JIS" as a descriptive term has become very vague and it often includes the extensions to `Shift_JIS` that are defined by various vendors.

For example, "SHIFT JIS" used in Japanese Windows environments is a Microsoft extension of `Shift_JIS` and its exact name is `Microsoft Windows Codepage : 932` or `cp932`. In addition to the characters supported by `Shift_JIS`, `cp932` supports extension characters such as NEC special characters, NEC selected — IBM extended characters, and IBM extended characters.

Many Japanese users have experienced problems using these extension characters. These problems stem from the following factors:

- MySQL automatically converts character sets.

- Character sets are converted via Unicode (`ucs2`).

- The `sjis` character set does not support the conversion of these extension characters.

- There are several conversion rules from so-called "SHIFT JIS" to Unicode, and some characters are converted to Unicode differently depending on the conversion rule. MySQL supports only one of these rules (described later).

The MySQL `cp932` character set is designed to solve these problems.

Because MySQL supports character set conversion, it is important to separate IANA `Shift_JIS` and `cp932` into two different character sets because they provide different conversion rules.

**How does `cp932` differ from `sjis`?**

The `cp932` character set differs from `sjis` in the following ways:

- `cp932` supports NEC special characters, NEC selected — IBM extended characters, and IBM selected characters.

- Some `cp932` characters have two different code points, both of which convert to the same Unicode code point. When converting from Unicode back to `cp932`, one of the code points must be selected. For this "round trip conversion," the rule recommended by Microsoft is used. (See http://support.microsoft.com/kb/170559/EN-US/.)

  The conversion rule works like this:

  - If the character is in both JIS X 0208 and NEC special characters, use the code point of JIS X 0208.

  - If the character is in both NEC special characters and IBM selected characters, use the code point of NEC special characters.

  - If the character is in both IBM selected characters and NEC selected — IBM extended characters, use the code point of IBM extended characters.

  The table shown at http://www.microsoft.com/globaldev/reference/dbcs/932.htm provides information about the Unicode values of `cp932` characters. For `cp932` table entries with characters under which a four-digit number appears, the number represents the corresponding Unicode (`ucs2`) encoding. For table entries with an underlined two-digit value appears, there is a range of `cp932` character values that begin with those two digits. Clicking such a table entry takes you to a page that displays the Unicode value for each of the `cp932` characters that begin with those digits.

  The following links are of special interest. They correspond to the encodings for the following sets of characters:

- NEC special characters:

  ```
  http://www.microsoft.com/globaldev/reference/dbcs/932/932_87.htm
  ```

- NEC selected — IBM extended characters:

```
http://www.microsoft.com/globaldev/reference/dbcs/932/932_ED.htm
http://www.microsoft.com/globaldev/reference/dbcs/932/932_EE.htm
```

- IBM selected characters:

```
http://www.microsoft.com/globaldev/reference/dbcs/932/932_FA.htm
http://www.microsoft.com/globaldev/reference/dbcs/932/932_FB.htm
http://www.microsoft.com/globaldev/reference/dbcs/932/932_FC.htm
```

- `cp932` supports conversion of user-defined characters in combination with `eucjpms`, and solves the problems with `sjis`/`ujis` conversion. For details, please refer to http://www.opengroup.or.jp/jvc/cde/sjis-euc-e.html.

For some characters, conversion to and from `ucs2` is different for `sjis` and `cp932`. The following tables illustrate these differences.

Conversion to `ucs2`:

| `sjis`/`cp932` Value | `sjis` -> `ucs2` Conversion | `cp932` -> `ucs2` Conversion |
| --- | --- | --- |
| 5C | 005C | 005C |
| 7E | 007E | 007E |
| 815C | 2015 | 2015 |
| 815F | 005C | FF3C |
| 8160 | 301C | FF5E |
| 8161 | 2016 | 2225 |
| 817C | 2212 | FF0D |
| 8191 | 00A2 | FFE0 |
| 8192 | 00A3 | FFE1 |
| 81CA | 00AC | FFE2 |

Conversion from `ucs2`:

| `ucs2` value | `ucs2` -> `sjis` Conversion | `ucs2` -> `cp932` Conversion |
| --- | --- | --- |
| 005C | 815F | 5C |
| 007E | 7E | 7E |
| 00A2 | 8191 | 3F |
| 00A3 | 8192 | 3F |
| 00AC | 81CA | 3F |
| 2015 | 815C | 815C |
| 2016 | 8161 | 3F |
| 2212 | 817C | 3F |
| 2225 | 3F | 8161 |
| 301C | 8160 | 3F |
| FF0D | 3F | 817C |
| FF3C | 3F | 815F |
| FF5E | 3F | 8160 |
| FFE0 | 3F | 8191 |
| FFE1 | 3F | 8192 |
| FFE2 | 3F | 81CA |

Users of any Japanese character sets should be aware that using `--character-set-client-handshake` (or `--skip-character-set-client-handshake`) has an important effect. See Server Command Options.

# Chapter 2. The Character Set Used for Data and Sorting

By default, MySQL uses the `latin1` (cp1252 West European) character set and the `latin1_swedish_ci` collation that sorts according to Swedish/Finnish rules. These defaults are suitable for the United States and most of Western Europe.

All MySQL binary distributions are compiled with `--with-extra-charsets=complex`. This adds code to all standard programs that enables them to handle `latin1` and all multi-byte character sets within the binary. Other character sets are loaded from a character-set definition file when needed.

The character set determines what characters are allowed in identifiers. The collation determines how strings are sorted by the `ORDER BY` and `GROUP BY` clauses of the `SELECT` statement.

You can change the default server character set and collation with the `--character-set-server` and `--collation-server` options when you start the server. The collation must be a legal collation for the default character set. (Use the `SHOW COLLATION` statement to determine which collations are available for each character set.) See Server Command Options.

The character sets available depend on the `--with-charset=charset_name` and `--with-extra-charsets=list-of-charsets | complex | all | none` options to `configure`, and the character set configuration files listed in `SHAREDIR/charsets/Index`. See Typical `configure` Options.

If you change the character set when running MySQL, that may also change the sort order. Consequently, you must run `myisamchk -r -q --set-collation=collation_name` on all `MyISAM` tables, or your indexes may not be ordered correctly.

When a client connects to a MySQL server, the server indicates to the client what the server's default character set is. The client switches to this character set for this connection.

You should use `mysql_real_escape_string()` when escaping strings for an SQL query. `mysql_real_escape_string()` is identical to the old `mysql_escape_string()` function, except that it takes the `MYSQL` connection handle as the first parameter so that the appropriate character set can be taken into account when escaping characters.

If the client is compiled with paths that differ from where the server is installed and the user who configured MySQL didn't include all character sets in the MySQL binary, you must tell the client where it can find the additional character sets it needs if the server runs with a different character set from the client. You can do this by specifying a `--character-sets-dir` option to indicate the path to the directory in which the dynamic MySQL character sets are stored. For example, you can put the following in an option file:

```
[client]
character-sets-dir=/usr/local/mysql/share/mysql/charsets
```

You can force the client to use specific character set as follows:

```
[client]
default-character-set=charset_name
```

This is normally unnecessary, however.

## 2.1. Using the German Character Set

In MySQL 6.0, character set and collation are specified separately. This means that if you want German sort order, you should select the `latin1` character set and either the `latin1_german1_ci` or `latin1_german2_ci` collation. For example, to start the server with the `latin1_german1_ci` collation, use the `--character-set-server=latin1` and `--collation-server=latin1_german1_ci` options.

For information on the differences between these two collations, see Section 1.13.2, "West European Character Sets".

# Chapter 3. Setting the Error Message Language

By default, `mysqld` produces error messages in English, but they can also be displayed in any of these other languages: Czech, Danish, Dutch, Estonian, French, German, Greek, Hungarian, Italian, Japanese, Korean, Norwegian, Norwegian-ny, Polish, Portuguese, Romanian, Russian, Slovak, Spanish, or Swedish.

To start `mysqld` with a particular language for error messages, use the `--language` or `-L` option. The option value can be a language name or the full path to the error message file. For example:

```
shell> mysqld --language=swedish
```

Or:

```
shell> mysqld --language=/usr/local/share/swedish
```

The language name should be specified in lowercase.

By default, the language files are located in the `share/LANGUAGE` directory under the MySQL base directory.

You can also change the content of the error messages produced by the server. Details can be found in the MySQL Internals manual, available at http://forge.mysql.com/wiki/MySQL_Internals_Error_Messages. If you upgrade to a newer version of MySQL after changing the error messages, remember to repeat your changes after the upgrade.

# Chapter 4. Adding a New Character Set

This section discusses the procedure for adding a new character set to MySQL. You must have a MySQL source distribution to use these instructions. The proper procedure depends on whether the character set is simple or complex:

- If the character set does not need to use special string collating routines for sorting and does not need multi-byte character support, it is simple.

- If the character set needs either of those features, it is complex.

For example, `greek` and `swe7` are simple character sets, whereas `big5` and `czech` are complex character sets.

In the following instructions, *MYSET* represents the name of the character set that you want to add.

1. Add a `<charset>` element for *MYSET* to the `sql/share/charsets/Index.xml` file. Use the existing contents in the file as a guide to adding new contents.

   The `<charset>` element must list all the collations for the character set. These must include at least a binary collation and a default collation. The default collation is usually named using a suffix of `general_ci` (general, case insensitive). It is possible for the binary collation to be the default collation, but usually they are different. The default collation should have a `primary` flag. The binary collation should have a `binary` flag.

   You must assign a unique ID number to each collation. As of MySQL 6.0.8, the range of IDs from 1024 to 2047 is reserved for user-defined collations. Before 6.0.8, the ID must be chosen from the range 1 to 254. To find the maximum of the currently used collation IDs, use this query:

   ```
   SELECT MAX(ID) FROM INFORMATION_SCHEMA.COLLATIONS;
   ```

2. This step depends on whether you are adding a simple or complex character set. A simple character set requires only a configuration file, whereas a complex character set requires C source file that defines collation functions, multi-byte functions, or both.

   For a simple character set, create a configuration file, *MYSET*`.xml`, that describes the character set properties. Create this file in the `sql/share/charsets` directory. (You can use a copy of `latin1.xml` as the basis for this file.) The syntax for the file is very simple:

   - Comments are written as ordinary XML comments (`<!-- text -->`).

   - Words within `<map>` array elements are separated by arbitrary amounts of whitespace.

   - Each word within `<map>` array elements must be a number in hexadecimal format.

   - The `<map>` array element for the `<ctype>` element has 257 words. The other `<map>` array elements after that have 256 words. See Section 4.1, "The Character Definition Arrays".

   - For each collation listed in the `<charset>` element for the character set in `Index.xml`, *MYSET*`.xml` must contain a `<collation>` element that defines the character ordering.

   For a complex character set, create a C source file that describes the character set properties and defines the support routines necessary to properly perform operations on the character set:

   a. Create the file `ctype-`*MYSET*`.c` in the `strings` directory. Look at one of the existing `ctype-*.c` files (such as `ctype-big5.c`) to see what needs to be defined. The arrays in your file must have names like `ctype_`*MYSET*, `to_lower_`*MYSET*, and so on. These correspond to the arrays for a simple character set. See Section 4.1, "The Character Definition Arrays".

   b. For each collation listed in the `<charset>` element for the character set in `Index.xml`, the `ctype-`*MYSET*`.c` file must provide an implementation of the collation.

   c. If you need string collating functions, see Section 4.2, "String Collating Support".

   d. If you need multi-byte character support, see Section 4.3, "Multi-Byte Character Support".

3. Follow these steps to modify the configuration information. Use the existing configuration information as a guide to adding information for *MYSYS*. The example here assumes that the character set has default and binary collations, but more lines will be needed if *MYSET* has additional collations.

a. Edit `mysys/charset-def.c`, and "register" the collations for the new character set.

Add these lines to the "declaration" section:

```
#ifdef HAVE_CHARSET_MYSET
extern CHARSET_INFO my_charset_MYSET_general_ci;
extern CHARSET_INFO my_charset_MYSET_bin;
#endif
```

Add these lines to the "registration" section:

```
#ifdef HAVE_CHARSET_MYSET
  add_compiled_collation(&my_charset_MYSET_general_ci);
  add_compiled_collation(&my_charset_MYSET_bin);
#endif
```

b. If the character set uses `ctype-MYSET.c`, edit `strings/Makefile.am` and add `ctype-MYSET.c` to each definition of the `CSRCS` variable, and to the `EXTRA_DIST` variable.

c. If the character set uses `ctype-MYSET.c`, edit `libmysql/Makefile.shared` and add `ctype-MYSET.lo` to the `mystringsobjects` definition.

d. Edit `config/ac-macros/character_sets.m4`:

   i. Add `MYSET` to one of the `define(CHARSETS_AVAILABLE...)` lines in alphabetic order.

   ii. Add `MYSET` to `CHARSETS_COMPLEX`. This is needed even for simple character sets, or `configure` will not recognize `--with-charset=MYSET`.

   iii. Add `MYSET` to the first `case` control structure. Omit the `USE_MB` and `USE_MB_IDENT` lines for 8-bit character sets.

```
MYSET)
    AC_DEFINE(HAVE_CHARSET_MYSET, 1, [Define to enable charset MYSET])
    AC_DEFINE([USE_MB], 1, [Use multi-byte character routines])
    AC_DEFINE(USE_MB_IDENT, 1)
    ;;
```

   iv. Add `MYSET` to the second `case` control structure:

```
MYSET)
    default_charset_default_collation="MYSET_general_ci"
    default_charset_collations="MYSET_general_ci MYSET_bin"
    ;;
```

4. Reconfigure, recompile, and test.

# 4.1. The Character Definition Arrays

Each simple character set has a configuration file located in the `sql/share/charsets` directory. The file is named `MYSET.xml`. It uses `<map>` array elements to list character set properties. `<map>` elements appear within these elements:

- `<ctype>` defines attributes for each character

- `<lower>` and `<upper>` list the lowercase and uppercase characters

- `<unicode>` maps 8-bit character values to Unicode values

- `<collation>` elements indicate character ordering for comparisons and sorts, one element per collation (binary collations need no `<map>` element because the character codes themselves provide the ordering)

For a complex character set as implemented in a `ctype-MYSET.c` file in the `strings` directory, there are corresponding arrays: `ctype_MYSET[]`, `to_lower_MYSET[]`, and so forth. Not every complex character set has all of the arrays. See the existing `ctype-*.c` files for examples. See the `CHARSET_INFO.txt` file in the `strings` directory for additional information.

The `ctype` array is indexed by character value + 1 and has 257 elements. This is an old legacy convention for handling `EOF`. The other arrays are indexed by character value and have 256 elements.

`ctype` array elements are bit values. Each element describes the attributes of a single character in the character set. Each attribute is associated with a bitmask, as defined in `include/m_ctype.h`:

```
#define _MY_U    01      /* Upper case */
#define _MY_L    02      /* Lower case */
#define _MY_NMR  04      /* Numeral (digit) */
#define _MY_SPC  010     /* Spacing character */
#define _MY_PNT  020     /* Punctuation */
#define _MY_CTR  040     /* Control character */
#define _MY_B    0100    /* Blank */
#define _MY_X    0200    /* heXadecimal digit */
```

The `ctype` value for a given character should be the union of the applicable bitmask values that describe the character. For example, `'A'` is an uppercase character (`_MY_U`) as well as a hexadecimal digit (`_MY_X`), so its `ctype` value should be defined like this:

```
ctype['A'+1] = _MY_U | _MY_X = 01 | 0200 = 0201
```

The bitmask values in `m_ctype.h` are octal values, but the elements of the `ctype` array in *MYSET*`.xml` should be written as hexadecimal values.

The `lower` and `upper` arrays hold the lowercase and uppercase characters corresponding to each member of the character set. For example:

```
lower['A'] should contain 'a'
upper['a'] should contain 'A'
```

Each `collation` array is a map indicating how characters should be ordered for comparison and sorting purposes. MySQL sorts characters based on the values of this information. In some cases, this is the same as the `upper` array, which means that sorting is case-insensitive. For more complicated sorting rules (for complex character sets), see the discussion of string collating in Section 4.2, "String Collating Support".

# 4.2. String Collating Support

For simple character sets, sorting rules are specified in the *MYSET*`.xml` configuration file using `<map>` array elements within `<collation>` elements. If the sorting rules for your language are too complex to be handled with simple arrays, you need to define string collating functions in the `ctype-`*MYSET*`.c` source file in the `strings` directory.

The existing character sets provide the best documentation and examples to show how these functions are implemented. Look at the `ctype-*.c` files in the `strings` directory, such as the files for the `big5`, `czech`, `gbk`, `sjis`, and `tis160` character sets. Take a look at the `MY_COLLATION_HANDLER` structures to see how they are used, and see the `CHARSET_INFO.txt` file in the `strings` directory for additional information.

# 4.3. Multi-Byte Character Support

If you want to add support for a new character set that includes multi-byte characters, you need to use multi-byte character functions in the `ctype-`*MYSET*`.c` source file in the `strings` directory.

The existing character sets provide the best documentation and examples to show how these functions are implemented. Look at the `ctype-*.c` files in the `strings` directory, such as the files for the `euc_kr`, `gb2312`, `gbk`, `sjis`, and `ujis` character sets. Take a look at the `MY_CHARSET_HANDLER` structures to see how they are used, and see the `CHARSET_INFO.txt` file in the `strings` directory for additional information.

# Chapter 5. How to Add a New Collation to a Character Set

A collation is a set of rules that defines how to compare and sort character strings. Each collation in MySQL belongs to a single character set. Every character set has at least one collation, and most have two or more collations.

A collation orders characters based on weights. Each character in a character set maps to a weight. Characters with equal weights compare as equal, and characters with unequal weights compare according to the relative magnitude of their weights.

The `WEIGHT_STRING()` function can be used to see the weights for the characters in a string. The value that it returns to indicate weights is a binary string, so it is convenient to use `HEX(WEIGHT_STRING(str))` to display the weights in printable form. The following example shows that weights do not differ for lettercase for the letters in `'AaBb'` it if is a nonbinary case-insensitive string, but do differ if it is a binary string:

```
mysql> SELECT HEX(WEIGHT_STRING('AaBb' COLLATE latin1_swedish_ci));
+-----------------------------------------------------+
| HEX(WEIGHT_STRING('AaBb' COLLATE latin1_swedish_ci)) |
+-----------------------------------------------------+
| 41414242                                            |
+-----------------------------------------------------+
mysql> SELECT HEX(WEIGHT_STRING(BINARY 'AaBb'));
+----------------------------------+
| HEX(WEIGHT_STRING(BINARY 'AaBb')) |
+----------------------------------+
| 41614262                         |
+----------------------------------+
```

MySQL supports several collation implementations, as discussed in Section 5.1, "Collation Implementation Types". Some of these can be added to MySQL without recompiling:

- Simple collations for 8-bit character sets

- UCA-based collations for Unicode character sets

- Binary (`xxx_bin`) collations

The following discussion describes how to add collations of the first two types to existing character sets. All existing character sets already have a binary collation, so there is no need here to describe how to add one.

Summary of the procedure for adding a new collation:

1. Choose a collation ID

2. Add configuration information that names the collation and describes the character-ordering rules

3. Restart the server

4. Verify that the collation is present

The instructions here cover only collations that can be added without recompiling MySQL. To add a collation that does require recompiling (as implemented by means of functions in a C source file), use the instructions in Chapter 4, *Adding a New Character Set*. However, instead of adding all the information required for a complete character set, just modify the appropriate files for an existing character set. That is, based on what is already present for the character set's current collations, add new data structures, functions, and configuration information for the new collation. For an example, see the MySQL Blog article in the following list of additional resources.

**Additional resources**

- The Unicode Collation Algorithm (UCA) specification: http://www.unicode.org/reports/tr10/

- The Locale Data Markup Language (LDML) specification: http://www.unicode.org/reports/tr35/

- MySQL University session "How to Add a Collation": http://forge.mysql.com/wiki/How_to_Add_a_Collation

- MySQL Blog article "Instructions for adding a new Unicode collation": http://blogs.mysql.com/peterg/2008/05/19/instructions-for-adding-a-new-unicode-collation/

# 5.1. Collation Implementation Types

MySQL implements several types of collations:

**Simple collations for 8-bit character sets**

This kind of collation is implemented using an array of 256 weights that defines a one-to-one mapping from character codes to weights. `latin1_swedish_ci` is an example. It is a case-insensitive collation, so the uppercase and lowercase versions of a character have the same weights and they compare as equal.

```
mysql> SET NAMES 'latin1' COLLATE 'latin1_swedish_ci';
Query OK, 0 rows affected (0.01 sec)
mysql> SELECT HEX(WEIGHT_STRING('a')), HEX(WEIGHT_STRING('A'));
+-------------------------+-------------------------+
| HEX(WEIGHT_STRING('a')) | HEX(WEIGHT_STRING('A')) |
+-------------------------+-------------------------+
| 41                      | 41                      |
+-------------------------+-------------------------+
1 row in set (0.01 sec)
mysql> SELECT 'a' = 'A';
+-----------+
| 'a' = 'A' |
+-----------+
|         1 |
+-----------+
1 row in set (0.12 sec)
```

**Complex collations for 8-bit character sets**

This kind of collation is implemented using functions in a C source file that define how to order characters, as described in Chapter 4, *Adding a New Character Set*.

**Collations for non-Unicode multi-byte character sets**

For this type of collation, 8-bit (single-byte) and multi-byte characters are handled differently. For 8-bit characters, character codes map to weights in case-insensitive fashion. (For example, the single-byte characters `'a'` and `'A'` both have a weight of `0x41`.) For multi-byte characters, there are two types of relationship between character codes and weights:

- Weights equal character codes. `sjis_japanese_ci` is an example of this kind of collation. The multi-byte character `'ぢ'` has a character code of `0x82C0`, and the weight is also `0x82C0`.

```
mysql> CREATE TABLE t1
    -> (c1 VARCHAR(2) CHARACTER SET sjis COLLATE sjis_japanese_ci);
Query OK, 0 rows affected (0.01 sec)
mysql> INSERT INTO t1 VALUES ('a'),('A'),(0x82C0);
Query OK, 3 rows affected (0.00 sec)
Records: 3  Duplicates: 0  Warnings: 0
mysql> SELECT c1, HEX(c1), HEX(WEIGHT_STRING(c1)) FROM t1;
+------+---------+-----------------------+
| c1   | HEX(c1) | HEX(WEIGHT_STRING(c1)) |
+------+---------+-----------------------+
| a    | 61      | 41                    |
| A    | 41      | 41                    |
| ぢ   | 82C0    | 82C0                  |
+------+---------+-----------------------+
3 rows in set (0.00 sec)
```

- Character codes map one-to-one to weights, but a code is not necessarily equal to the weight. `gbk_chinese_ci` is an example of this kind of collation. The multi-byte character `'腤'` has a character code of `0x81B0` but a weight of `0xC286`.

```
mysql> CREATE TABLE t1
    -> (c1 VARCHAR(2) CHARACTER SET gbk COLLATE gbk_chinese_ci);
Query OK, 0 rows affected (0.33 sec)
mysql> INSERT INTO t1 VALUES ('a'),('A'),(0x81B0);
Query OK, 3 rows affected (0.00 sec)
Records: 3  Duplicates: 0  Warnings: 0
mysql> SELECT c1, HEX(c1), HEX(WEIGHT_STRING(c1)) FROM t1;
+------+---------+-----------------------+
| c1   | HEX(c1) | HEX(WEIGHT_STRING(c1)) |
+------+---------+-----------------------+
| a    | 61      | 41                    |
| A    | 41      | 41                    |
| 腤   | 81B0    | C286                  |
+------+---------+-----------------------+
3 rows in set (0.00 sec)
```

**Collations for Unicode multi-byte character sets**

Some of these collations are based on the Unicode Collation Algorithm (UCA), others are not.

Non-UCA collations have a one-to-one mapping from character code to weight. In MySQL, such collations are case insensitive and accent insensitive. `utf8_general_ci` is an example: `'a'`, `'A'`, `'Ã#'`, and `'Ã¡'` each have different character codes but all have a weight of `0x0041` and compare as equal.

```
mysql> SET NAMES 'utf8' COLLATE 'utf8_general_ci';
Query OK, 0 rows affected (0.00 sec)
mysql> CREATE TABLE t1
    -> (c1 CHAR(1) CHARACTER SET UTF8 COLLATE utf8_general_ci);
Query OK, 0 rows affected (0.01 sec)
mysql> INSERT INTO t1 VALUES ('a'),('A'),('Ã#'),('Ã¡');
Query OK, 4 rows affected (0.00 sec)
Records: 4  Duplicates: 0  Warnings: 0
mysql> SELECT c1, HEX(c1), HEX(WEIGHT_STRING(c1)) FROM t1;
+------+---------+-----------------------+
| c1   | HEX(c1) | HEX(WEIGHT_STRING(c1)) |
+------+---------+-----------------------+
| a    | 61      | 0041                  |
| A    | 41      | 0041                  |
| Ã#   | C380    | 0041                  |
| Ã¡   | C3A1    | 0041                  |
+------+---------+-----------------------+
4 rows in set (0.00 sec)
```

UCA-based collations in MySQL have these properties:

- If a character has weights, each weight uses 2 bytes (16 bits)

- A character may have zero weights (or an empty weight). In this case, the character is ignorable. Example: "U+0000 NULL" does not have a weight and is ignorable.

- A character may have one weight. Example: `'a'` has a weight of `0x0E33`.

```
mysql> SET NAMES 'utf8' COLLATE 'utf8_unicode_ci';
Query OK, 0 rows affected (0.05 sec)
mysql> SELECT HEX('a'), HEX(WEIGHT_STRING('a'));
+----------+------------------------+
| HEX('a') | HEX(WEIGHT_STRING('a')) |
+----------+------------------------+
| 61       | 0E33                   |
+----------+------------------------+
1 row in set (0.02 sec)
```

- A character may have many weights. This is an expansion. Example: The German letter `'Ã#'` (SZ LEAGUE, or SHARP S) has a weight of `0x0FEA0FEA`.

```
mysql> SET NAMES 'utf8' COLLATE 'utf8_unicode_ci';
Query OK, 0 rows affected (0.11 sec)
mysql> SELECT HEX('Ã#'), HEX(WEIGHT_STRING('Ã#'));
+-----------+-------------------------+
| HEX('Ã#') | HEX(WEIGHT_STRING('Ã#')) |
+-----------+-------------------------+
| C39F      | 0FEA0FEA                |
+-----------+-------------------------+
1 row in set (0.00 sec)
```

- Many characters may have one weight. This is a contraction. Example: `'ch'` is a single letter in Czech and has a weight of `0x0EE2`.

```
mysql> SET NAMES 'utf8' COLLATE 'utf8_czech_ci';
Query OK, 0 rows affected (0.09 sec)
mysql> SELECT HEX('ch'), HEX(WEIGHT_STRING('ch'));
+-----------+-------------------------+
| HEX('ch') | HEX(WEIGHT_STRING('ch')) |
+-----------+-------------------------+
| 6368      | 0EE2                    |
+-----------+-------------------------+
1 row in set (0.00 sec)
```

A many-characters-to-many-weights mapping is also possible (this is contraction with expansion), but is not supported by MySQL.

**Miscellaneous collations**

There are also a few collations that do not fall into any of the previous categories.

# 5.2. Choosing a Collation ID

Each collation must have a unique ID. To add a new collation, you must choose an ID value that is not currently used. As of

MySQL 6.0.8, the range of IDs from 1024 to 2047 is reserved for user-defined collations. Before 6.0.8, the ID must be chosen from the range 1 to 254. The collation ID that you choose will show up in these contexts:

- The `Id` column of `SHOW COLLATION` output

- The `ID` column of the `INFORMATION_SCHEMA.COLLATIONS` table

- The `charsetnr` member of the `MYSQL_FIELD` C API data structure

- The `number` member of the `MY_CHARSET_INFO` data structure returned by the `mysql_get_character_set_info()` C API function

To determine the largest currently used ID, issue the following statement:

```
mysql> SELECT MAX(ID) FROM INFORMATION_SCHEMA.COLLATIONS;
+---------+
| MAX(ID) |
+---------+
|     210 |
+---------+
```

To display a list of all currently used IDs, issue this statement:

```
mysql> SELECT ID FROM INFORMATION_SCHEMA.COLLATIONS ORDER BY ID;
+-----+
| ID  |
+-----+
|   1 |
|   2 |
| ... |
|  52 |
|  53 |
|  57 |
|  58 |
| ... |
|  98 |
|  99 |
| 128 |
| 129 |
| ... |
| 210 |
+-----+
```

> **Warning**
>
> Before MySQL 6.0.8, which provides for a range of user-defined collation IDs, you must choose an ID in the range from 1 to 254. In this case, if you upgrade MySQL, you may find that the collation ID you choose has been assigned to a collation included in the new MySQL distribution. In this case, you will need to choose a new value for your own collation.
>
> In addition, before upgrading, you should save the configuration files that you change. If you upgrade in place, the process will replace the your modified files.

# 5.3. Adding a Simple Collation to an 8-Bit Character Set

To add a simple collation for an 8-bit character set without recompiling MySQL, use the following procedure. The example adds a collation named `latin1_test_ci` to the `latin1` character set.

1. Choose a collation ID, as shown in Section 5.2, "Choosing a Collation ID". The following steps use an ID of 1024.

2. You will need to modify the `Index.xml` and `latin1.xml` configuration files. These files will be located in the directory named by the `character_sets_dir` system variable. You can check the variable value as follows, although the path name might be different on your system:

```
mysql> SHOW VARIABLES LIKE 'character_sets_dir';
+--------------------+----------------------------------------+
| Variable_name      | Value                                  |
+--------------------+----------------------------------------+
| character_sets_dir | /user/local/mysql/share/mysql/charsets/ |
+--------------------+----------------------------------------+
```

3. Choose a name for the collation and list it in the `Index.xml` file. Find the `<charset>` element for the character set to which the collation is being added, and add a `<collation>` element that indicates the collation name and ID. For example:

```
<charset name="latin1">
  ...
  <!-- associate collation name with its ID -->
  <collation name="latin1_test_ci" id="1024"/>
  ...
</charset>
```

4. In the `latin1.xml` configuration file, add a `<collation>` element that names the collation and that contains a `<map>` element that defines a character code-to-weight mapping table. Each word within the `<map>` element must be a number in hexadecimal format.

```
<collation name="latin1_test_ci">
<map>
 00 01 02 03 04 05 06 07 08 09 0A 0B 0C 0D 0E 0F
 10 11 12 13 14 15 16 17 18 19 1A 1B 1C 1D 1E 1F
 20 21 22 23 24 25 26 27 28 29 2A 2B 2C 2D 2E 2F
 30 31 32 33 34 35 36 37 38 39 3A 3B 3C 3D 3E 3F
 40 41 42 43 44 45 46 47 48 49 4A 4B 4C 4D 4E 4F
 50 51 52 53 54 55 56 57 58 59 5A 5B 5C 5D 5E 5F
 60 41 42 43 44 45 46 47 48 49 4A 4B 4C 4D 4E 4F
 50 51 52 53 54 55 56 57 58 59 5A 7B 7C 7D 7E 7F
 80 81 82 83 84 85 86 87 88 89 8A 8B 8C 8D 8E 8F
 90 91 92 93 94 95 96 97 98 99 9A 9B 9C 9D 9E 9F
 A0 A1 A2 A3 A4 A5 A6 A7 A8 A9 AA AB AC AD AE AF
 B0 B1 B2 B3 B4 B5 B6 B7 B8 B9 BA BB BC BD BE BF
 41 41 41 41 5B 5D 5B 43 45 45 45 45 49 49 49 49
 44 4E 4F 4F 4F 4F 5C D7 5C 55 55 55 59 59 DE DF
 41 41 41 41 5B 5D 5B 43 45 45 45 45 49 49 49 49
 44 4E 4F 4F 4F 4F 5C F7 5C 55 55 55 59 59 DE FF
</map>
</collation>
```

5. Restart the server and use this statement to verify that the collation is present:

```
mysql> SHOW COLLATION LIKE 'latin1_test_ci';
+----------------+---------+------+---------+----------+---------+
| Collation      | Charset | Id   | Default | Compiled | Sortlen |
+----------------+---------+------+---------+----------+---------+
| latin1_test_ci | latin1  | 1024 |         |          |       1 |
+----------------+---------+------+---------+----------+---------+
```

# 5.4. Adding a UCA Collation to a Unicode Character Set

UCA collations for Unicode character sets can be added to MySQL without recompiling by using a subset of the Locale Data Markup Language (LDML), which is available at http://www.unicode.org/reports/tr35/. In 6.0, this method of adding collations is supported as of MySQL 6.0.4. With this method, you begin with an existing "base" collation. Then you describe the new collation in terms of how it differs from the base collation, rather than defining the entire collation. The following table lists the base collations for the Unicode character sets.

| Character Set | Base Collation |
|---|---|
| utf8 | utf8_unicode_ci |
| ucs2 | ucs2_unicode_ci |
| utf16 | utf16_unicode_ci |
| utf32 | utf32_unicode_ci |

The following brief summary describes the LDML characteristics required for understanding the procedure for adding a collation given later in this section:

- LDML has reset, shift, and identity rules.

- Characters named in these rules can be written in `\u`*nnnn* format, where *nnnn* is the hexadecimal Unicode code point value. Basic Latin letters `A-Z` and `a-z` can also be written literally (this is a MySQL limitation; the LDML specification allows literal non-Latin1 characters in the rules). Only characters in the Basic Multilingual Plane can be specified. This notation does not apply to characters outside the BMP range of `0000` to `FFFF`.

- A reset rule does not specify any ordering in and of itself. Instead, it "resets" the ordering for subsequent shift rules to cause them to be taken in relation to a given character. Either of the following rules resets subsequent shift rules to be taken in relation to the letter `'A'`:

```
<reset>A</reset>
<reset>\u0041</reset>
```

- Shift rules define primary, secondary, and tertiary differences of a character from another character. They are specified using `<p>`, `<s>`, and `<t>` elements. Either of the following rules specifies a primary shift rule for the `'G'` character:

```
<p>G</p>
<p>\u0047</p>
```

  - Use primary differences to distinguish separate letters.

  - Use secondary differences to distinguish accent variations.

  - Use tertiary differences to distinguish lettercase variations.

- Identity rules indicate that one character sorts identically to another. The following rules cause `'b'` sort the same as `'a'`:

```
<reset>a</reset>
<i>b</i>
```

  Identity rules are supported as of MySQL 6.0.9. Prior to 6.0.9, use `<s> ... </s>` instead.

To add a UCA collation for a Unicode character set without recompiling MySQL, use the following procedure. The example adds a collation named `utf8_phone_ci` to the `utf8` character set. The collation is designed for a scenario involving a Web application for which users post their names and phone numbers. Phone numbers can be given in very different formats:

```
+7-12345-67
+7-12-345-67
+7 12 345 67
+7 (12) 345 67
+71234567
```

The problem raised by dealing with these kinds of values is that the varying allowable formats make searching for a specific phone number very difficult. The solution is to define a new collation that reorders punctuation characters, making them ignorable.

1. Choose a collation ID, as shown in Section 5.2, "Choosing a Collation ID". The following steps use an ID of 1029.

2. You will need to modify the `Index.xml` configuration file. This file will be located in the directory named by the `character_sets_dir` system variable. You can check the variable value as follows, although the path name might be different on your system:

```
mysql> SHOW VARIABLES LIKE 'character_sets_dir';
+--------------------+-----------------------------------------+
| Variable_name      | Value                                   |
+--------------------+-----------------------------------------+
| character_sets_dir | /user/local/mysql/share/mysql/charsets/ |
+--------------------+-----------------------------------------+
```

3. Choose a name for the collation and list it in the `Index.xml` file. In addition, you'll need to provide the collation ordering rules. Find the `<charset>` element for the character set to which the collation is being added, and add a `<collation>` element that indicates the collation name and ID. Within the `<collation>` element, provide a `<rules>` element containing the ordering rules:

```
<charset name="utf8">
  ...
  <!-- associate collation name with its ID -->
  <collation name="utf8_phone_ci" id="1029">
    <rules>
      <reset>\u0000</reset>
        <i>\u0020</i> <!-- space -->
        <i>\u0028</i> <!-- left parenthesis -->
        <i>\u0029</i> <!-- right parenthesis -->
        <i>\u002B</i> <!-- plus -->
        <i>\u002D</i> <!-- hyphen -->
    </rules>
  </collation>
  ...
</charset>
```

4. If you want a similar collation for other Unicode character sets, add other `<collation>` elements. For example, to define `ucs2_phone_ci`, add a `<collation>` element to the `<charset name="ucs2">` element. Remember that each collation must have its own unique ID.

5. Restart the server and use this statement to verify that the collation is present:

```
mysql> SHOW COLLATION LIKE 'utf8_phone_ci';
+---------------+---------+------+---------+----------+---------+
| Collation     | Charset | Id   | Default | Compiled | Sortlen |
+---------------+---------+------+---------+----------+---------+
| utf8_phone_ci | utf8    | 1029 |         |          |       8 |
+---------------+---------+------+---------+----------+---------+
```

Now we can test the collation to make sure that it has the desired properties.

Create a table containing some sample phone numbers using the new collation:

```
mysql> CREATE TABLE phonebook (
    ->    name VARCHAR(64),
    ->    phone VARCHAR(64) CHARACTER SET utf8 COLLATE utf8_phone_ci
    -> );
Query OK, 0 rows affected (0.09 sec)
mysql> INSERT INTO phonebook VALUES ('Svoj','+7 912 800 80 02');
Query OK, 1 row affected (0.00 sec)
mysql> INSERT INTO phonebook VALUES ('Hf','+7 (912) 800 80 04');
Query OK, 1 row affected (0.00 sec)
mysql> INSERT INTO phonebook VALUES ('Bar','+7-912-800-80-01');
Query OK, 1 row affected (0.00 sec)
mysql> INSERT INTO phonebook VALUES ('Ramil','(7912) 800 80 03');
Query OK, 1 row affected (0.00 sec)
mysql> INSERT INTO phonebook VALUES ('Sanja','+380 (912) 8008005');
Query OK, 1 row affected (0.00 sec)
```

Run some queries to see whether the ignored punctuation characters are in fact ignored for sorting and comparisons:

```
mysql> SELECT * FROM phonebook ORDER BY phone;
+-------+--------------------+
| name  | phone              |
+-------+--------------------+
| Sanja | +380 (912) 8008005 |
| Bar   | +7-912-800-80-01   |
| Svoj  | +7 912 800 80 02   |
| Ramil | (7912) 800 80 03   |
| Hf    | +7 (912) 800 80 04 |
+-------+--------------------+
5 rows in set (0.00 sec)
mysql> SELECT * FROM phonebook WHERE phone='+7(912)800-80-01';
+------+------------------+
| name | phone            |
+------+------------------+
| Bar  | +7-912-800-80-01 |
+------+------------------+
1 row in set (0.00 sec)
mysql> SELECT * FROM phonebook WHERE phone='79128008001';
+------+------------------+
| name | phone            |
+------+------------------+
| Bar  | +7-912-800-80-01 |
+------+------------------+
1 row in set (0.00 sec)
mysql> SELECT * FROM phonebook WHERE phone='7 9 1 2 8 0 0 8 0 0 1';
+------+------------------+
| name | phone            |
+------+------------------+
| Bar  | +7-912-800-80-01 |
+------+------------------+
1 row in set (0.00 sec)
```

# Chapter 6. Problems With Character Sets

If you try to use a character set that is not compiled into your binary, you might run into the following problems:

- Your program uses an incorrect path to determine where the character sets are stored (which is typically the `share/mysql/charsets` or `share/charsets` directory under the MySQL installation directory). This can be fixed by using the `--character-sets-dir` option when you run the program in question. For example, to specify a directory to be used by MySQL client programs, list it in the `[client]` group of your option file. The examples given here show what the setting might look like for Unix or Windows, respectively:

```
[client]
character-sets-dir=/usr/local/mysql/share/mysql/charsets
[client]
character-sets-dir="C:/Program Files/MySQL/MySQL Server 6.0/share/charsets"
```

- The character set is a complex character set that cannot be loaded dynamically. In this case, you must recompile the program with support for the character set.

  For Unicode character sets, you can define collations without recompiling by using LDML notation. See Section 5.4, "Adding a UCA Collation to a Unicode Character Set".

- The character set is a dynamic character set, but you do not have a configuration file for it. In this case, you should install the configuration file for the character set from a new MySQL distribution.

- If your character set index file does not contain the name for the character set, your program displays an error message. The file is named `Index.xml` and the message is:

```
Character set 'charset_name' is not a compiled character set and is not
specified in the '/usr/share/mysql/charsets/Index.xml' file
```

  To solve this problem, you should either get a new index file or manually add the name of any missing character sets to the current file.

For `MyISAM` tables, you can check the character set name and number for a table with `myisamchk -dvv tbl_name`.

# Chapter 7. MySQL Server Time Zone Support

The MySQL server maintains several time zone settings:

- The system time zone. When the server starts, it attempts to determine the time zone of the host machine and uses it to set the `system_time_zone` system variable. The value does not change thereafter.

  You can set the system time zone for MySQL Server at startup with the `--timezone=timezone_name` option to `mysqld_safe`. You can also set it by setting the `TZ` environment variable before you start `mysqld`. The allowable values for `--timezone` or `TZ` are system-dependent. Consult your operating system documentation to see what values are acceptable.

- The server's current time zone. The global `time_zone` system variable indicates the time zone the server currently is operating in. The initial value for `time_zone` is `'SYSTEM'`, which indicates that the server time zone is the same as the system time zone.

  The initial global server time zone value can be specified explicitly at startup with the `--default-time-zone=timezone` option on the command line, or you can use the following line in an option file:

  ```
  default-time-zone='timezone'
  ```

  If you have the `SUPER` privilege, you can set the global server time zone value at runtime with this statement:

  ```
  mysql> SET GLOBAL time_zone = timezone;
  ```

- Per-connection time zones. Each client that connects has its own time zone setting, given by the session `time_zone` variable. Initially, the session variable takes its value from the global `time_zone` variable, but the client can change its own time zone with this statement:

  ```
  mysql> SET time_zone = timezone;
  ```

The current session time zone setting affects display and storage of time values that are zone-sensitive. This includes the values displayed by functions such as `NOW()` or `CURTIME()`, and values stored in and retrieved from `TIMESTAMP` columns. Values for `TIMESTAMP` columns are converted from the current time zone to UTC for storage, and from UTC to the current time zone for retrieval.

The current time zone setting does not affect values displayed by functions such as `UTC_TIMESTAMP()` or values in `DATE`, `TIME`, or `DATETIME` columns. Nor are values in those data types stored in UTC; the time zone applies for them only when converting from `TIMESTAMP` values. If you want locale-specific arithmetic for `DATE`, `TIME`, or `DATETIME` values, convert them to UTC, perform the arithmetic, and then convert back.

The current values of the global and client-specific time zones can be retrieved like this:

```
mysql> SELECT @@global.time_zone, @@session.time_zone;
```

`timezone` values can be given in several formats, none of which are case sensitive:

- The value `'SYSTEM'` indicates that the time zone should be the same as the system time zone.

- The value can be given as a string indicating an offset from UTC, such as `'+10:00'` or `'-6:00'`.

- The value can be given as a named time zone, such as `'Europe/Helsinki'`, `'US/Eastern'`, or `'MET'`. Named time zones can be used only if the time zone information tables in the `mysql` database have been created and populated.

The MySQL installation procedure creates the time zone tables in the `mysql` database, but does not load them. You must do so manually using the following instructions. (If you are upgrading to MySQL 4.1.3 or later from an earlier version, you can create the tables by upgrading your `mysql` database. Use the instructions in `mysql_upgrade`. After creating the tables, you can load them.)

> **Note**
>
> Loading the time zone information is not necessarily a one-time operation because the information changes occasionally. For example, the rules for Daylight Saving Time in the United States, Mexico, and parts of Canada changed in 2007. When such changes occur, applications that use the old rules become out of date and you may find it necessary to reload the time zone tables to keep the information used by your MySQL server current. See the notes at the end of

this section.

If your system has its own *zoneinfo* database (the set of files describing time zones), you should use the `mysql_tzinfo_to_sql` program for filling the time zone tables. Examples of such systems are Linux, FreeBSD, Sun Solaris, and Mac OS X. One likely location for these files is the `/usr/share/zoneinfo` directory. If your system does not have a zoneinfo database, you can use the downloadable package described later in this section.

The `mysql_tzinfo_to_sql` program is used to load the time zone tables. On the command line, pass the zoneinfo directory path name to `mysql_tzinfo_to_sql` and send the output into the `mysql` program. For example:

```
shell> mysql_tzinfo_to_sql /usr/share/zoneinfo | mysql -u root mysql
```

`mysql_tzinfo_to_sql` reads your system's time zone files and generates SQL statements from them. `mysql` processes those statements to load the time zone tables.

`mysql_tzinfo_to_sql` also can be used to load a single time zone file or to generate leap second information:

- To load a single time zone file `tz_file` that corresponds to a time zone name `tz_name`, invoke `mysql_tzinfo_to_sql` like this:

  ```
  shell> mysql_tzinfo_to_sql tz_file tz_name | mysql -u root mysql
  ```

  With this approach, you must execute a separate command to load the time zone file for each named zone that the server needs to know about.

- If your time zone needs to account for leap seconds, initialize the leap second information like this, where `tz_file` is the name of your time zone file:

  ```
  shell> mysql_tzinfo_to_sql --leap tz_file | mysql -u root mysql
  ```

- After running `mysql_tzinfo_to_sql`, it is best to restart the server so that it does not continue to use any previously cached time zone data.

If your system is one that has no zoneinfo database (for example, Windows or HP-UX), you can use the package of pre-built time zone tables that is available for download at the MySQL Developer Zone:

http://dev.mysql.com/downloads/timezones.html

This time zone package contains `.frm`, `.MYD`, and `.MYI` files for the `MyISAM` time zone tables. These tables should be part of the `mysql` database, so you should place the files in the `mysql` subdirectory of your MySQL server's data directory. The server should be stopped while you do this and restarted afterward.

> **Warning**
>
> Do not use the downloadable package if your system has a zoneinfo database. Use the `mysql_tzinfo_to_sql` utility instead. Otherwise, you may cause a difference in datetime handling between MySQL and other applications on your system.

For information about time zone settings in replication setup, please see Replication Features and Issues.

# 7.1. Staying Current with Time Zone Changes

As mentioned earlier, when the time zone rules change, applications that use the old rules become out of date. To stay current, it is necessary to make sure that your system uses current time zone information is used. For MySQL, there are two factors to consider in staying current:

- The operating system time affects the value that the MySQL server uses for times if its time zone is set to `SYSTEM`. Make sure that your operating system is using the latest time zone information. For most operating systems, the latest update or service pack prepares your system for the time changes. Check the Web site for your operating system vendor for an update that addresses the time changes.

- If you replace the system's `/etc/localtime` timezone file with a version that uses rules differing from those in effect at `mysqld` startup, you should restart `mysqld` so that it uses the updated rules. Otherwise, `mysqld` might not notice when the system changes its time.

- If you use named time zones with MySQL, make sure that the time zone tables in the `mysql` database are up to date. If your system has its own zoneinfo database, you should reload the MySQL time zone tables whenever the zoneinfo database is updated, using the instructions given earlier in this section. For systems that do not have their own zoneinfo database, check the MySQL Developer Zone for updates. When a new update is available, download it and use it to replace your current time zone tables. `mysqld` caches time zone information that it looks up, so after replacing the time zone tables, you should restart `mysqld` to make sure that it does not continue to serve outdated time zone data.

If you are uncertain whether named time zones are available, for use either as the server's time zone setting or by clients that set their own time zone, check whether your time zone tables are empty. The following query determines whether the table that contains time zone names has any rows:

```
mysql> SELECT COUNT(*) FROM mysql.time_zone_name;
+----------+
| COUNT(*) |
+----------+
|        0 |
+----------+
```

A count of zero indicates that the table is empty. In this case, no one can be using named time zones, and you don't need to update the tables. A count greater than zero indicates that the table is not empty and that its contents are available to be used for named time zone support. In this case, you should be sure to reload your time zone tables so that anyone who uses named time zones will get correct query results.

To check whether your MySQL installation is updated properly for a change in Daylight Saving Time rules, use a test like the one following. The example uses values that are appropriate for the 2007 DST 1-hour change that occurs in the United States on March 11 at 2 a.m.

The test uses these two queries:

```
SELECT CONVERT_TZ('2007-03-11 2:00:00','US/Eastern','US/Central');
SELECT CONVERT_TZ('2007-03-11 3:00:00','US/Eastern','US/Central');
```

The two time values indicate the times at which the DST change occurs, and the use of named time zones requires that the time zone tables be used. The desired result is that both queries return the same result (the input time, converted to the equivalent value in the 'US/Central' time zone).

Before updating the time zone tables, you would see an incorrect result like this:

```
mysql> SELECT CONVERT_TZ('2007-03-11 2:00:00','US/Eastern','US/Central');
+-----------------------------------------------------------+
| CONVERT_TZ('2007-03-11 2:00:00','US/Eastern','US/Central') |
+-----------------------------------------------------------+
| 2007-03-11 01:00:00                                       |
+-----------------------------------------------------------+
mysql> SELECT CONVERT_TZ('2007-03-11 3:00:00','US/Eastern','US/Central');
+-----------------------------------------------------------+
| CONVERT_TZ('2007-03-11 3:00:00','US/Eastern','US/Central') |
+-----------------------------------------------------------+
| 2007-03-11 02:00:00                                       |
+-----------------------------------------------------------+
```

After updating the tables, you should see the correct result:

```
mysql> SELECT CONVERT_TZ('2007-03-11 2:00:00','US/Eastern','US/Central');
+-----------------------------------------------------------+
| CONVERT_TZ('2007-03-11 2:00:00','US/Eastern','US/Central') |
+-----------------------------------------------------------+
| 2007-03-11 01:00:00                                       |
+-----------------------------------------------------------+
mysql> SELECT CONVERT_TZ('2007-03-11 3:00:00','US/Eastern','US/Central');
+-----------------------------------------------------------+
| CONVERT_TZ('2007-03-11 3:00:00','US/Eastern','US/Central') |
+-----------------------------------------------------------+
| 2007-03-11 01:00:00                                       |
+-----------------------------------------------------------+
```

# 7.2. Time Zone Leap Second Support

Before MySQL 6.0.9, if the operating system is configured to return leap seconds from OS time calls or if the MySQL server uses a time zone definition that has leap seconds, functions such as `NOW()` could return a value having a time part that ends with `:59:60` or `:59:61`. If such values are inserted into a table, they would be dumped as is by `mysqldump` but considered invalid when reloaded, leading to backup/restore problems.

As of MySQL 6.0.9, leap second values are returned with a time part that ends with `:59:59`. This means that a function such as `NOW()` can return the same value for two or three consecutive seconds during the leap second. It remains true that literal temporal

values having a time part that ends with `:59:60` or `:59:61` are considered invalid.

If it is necessary to search for `TIMESTAMP` values one second before the leap second, anomalous results may be obtained if you use a comparison with `'YYYY-MM-DD hh:mm:ss'` values:

```
mysql> CREATE TABLE t1 (a INT, ts TIMESTAMP DEFAULT NOW(), PRIMARY KEY (ts));
Query OK, 0 rows affected (0.11 sec)
mysql> # Simulate NOW() = '2009-01-01 02:59:59'
mysql> SET timestamp = 1230768022;
Query OK, 0 rows affected (0.00 sec)
mysql> INSERT INTO t1 (a) VALUES (1);
Query OK, 1 row affected (0.07 sec)
mysql> # Simulate NOW() = '2009-01-01 02:59:60'
mysql> SET timestamp = 1230768023;
Query OK, 0 rows affected (0.00 sec)
mysql> INSERT INTO t1 (a) VALUES (2);
Query OK, 1 row affected (0.02 sec)
mysql> SELECT * FROM t1;
+------+---------------------+
| a    | ts                  |
+------+---------------------+
|    1 | 2008-12-31 18:00:22 |
|    2 | 2008-12-31 18:00:23 |
+------+---------------------+
2 rows in set (0.02 sec)
mysql> SELECT * FROM t1 WHERE ts = '2009-01-01 02:59:59';
Empty set (0.03 sec)
```

To work around this, you can use a comparison based on the UTC value actually stored in column, which has the leap second correction applied:

```
mysql> SELECT * FROM t1 WHERE UNIX_TIMESTAMP(ts) = 1230768023;
+------+---------------------+
| a    | ts                  |
+------+---------------------+
|    2 | 2008-12-31 18:00:23 |
+------+---------------------+
1 row in set (0.02 sec)
```

# Chapter 8. MySQL Server Locale Support

The locale indicated by the `lc_time_names` system variable controls the language used to display day and month names and abbreviations. This variable affects the output from the `DATE_FORMAT()`, `DAYNAME()` and `MONTHNAME()` functions.

Locale names are POSIX-style values such as `'ja_JP'` or `'pt_BR'`. The default value is `'en_US'` regardless of your system's locale setting, but you can set the value at server startup or set the `GLOBAL` value if you have the `SUPER` privilege. Any client can examine the value of `lc_time_names` or set its `SESSION` value to affect the locale for its own connection.

```
mysql> SET NAMES 'utf8';
Query OK, 0 rows affected (0.09 sec)
mysql> SELECT @@lc_time_names;
+-----------------+
| @@lc_time_names |
+-----------------+
| en_US           |
+-----------------+
1 row in set (0.00 sec)
mysql> SELECT DAYNAME('2010-01-01'), MONTHNAME('2010-01-01');
+-----------------------+-------------------------+
| DAYNAME('2010-01-01') | MONTHNAME('2010-01-01') |
+-----------------------+-------------------------+
| Friday                | January                 |
+-----------------------+-------------------------+
1 row in set (0.00 sec)
mysql> SELECT DATE_FORMAT('2010-01-01','%W %a %M %b');
+----------------------------------------+
| DATE_FORMAT('2010-01-01','%W %a %M %b') |
+----------------------------------------+
| Friday Fri January Jan                 |
+----------------------------------------+
1 row in set (0.00 sec)
mysql> SET lc_time_names = 'es_MX';
Query OK, 0 rows affected (0.00 sec)
mysql> SELECT @@lc_time_names;
+-----------------+
| @@lc_time_names |
+-----------------+
| es_MX           |
+-----------------+
1 row in set (0.00 sec)
mysql> SELECT DAYNAME('2010-01-01'), MONTHNAME('2010-01-01');
+-----------------------+-------------------------+
| DAYNAME('2010-01-01') | MONTHNAME('2010-01-01') |
+-----------------------+-------------------------+
| viernes               | enero                   |
+-----------------------+-------------------------+
1 row in set (0.00 sec)
mysql> SELECT DATE_FORMAT('2010-01-01','%W %a %M %b');
+----------------------------------------+
| DATE_FORMAT('2010-01-01','%W %a %M %b') |
+----------------------------------------+
| viernes vie enero ene                  |
+----------------------------------------+
1 row in set (0.00 sec)
```

The day or month name for each of the affected functions is converted from `utf8` to the character set indicated by the `character_set_connection` system variable.

`lc_time_names` may be set to any of the following locale values.

| | |
|---|---|
| ar_AE: Arabic - United Arab Emirates | ar_BH: Arabic - Bahrain |
| ar_DZ: Arabic - Algeria | ar_EG: Arabic - Egypt |
| ar_IN: Arabic - Iran | ar_IQ: Arabic - Iraq |
| ar_JO: Arabic - Jordan | ar_KW: Arabic - Kuwait |
| ar_LB: Arabic - Lebanon | ar_LY: Arabic - Libya |
| ar_MA: Arabic - Morocco | ar_OM: Arabic - Oman |
| ar_QA: Arabic - Qatar | ar_SA: Arabic - Saudi Arabia |
| ar_SD: Arabic - Sudan | ar_SY: Arabic - Syria |
| ar_TN: Arabic - Tunisia | ar_YE: Arabic - Yemen |
| be_BY: Belarusian - Belarus | bg_BG: Bulgarian - Bulgaria |
| ca_ES: Catalan - Catalan | cs_CZ: Czech - Czech Republic |
| da_DK: Danish - Denmark | de_AT: German - Austria |
| de_BE: German - Belgium | de_CH: German - Switzerland |
| de_DE: German - Germany | de_LU: German - Luxembourg |

| | |
|---|---|
| `EE`: Estonian - Estonia | `en_AU`: English - Australia |
| `en_CA`: English - Canada | `en_GB`: English - United Kingdom |
| `en_IN`: English - India | `en_NZ`: English - New Zealand |
| `en_PH`: English - Philippines | `en_US`: English - United States |
| `en_ZA`: English - South Africa | `en_ZW`: English - Zimbabwe |
| `es_AR`: Spanish - Argentina | `es_BO`: Spanish - Bolivia |
| `es_CL`: Spanish - Chile | `es_CO`: Spanish - Columbia |
| `es_CR`: Spanish - Costa Rica | `es_DO`: Spanish - Dominican Republic |
| `es_EC`: Spanish - Ecuador | `es_ES`: Spanish - Spain |
| `es_GT`: Spanish - Guatemala | `es_HN`: Spanish - Honduras |
| `es_MX`: Spanish - Mexico | `es_NI`: Spanish - Nicaragua |
| `es_PA`: Spanish - Panama | `es_PE`: Spanish - Peru |
| `es_PR`: Spanish - Puerto Rico | `es_PY`: Spanish - Paraguay |
| `es_SV`: Spanish - El Salvador | `es_US`: Spanish - United States |
| `es_UY`: Spanish - Uruguay | `es_VE`: Spanish - Venezuela |
| `eu_ES`: Basque - Basque | `fi_FI`: Finnish - Finland |
| `fo_FO`: Faroese - Faroe Islands | `fr_BE`: French - Belgium |
| `fr_CA`: French - Canada | `fr_CH`: French - Switzerland |
| `fr_FR`: French - France | `fr_LU`: French - Luxembourg |
| `gl_ES`: Galician - Galician | `gu_IN`: Gujarati - India |
| `he_IL`: Hebrew - Israel | `hi_IN`: Hindi - India |
| `hr_HR`: Croatian - Croatia | `hu_HU`: Hungarian - Hungary |
| `id_ID`: Indonesian - Indonesia | `is_IS`: Icelandic - Iceland |
| `it_CH`: Italian - Switzerland | `it_IT`: Italian - Italy |
| `ja_JP`: Japanese - Japan | `ko_KR`: Korean - Korea |
| `lt_LT`: Lithuanian - Lithuania | `lv_LV`: Latvian - Latvia |
| `mk_MK`: Macedonian - FYROM | `mn_MN`: Mongolia - Mongolian |
| `ms_MY`: Malay - Malaysia | `nb_NO`: Norwegian(Bokml) - Norway |
| `nl_BE`: Dutch - Belgium | `nl_NL`: Dutch - The Netherlands |
| `no_NO`: Norwegian - Norway | `pl_PL`: Polish - Poland |
| `pt_BR`: Portugese - Brazil | `pt_PT`: Portugese - Portugal |
| `ro_RO`: Romanian - Romania | `ru_RU`: Russian - Russia |
| `ru_UA`: Russian - Ukraine | `sk_SK`: Slovak - Slovakia |
| `sl_SI`: Slovenian - Slovenia | `sq_AL`: Albanian - Albania |
| `sr_YU`: Serbian - Yugoslavia | `sv_FI`: Swedish - Finland |
| `sv_SE`: Swedish - Sweden | `ta_IN`: Tamil - India |
| `te_IN`: Telugu - India | `th_TH`: Thai - Thailand |
| `tr_TR`: Turkish - Turkey | `uk_UA`: Ukrainian - Ukraine |
| `ur_PK`: Urdu - Pakistan | `vi_VN`: Vietnamese - Vietnam |
| `zh_CN`: Chinese - Peoples Republic of China | `zh_HK`: Chinese - Hong Kong SAR |
| `zh_TW`: Chinese - Taiwan | |

`lc_time_names` currently does not affect the `STR_TO_DATE()` or `GET_FORMAT()` function.

# Chapter 9. MySQL 6.0 FAQ — MySQL Chinese, Japanese, and Korean Character Sets

This set of Frequently Asked Questions derives from the experience of MySQL's Support and Development groups in handling many inquiries about CJK (Chinese-Japanese-Korean) issues.

**Questions**

- 9.1: What CJK character sets are available in MySQL?

- 9.2: I have inserted CJK characters into my table. Why does `SELECT` display them as "?" characters?

- 9.3: What problems should I be aware of when working with the Big5 Chinese character set?

- 9.4: Why do Japanese character set conversions fail?

- 9.5: What should I do if I want to convert SJIS `81CA` to `cp932`?

- 9.6: How does MySQL represent the Yen (Â¥) sign?

- 9.7: Do MySQL plan to make a separate character set where `5C` is the Yen sign, as at least one other major DBMS does?

- 9.8: Of what issues should I be aware when working with Korean character sets in MySQL?

- 9.9: Why do I get `DATA TRUNCATED` error messages?

- 9.10: Why does my GUI front end or browser not display CJK characters correctly in my application using Access, PHP, or another API?

- 9.11: I've upgraded to MySQL 6.0. How can I revert to behavior like that in MySQL 4.0 with regard to character sets?

- 9.12: Why do some `LIKE` and `FULLTEXT` searches with CJK characters fail?

- 9.13: How do I know whether character *X* is available in all character sets?

- 9.14: Why don't CJK strings sort correctly in Unicode? (I)

- 9.15: Why don't CJK strings sort correctly in Unicode? (II)

- 9.16: Why are my supplementary characters rejected by MySQL?

- 9.17: Shouldn't it be "CJKV"?

- 9.18: Does MySQL allow CJK characters to be used in database and table names?

- 9.19: Where can I find translations of the MySQL Manual into Chinese, Japanese, and Korean?

- 9.20: Where can I get help with CJK and related issues in MySQL?

**Questions and Answers**

**9.1: What CJK character sets are available in MySQL?**

The list of CJK character sets may vary depending on your MySQL version. For example, the `eucjpms` character set was not supported prior to MySQL 5.0.3. However, since the name of the applicable language appears in the `DESCRIPTION` column for every entry in the `INFORMATION_SCHEMA.CHARACTER_SETS` table, you can obtain a current list of all the non-Unicode CJK character sets using this query:

```
mysql> SELECT CHARACTER_SET_NAME, DESCRIPTION
    -> FROM INFORMATION_SCHEMA.CHARACTER_SETS
    -> WHERE DESCRIPTION LIKE '%Chinese%'
    -> OR DESCRIPTION LIKE '%Japanese%'
    -> OR DESCRIPTION LIKE '%Korean%'
    -> ORDER BY CHARACTER_SET_NAME;
+--------------------+-------------------------+
| CHARACTER_SET_NAME | DESCRIPTION             |
+--------------------+-------------------------+
| big5               | Big5 Traditional Chinese |
| cp932              | SJIS for Windows Japanese |
| eucjpms            | UJIS for Windows Japanese |
| euckr              | EUC-KR Korean           |
| gb2312             | GB2312 Simplified Chinese |
```

```
| gbk               | GBK Simplified Chinese   |
| sjis              | Shift-JIS Japanese       |
| ujis              | EUC-JP Japanese          |
+-------------------+--------------------------+
8 rows in set (0.01 sec)
```

(See The INFORMATION_SCHEMA CHARACTER_SETS Table, for more information.)

MySQL supports the two common variants of the *GB* (*Guojia Biaozhun*, or *National Standard*, or *Simplified Chinese*) character sets which are official in the People's Republic of China: gb2312 and gbk. Sometimes people try to insert gbk characters into gb2312, and it works most of the time because gbk is a superset of gb2312 — but eventually they try to insert a rarer Chinese character and it doesn't work. (See Bug#16072 for an example).

Here, we try to clarify exactly what characters are legitimate in gb2312 or gbk, with reference to the official documents. Please check these references before reporting gb2312 or gbk bugs.

- For a complete listing of the gb2312 characters, ordered according to the gb2312_chinese_ci collation: gb2312

- MySQL's gbk is in reality "Microsoft code page 936". This differs from the official gbk for characters A1A4 (middle dot), A1AA (em dash), A6E0-A6F5, and A8BB-A8C0. For a listing of the differences, see http://recode.progiciels-bpi.ca/showfile.html?name=dist/libiconv/gbk.h.

- For a listing of gbk/Unicode mappings, see http://www.unicode.org/Public/MAPPINGS/VENDORS/MICSFT/WINDOWS/CP936.TXT.

- For MySQL's listing of gbk characters, see gbk.

**9.2: I have inserted CJK characters into my table. Why does SELECT display them as "?" characters?**

This problem is usually due to a setting in MySQL that doesn't match the settings for the application program or the operating system. Here are some common steps for correcting these types of issues:

- *Be certain of what MySQL version you are using*.

  Use the statement SELECT VERSION(); to determine this.

- *Make sure that the database is actually using the desired character set*.

  People often think that the client character set is always the same as either the server character set or the character set used for display purposes. However, both of these are false assumptions. You can make sure by checking the result of SHOW CREATE TABLE *tablename* or — better — yet by using this statement:

```
SELECT character_set_name, collation_name
    FROM information_schema.columns
    WHERE table_schema = your_database_name
        AND table_name = your_table_name
        AND column_name = your_column_name;
```

- *Determine the hexadecimal value of the character or characters that are not being displayed correctly*.

  You can obtain this information for a column *column_name* in the table *table_name* using the following query:

```
SELECT HEX(column_name)
FROM table_name;
```

  3F is the encoding for the ? character; this means that ? is the character actually stored in the column. This most often happens because of a problem converting a particular character from your client character set to the target character set.

- *Make sure that a round trip possible — that is, when you select* literal *(or* _introducer hexadecimal-value*), you obtain* literal *as a result*.

  For example, the Japanese *Katakana* character *Pe* (ã##') exists in all CJK character sets, and has the code point value (hexadecimal coding) 0x30da. To test a round trip for this character, use this query:

```
SELECT 'ã##' AS `ã##`;          /* or SELECT _ucs2 0x30da; */
```

  If the result is not also ã##, then the round trip has failed.

  For bug reports regarding such failures, we might ask you to follow up with SELECT HEX('ã##');. Then we can determine whether the client encoding is correct.

- *Make sure that the problem is not with the browser or other application, rather than with MySQL*.

Use the `mysql` client program (on Windows: `mysql.exe`) to accomplish this task. If `mysql` displays correctly but your application doesn't, then your problem is probably due to system settings.

To find out what your settings are, use the `SHOW VARIABLES` statement, whose output should resemble what is shown here:

```
mysql> SHOW VARIABLES LIKE 'char%';
+--------------------------+--------------------------------------+
| Variable_name            | Value                                |
+--------------------------+--------------------------------------+
| character_set_client     | utf8                                 |
| character_set_connection | utf8                                 |
| character_set_database   | latin1                               |
| character_set_filesystem | binary                               |
| character_set_results    | utf8                                 |
| character_set_server     | latin1                               |
| character_set_system     | utf8                                 |
| character_sets_dir       | /usr/local/mysql/share/mysql/charsets/ |
+--------------------------+--------------------------------------+
8 rows in set (0.03 sec)
```

These are typical character-set settings for an international-oriented client (notice the use of `utf8` Unicode) connected to a server in the West (`latin1` is a West Europe character set and a default for MySQL).

Although Unicode (usually the `utf8` variant on Unix, and the `ucs2` variant on Windows) is preferable to Latin, it is often not what your operating system utilities support best. Many Windows users find that a Microsoft character set, such as `cp932` for Japanese Windows, is suitable.

If you cannot control the server settings, and you have no idea what your underlying computer is, then try changing to a common character set for the country that you're in (`euckr` = Korea; `gb2312` or `gbk` = People's Republic of China; `big5` = Taiwan; `sjis`, `ujis`, `cp932`, or `eucjpms` = Japan; `ucs2` or `utf8` = anywhere). Usually it is necessary to change only the client and connection and results settings. There is a simple statement which changes all three at once: `SET NAMES`. For example:

```
SET NAMES 'big5';
```

Once the setting is correct, you can make it permanent by editing `my.cnf` or `my.ini`. For example you might add lines looking like these:

```
[mysqld]
character-set-server=big5
[client]
default-character-set=big5
```

It is also possible that there are issues with the API configuration setting being used in your application; see *Why does my GUI front end or browser not display CJK characters correctly...?* for more information.

**9.3: What problems should I be aware of when working with the Big5 Chinese character set?**

MySQL supports the Big5 character set which is common in Hong Kong and Taiwan (Republic of China). MySQL's `big5` is in reality Microsoft code page 950, which is very similar to the original `big5` character set. We changed to this character set starting with MySQL version 4.1.16 / 5.0.16 (as a result of Bug#12476). For example, the following statements work in current versions of MySQL, but not in old versions:

```
mysql> CREATE TABLE big5 (BIG5 CHAR(1) CHARACTER SET BIG5);
Query OK, 0 rows affected (0.13 sec)
mysql> INSERT INTO big5 VALUES (0xf9dc);
Query OK, 1 row affected (0.00 sec)
mysql> SELECT * FROM big5;
+------+
| big5 |
+------+
| 許   |
+------+
1 row in set (0.02 sec)
```

A feature request for adding HKSCS extensions has been filed. People who need this extension may find the suggested patch for Bug#13577 to be of interest.

**9.4: Why do Japanese character set conversions fail?**

MySQL supports the `sjis`, `ujis`, `cp932`, and eucjpms character sets, as well as Unicode. A common need is to convert between character sets. For example, there might be a Unix server (typically with `sjis` or `ujis`) and a Windows client (typically with `cp932`).

In the following conversion table, the `ucs2` column represents the source, and the `sjis`, `cp932`, `ujis`, and `eucjpms` columns represent the destinations — that is, the last 4 columns provide the hexadecimal result when we use `CONVERT(ucs2)` or we as-

sign a `ucs2` column containing the value to an `sjis`, `cp932`, `ujis`, or `eucjpms` column.

| Character Name | ucs2 | sjis | cp932 | ujis | eucjpms |
|---|---|---|---|---|---|
| BROKEN BAR | 00A6 | 3F | 3F | 8FA2C3 | 3F |
| FULLWIDTH BROKEN BAR | FFE4 | 3F | FA55 | 3F | 8FA2 |
| YEN SIGN | 00A5 | 3F | 3F | 20 | 3F |
| FULLWIDTH YEN SIGN | FFE5 | 818F | 818F | A1EF | 3F |
| TILDE | 007E | 7E | 7E | 7E | 7E |
| OVERLINE | 203E | 3F | 3F | 20 | 3F |
| HORIZONTAL BAR | 2015 | 815C | 815C | A1BD | A1BD |
| EM DASH | 2014 | 3F | 3F | 3F | 3F |
| REVERSE SOLIDUS | 005C | 815F | 5C | 5C | 5C |
| FULLWIDTH "" | FF3C | 3F | 815F | 3F | A1C0 |
| WAVE DASH | 301C | 8160 | 3F | A1C1 | 3F |
| FULLWIDTH TILDE | FF5E | 3F | 8160 | 3F | A1C1 |
| DOUBLE VERTICAL LINE | 2016 | 8161 | 3F | A1C2 | 3F |
| PARALLEL TO | 2225 | 3F | 8161 | 3F | A1C2 |
| MINUS SIGN | 2212 | 817C | 3F | A1DD | 3F |
| FULLWIDTH HYPHEN-MINUS | FF0D | 3F | 817C | 3F | A1DD |
| CENT SIGN | 00A2 | 8191 | 3F | A1F1 | 3F |
| FULLWIDTH CENT SIGN | FFE0 | 3F | 8191 | 3F | A1F1 |
| POUND SIGN | 00A3 | 8192 | 3F | A1F2 | 3F |
| FULLWIDTH POUND SIGN | FFE1 | 3F | 8192 | 3F | A1F2 |
| NOT SIGN | 00AC | 81CA | 3F | A2CC | 3F |
| FULLWIDTH NOT SIGN | FFE2 | 3F | 81CA | 3F | A2CC |

Now consider the following portion of the table.

| | ucs2 | sjis | cp932 |
|---|---|---|---|
| NOT SIGN | 00AC | 81CA | 3F |
| FULLWIDTH NOT SIGN | FFE2 | 3F | 81CA |

This means that MySQL converts the `NOT SIGN` (Unicode `U+00AC`) to `sjis` code point `0x81CA` and to `cp932` code point `3F`. (`3F` is the question mark ("?") — this is what is always used when the conversion cannot be performed.

**9.5: What should I do if I want to convert SJIS `81CA` to `cp932`?**

Our answer is: "?". There are serious complaints about this: many people would prefer a "loose" conversion, so that `81CA (NOT SIGN)` in `sjis` becomes `81CA (FULLWIDTH NOT SIGN)` in `cp932`. We are considering a change to this behavior.

**9.6: How does MySQL represent the Yen (â¥) sign?**

A problem arises because some versions of Japanese character sets (both `sjis` and `euc`) treat `5C` as a *reverse solidus* (\ — also known as a backslash), and others treat it as a yen sign (â¥).

MySQL follows only one version of the JIS (Japanese Industrial Standards) standard description. In MySQL, *5C is always the reverse solidus (\)*.

**9.7: Do MySQL plan to make a separate character set where `5C` is the Yen sign, as at least one other major DBMS does?**

This is one possible solution to the Yen sign issue; however, this will not happen in MySQL 5.1 or 6.0.

**9.8: Of what issues should I be aware when working with Korean character sets in MySQL?**

In theory, while there have been several versions of the `euckr` (*Extended Unix Code Korea*) character set, only one problem has been noted.

We use the "ASCII" variant of EUC-KR, in which the code point `0x5c` is REVERSE SOLIDUS, that is \, instead of the "KS-Roman" variant of EUC-KR, in which the code point `0x5c` is `WON SIGN`(â#©). This means that you cannot convert Unicode `U+20A9` to `euckr`:

```
mysql> SELECT
    ->     CONVERT('â#©' USING euckr) AS euckr,
    ->     HEX(CONVERT('â#©' USING euckr)) AS hexeuckr;
+-------+----------+
| euckr | hexeuckr |
+-------+----------+
| ?     | 3F       |
+-------+----------+
1 row in set (0.00 sec)
```

MySQL's graphic Korean chart is here: euckr.

### 9.9: Why do I get DATA TRUNCATED error messages?

For illustration, we'll create a table with one Unicode (ucs2) column and one Chinese (gb2312) column.

```
mysql> CREATE TABLE ch
    -> (ucs2 CHAR(3) CHARACTER SET ucs2,
    -> gb2312 CHAR(3) CHARACTER SET gb2312);
Query OK, 0 rows affected (0.05 sec)
```

We'll try to place the rare character æ±# in both columns.

```
mysql> INSERT INTO ch VALUES ('Aæ±#B','Aæ±#B');
Query OK, 1 row affected, 1 warning (0.00 sec)
```

Ah, there is a warning. Use the following statement to see what it is:

```
mysql> SHOW WARNINGS;
+---------+------+-------------------------------------------+
| Level   | Code | Message                                   |
+---------+------+-------------------------------------------+
| Warning | 1265 | Data truncated for column 'gb2312' at row 1 |
+---------+------+-------------------------------------------+
1 row in set (0.00 sec)
```

So it is a warning about the gb2312 column only.

```
mysql> SELECT ucs2,HEX(ucs2),gb2312,HEX(gb2312) FROM ch;
+-------+--------------+--------+-------------+
| ucs2  | HEX(ucs2)    | gb2312 | HEX(gb2312) |
+-------+--------------+--------+-------------+
| Aæ±#B | 00416C4C0042 | A?B    | 413F42      |
+-------+--------------+--------+-------------+
1 row in set (0.00 sec)
```

There are several things that need explanation here.

1.  The fact that it is a "warning" rather than an "error" is characteristic of MySQL. We like to try to do what we can, to get the best fit, rather than give up.

2.  The æ±# character isn't in the gb2312 character set. We described that problem earlier.

3.  Admittedly the message is misleading. We didn't "truncate" in this case, we replaced with a question mark. We've had a complaint about this message (See Bug#9337). But until we come up with something better, just accept that error/warning code 2165 can mean a variety of things.

4.  With SQL_MODE=TRADITIONAL, there would be an error message, but instead of error 2165 you would see: ERROR 1406 (22001): Data too long for column 'gb2312' at row 1.

### 9.10: Why does my GUI front end or browser not display CJK characters correctly in my application using Access, PHP, or another API?

Obtain a direct connection to the server using the mysql client (Windows: mysql.exe), and try the same query there. If mysql responds correctly, then the trouble may be that your application interface requires initialization. Use mysql to tell you what character set or sets it uses with the statement SHOW VARIABLES LIKE 'char%';. If you are using Access, then you are most likely connecting with MyODBC. In this case, you should check Connector/ODBC Configuration. If, for instance, you use big5, you would enter SET NAMES 'big5'. (Note that no ; is required in this case). If you are using ASP, you might need to add SET NAMES in the code. Here is an example that has worked in the past:

```
<%
Session.CodePage=0
Dim strConnection
Dim Conn
strConnection="driver={MySQL ODBC 3.51 Driver};server=server;uid=username;" \
          & "pwd=password;database=database;stmt=SET NAMES 'big5';"
Set Conn = Server.CreateObject("ADODB.Connection")
Conn.Open strConnection
%>
```

In much the same way, if you are using any character set other than latin1 with Connector/NET, then you must specify the char-

acter set in the connection string. See Connecting to MySQL Using Connector/NET, for more information.

If you are using PHP, try this:

```
<?php
  $link = mysql_connect($host, $usr, $pwd);
  mysql_select_db($db);
  if( mysql_error() ) { print "Database ERROR: " . mysql_error(); }
  mysql_query("SET NAMES 'utf8'", $link);
?>
```

In this case, we used SET NAMES to change character_set_client and character_set_connection and character_set_results.

We encourage the use of the newer mysqli extension, rather than mysql. Using mysqli, the previous example could be rewritten as shown here:

```
<?php
  $link = new mysqli($host, $usr, $pwd, $db);
  if( mysqli_connect_errno() )
  {
    printf("Connect failed: %s\n", mysqli_connect_error());
    exit();
  }
  $link->query("SET NAMES 'utf8'");
?>
```

Another issue often encountered in PHP applications has to do with assumptions made by the browser. Sometimes adding or changing a <meta> tag suffices to correct the problem: for example, to insure that the user agent interprets page content as UTF-8, you should include <meta http-equiv="Content-Type" content="text/html; charset=utf-8"> in the <head> of the HTML page.

If you are using Connector/J, see Using Character Sets and Unicode.

**9.11: I've upgraded to MySQL 6.0. How can I revert to behavior like that in MySQL 4.0 with regard to character sets?**

In MySQL Version 4.0, there was a single "global" character set for both server and client, and the decision as to which character to use was made by the server administrator. This changed starting with MySQL Version 4.1. What happens now is a "handshake", as described in Section 1.4, "Connection Character Sets and Collations":

> When a client connects, it sends to the server the name of the character set that it wants to use. The server uses the name to set the character_set_client, character_set_results, and character_set_connection system variables. In effect, the server performs a SET NAMES operation using the character set name.

The effect of this is that you cannot control the client character set by starting mysqld with --character-set-server=utf8. However, some of our Asian customers have said that they prefer the MySQL 4.0 behavior. To make it possible to retain this behavior, we added a mysqld switch, --character-set-client-handshake, which can be turned off with --skip-character-set-client-handshake. If you start mysqld with --skip-character-set-client-handshake, then, when a client connects, it sends to the server the name of the character set that it wants to use — however, *the server ignores this request from the client*.

By way of example, suppose that your favorite server character set is latin1 (unlikely in a CJK area, but this is the default value). Suppose further that the client uses utf8 because this is what the client's operating system supports. Now, start the server with latin1 as its default character set:

```
mysqld --character-set-server=latin1
```

And then start the client with the default character set utf8:

```
mysql --default-character-set=utf8
```

The current settings can be seen by viewing the output of SHOW VARIABLES:

```
mysql> SHOW VARIABLES LIKE 'char%';
+--------------------------+------------------------------------+
| Variable_name            | Value                              |
+--------------------------+------------------------------------+
| character_set_client     | utf8                               |
| character_set_connection | utf8                               |
| character_set_database   | latin1                             |
| character_set_filesystem | binary                             |
| character_set_results    | utf8                               |
| character_set_server     | latin1                             |
| character_set_system     | utf8                               |
| character_sets_dir       | /usr/local/mysql/share/mysql/charsets/ |
+--------------------------+------------------------------------+
8 rows in set (0.01 sec)
```

Now stop the client, and then stop the server using mysqladmin. Then start the server again, but this time tell it to skip the handshake like so:

```
mysqld --character-set-server=utf8 --skip-character-set-client-handshake
```

Start the client with `utf8` once again as the default character set, then display the current settings:

```
mysql> SHOW VARIABLES LIKE 'char%';
+--------------------------+------------------------------------+
| Variable_name            | Value                              |
+--------------------------+------------------------------------+
| character_set_client     | latin1                             |
| character_set_connection | latin1                             |
| character_set_database   | latin1                             |
| character_set_filesystem | binary                             |
| character_set_results    | latin1                             |
| character_set_server     | latin1                             |
| character_set_system     | utf8                               |
| character_sets_dir       | /usr/local/mysql/share/mysql/charsets/ |
+--------------------------+------------------------------------+
8 rows in set (0.01 sec)
```

As you can see by comparing the differing results from SHOW VARIABLES, the server ignores the client's initial settings if the `--skip-character-set-client-handshake` is used.

**9.12: Why do some `LIKE` and `FULLTEXT` searches with CJK characters fail?**

There is a very simple problem with LIKE searches on BINARY and BLOB columns: we need to know the end of a character. With multi-byte character sets, different characters might have different octet lengths. For example, in `utf8`, A requires one byte but ã## requires three bytes, as shown here:

```
+------------------------+--------------------------+
| OCTET_LENGTH(_utf8 'A') | OCTET_LENGTH(_utf8 'ã##') |
+------------------------+--------------------------+
|                      1 |                        3 |
+------------------------+--------------------------+
1 row in set (0.00 sec)
```

If we don't know where the first character ends, then we don't know where the second character begins, in which case even very simple searches such as LIKE `'_A%'` fail. The solution is to use a regular CJK character set in the first place, or to convert to a CJK character set before comparing.

This is one reason why MySQL cannot allow encodings of non-existent characters. If it is not strict about rejecting bad input, then it has no way of knowing where characters end.

For FULLTEXT searches, we need to know where words begin and end. With Western languages, this is rarely a problem because most (if not all) of these use an easy-to-identify word boundary — the space character. However, this is not usually the case with Asian writing. We could use arbitrary halfway measures, like assuming that all Han characters represent words, or (for Japanese) depending on changes from Katakana to Hiragana due to grammatical endings. However, the only sure solution requires a comprehensive word list, which means that we would have to include a dictionary in the server for each Asian language supported. This is simply not feasible.

**9.13: How do I know whether character _X_ is available in all character sets?**

The majority of simplified Chinese and basic non-halfwidth Japanese *Kana* characters appear in all CJK character sets. This stored procedure accepts a `UCS-2` Unicode character, converts it to all other character sets, and displays the results in hexadecimal.

```
DELIMITER //
CREATE PROCEDURE p_convert(ucs2_char CHAR(1) CHARACTER SET ucs2)
BEGIN
CREATE TABLE tj
            (ucs2 CHAR(1) character set ucs2,
             utf8 CHAR(1) character set utf8,
             big5 CHAR(1) character set big5,
             cp932 CHAR(1) character set cp932,
             eucjpms CHAR(1) character set eucjpms,
             euckr CHAR(1) character set euckr,
             gb2312 CHAR(1) character set gb2312,
             gbk CHAR(1) character set gbk,
             sjis CHAR(1) character set sjis,
             ujis CHAR(1) character set ujis);
INSERT INTO tj (ucs2) VALUES (ucs2_char);
UPDATE tj SET utf8=ucs2,
             big5=ucs2,
             cp932=ucs2,
             eucjpms=ucs2,
             euckr=ucs2,
             gb2312=ucs2,
             gbk=ucs2,
             sjis=ucs2,
             ujis=ucs2;
/* If there is a conversion problem, UPDATE will produce a warning. */
SELECT hex(ucs2) AS ucs2,
       hex(utf8) AS utf8,
       hex(big5) AS big5,
       hex(cp932) AS cp932,
       hex(eucjpms) AS eucjpms,
       hex(euckr) AS euckr,
       hex(gb2312) AS gb2312,
```

```
        hex(gbk) AS gbk,
        hex(sjis) AS sjis,
        hex(ujis) AS ujis
FROM tj;
DROP TABLE tj;
END//
```

The input can be any single `ucs2` character, or it can be the code point value (hexadecimal representation) of that character. For example, from Unicode's list of `ucs2` encodings and names (http://www.unicode.org/Public/UNIDATA/UnicodeData.txt), we know that the *Katakana* character *Pe* appears in all CJK character sets, and that its code point value is `0x30da`. If we use this value as the argument to `p_convert()`, the result is as shown here:

```
mysql> CALL p_convert(0x30da)//
+------+--------+------+-------+---------+-------+--------+------+------+------+
| ucs2 | utf8   | big5 | cp932 | eucjpms | euckr | gb2312 | gbk  | sjis | ujis |
+------+--------+------+-------+---------+-------+--------+------+------+------+
| 30DA | E3839A | C772 | 8379  | A5DA    | ABDA  | A5DA   | A5DA | 8379 | A5DA |
+------+--------+------+-------+---------+-------+--------+------+------+------+
1 row in set (0.04 sec)
```

Since none of the column values is `3F` — that is, the question mark character (`?`) — we know that every conversion worked.

### 9.14: Why don't CJK strings sort correctly in Unicode? (I)

Sometimes people observe that the result of a `utf8_unicode_ci` or `ucs2_unicode_ci` search, or of an `ORDER BY` sort is not what they think a native would expect. Although we never rule out the possibility that there is a bug, we have found in the past that many people do not read correctly the standard table of weights for the Unicode Collation Algorithm. MySQL uses the table found at http://www.unicode.org/Public/UCA/4.0.0/allkeys-4.0.0.txt. This is not the first table you will find by navigating from the unicode.org home page, because MySQL uses the older 4.0.0 "allkeys" table, rather than the more recent 4.1.0 table. This is because we are very wary about changing ordering which affects indexes, lest we bring about situations such as that reported in Bug#16526, illustrated as follows:

```
mysql< CREATE TABLE tj (s1 CHAR(1) CHARACTER SET utf8 COLLATE utf8_unicode_ci);
Query OK, 0 rows affected (0.05 sec)
mysql> INSERT INTO tj VALUES ('ã##'),('ã##');
Query OK, 2 rows affected (0.00 sec)
Records: 2  Duplicates: 0  Warnings: 0
mysql> SELECT * FROM tj WHERE s1 = 'ã##';
+------+
| s1   |
+------+
| ã##  |
| ã##  |
+------+
2 rows in set (0.00 sec)
```

The character in the first result row is not the one that we searched for. Why did MySQL retrieve it? First we look for the Unicode code point value, which is possible by reading the hexadecimal number for the `ucs2` version of the characters:

```
mysql> SELECT s1, HEX(CONVERT(s1 USING ucs2)) FROM tj;
+------+-----------------------------+
| s1   | HEX(CONVERT(s1 USING ucs2)) |
+------+-----------------------------+
| ã##  | 304C                        |
| ã##  | 304B                        |
+------+-----------------------------+
2 rows in set (0.03 sec)
```

Now we search for `304B` and `304C` in the `4.0.0 allkeys` table, and find these lines:

```
304B  ; [.1E57.0020.000E.304B] # HIRAGANA LETTER KA
304C  ; [.1E57.0020.000E.304B][.0000.0140.0002.3099] # HIRAGANA LETTER GA; QQCM
```

The official Unicode names (following the "#" mark) tell us the Japanese syllabary (Hiragana), the informal classification (letter, digit, or punctuation mark), and the Western identifier (`KA` or `GA`, which happen to be voiced and unvoiced components of the same letter pair). More importantly, the *primary weight* (the first hexadecimal number inside the square brackets) is `1E57` on both lines. For comparisons in both searching and sorting, MySQL pays attention to the primary weight only, ignoring all the other numbers. This means that we are sorting `ã##` and `ã##` correctly according to the Unicode specification. If we wanted to distinguish them, we'd have to use a non-UCA (Unicode Collation Algorithm) collation (`utf8_bin` or `utf8_general_ci`), or to compare the `HEX()` values, or use `ORDER BY CONVERT(s1 USING sjis)`. Being correct "according to Unicode" isn't enough, of course: the person who submitted the bug was equally correct. We plan to add another collation for Japanese according to the JIS X 4061 standard, in which voiced/unvoiced letter pairs like `KA`/`GA` are distinguishable for ordering purposes.

### 9.15: Why don't CJK strings sort correctly in Unicode? (II)

If you are using Unicode (`ucs2` or `utf8`), and you know what the Unicode sort order is (see Chapter 9, *MySQL 6.0 FAQ — MySQL Chinese, Japanese, and Korean Character Sets*), but MySQL still seems to sort your table incorrectly, then you should first verify the table character set:

```
mysql> SHOW CREATE TABLE t\G
******************** 1. row ******************
Table: t
Create Table: CREATE TABLE `t` (
```

```
`s1` char(1) CHARACTER SET ucs2 DEFAULT NULL
) ENGINE=MyISAM DEFAULT CHARSET=latin1
1 row in set (0.00 sec)
```

Since the character set appears to be correct, let's see what information the `INFORMATION_SCHEMA.COLUMNS` table can provide about this column:

```
mysql> SELECT COLUMN_NAME, CHARACTER_SET_NAME, COLLATION_NAME
    -> FROM INFORMATION_SCHEMA.COLUMNS
    -> WHERE COLUMN_NAME = 's1'
    -> AND TABLE_NAME = 't';
+-------------+--------------------+-----------------+
| COLUMN_NAME | CHARACTER_SET_NAME | COLLATION_NAME  |
+-------------+--------------------+-----------------+
| s1          | ucs2               | ucs2_general_ci |
+-------------+--------------------+-----------------+
1 row in set (0.01 sec)
```

(See The `INFORMATION_SCHEMA COLUMNS` Table, for more information.)

You can see that the collation is `ucs2_general_ci` instead of `ucs2_unicode_ci`. The reason why this is so can be found using `SHOW CHARSET`, as shown here:

```
mysql> SHOW CHARSET LIKE 'ucs2%';
+---------+---------------+-------------------+--------+
| Charset | Description   | Default collation | Maxlen |
+---------+---------------+-------------------+--------+
| ucs2    | UCS-2 Unicode | ucs2_general_ci   |      2 |
+---------+---------------+-------------------+--------+
1 row in set (0.00 sec)
```

For `ucs2` and `utf8`, the default collation is "general". To specify a Unicode collation, use `COLLATE ucs2_unicode_ci`.

**9.16: Why are my supplementary characters rejected by MySQL?**

Before MySQL 6.0.4, MySQL does not support supplementary characters — that is, characters which need more than 3 bytes — for `UTF-8`. We support only what Unicode calls the *Basic Multilingual Plane / Plane 0*. Only a few very rare Han characters are supplementary; support for them is uncommon. This has led to reports such as that found in Bug#12600, which we rejected as "not a bug". With `utf8`, we must truncate an input string when we encounter bytes that we don't understand. Otherwise, we wouldn't know how long the bad multi-byte character is.

One possible workaround is to use `ucs2` instead of `utf8`, in which case the "bad" characters are changed to question marks; however, no truncation takes place. You can also change the data type to `BLOB` or `BINARY`, which perform no validity checking.

As of MySQL 6.0.4, Unicode support is extended to include supplementary characters by means of additional Unicode character sets: `utf16`, `utf32`, and 4-byte `utf8`. These character sets support supplementary Unicode characters outside the Basic Multilingual Plane (BMP).

**9.17: Shouldn't it be "CJKV"?**

No. The term "CJKV" (*Chinese Japanese Korean Vietnamese*) refers to Vietnamese character sets which contain Han (originally Chinese) characters. MySQL has no plan to support the old Vietnamese script using Han characters. MySQL does of course support the modern Vietnamese script with Western characters.

Bug#4745 is a request for a specialized Vietnamese collation, which we might add in the future if there is sufficient demand for it.

**9.18: Does MySQL allow CJK characters to be used in database and table names?**

This issue was fixed in MySQL 5.1, by automatically rewriting the names of the corresponding directories and files.

For example, if you create a database named 株価 on a server whose operating system does not support CJK in directory names, MySQL creates a directory named @0w@00a5@00ae. which is just a fancy way of encoding E6A5AE — that is, the Unicode hexadecimal representation for the 株価 character. However, if you run a `SHOW DATABASES` statement, you can see that the database is listed as 株価.

**9.19: Where can I find translations of the MySQL Manual into Chinese, Japanese, and Korean?**

A Simplified Chinese version of the Manual, current for MySQL 5.1.12, can be found at http://dev.mysql.com/doc/. The Japanese translation of the MySQL 4.1 manual can be downloaded from http://dev.mysql.com/doc/.

**9.20: Where can I get help with CJK and related issues in MySQL?**

The following resources are available:

• A listing of MySQL user groups can be found at http://dev.mysql.com/user-groups/.

• You can contact a sales engineer at the MySQL KK Japan office using any of the following:

```
Tel: +81(0)3-5326-3133
Fax: +81(0)3-5326-3001
Email: dsaito@mysql.com
```

- View feature requests relating to character set issues at http://tinyurl.com/y6xcuf.

- Visit the MySQL Character Sets, Collation, Unicode Forum. We are also in the process of adding foreign-language forums at http://forums.mysql.com/.