

# Business Intelligence Competency Center (BICC) Core System Documentation: Encoding & Unicode Considerations

This document describes in detail what you need to know when dealing with Encoding and Unicode in applications built on top of Microsoft solutions. To ensure the highest compatibility, use Unicode in international contexts and apply the best practices discussed in this document.

*This document is provided "as-is". Information and views expressed in this document, including URL and other Internet Web site references, may change without notice.*

*Some examples depicted herein are provided for illustration only and are fictitious. No real association or connection is intended or should be inferred.*

*This document does not provide you with any legal rights to any intellectual property in any Microsoft product. You may copy and use this document for your internal, reference purposes. You may modify this document for your internal, reference purposes.*

*© 2011 Microsoft. All rights reserved.*

# Main take-outs and Best Practices

The following statements provide guidance regarding Encodings and Unicode in your projects.

The best choice for a code page-specific server is to communicate only with the clients using the same code page. If you must communicate with clients using different code pages, the supported solution is to store your data in Unicode columns. If any one of these options is not feasible, the other alternative is to store the data in binary columns using the binary, varbinary, or varbinary (max) data types.

In international databases, character data must always use Unicode nchar and nvarchar data types in place of their non-Unicode equivalents (char and varchar). Nchar/nvarchar type is the Unicode data type used by Microsoft® SQL Server™, which can store any characters defined by Unicode.org. Char/varchar/text type is always associated with a code page, so the number of supported characters is limited. For example:

- Starting from Windows® 2000, Windows has full Unicode support implemented by using UTF-16 Encoding. Windows APIs take WCHAR\* as input which represent a Unicode String.
- Virtual C++ has WCHAR, which is Unicode char type.
- .Net String is Unicode String only encoded in UTF-16.
- Java String is Unicode String only encoded in UTF-16.
- Char type might have a data corruption issue if your client locale is different from the server locale.

The only benefit of using char type is the space saving for single byte code page. You can use the Data Compression feature in SQL Server 2008 if you care about disk space, which can save more space than using char type.

Put N' for your string literal in T-SQL. String literal N'abc' is a Unicode nvarchar string literal. Literal 'ab©' is a varchar literal, it always associates with a code page (string literal always use current database's collation). Suppose char © is not in the code page, you will get a question mark (?) when inserting the literal into a table, even if the column is nvarchar type.

Use nvarchar instead of nchar. The difference between nvarchar and nchar is the storage. A column with nchar(30) type always takes 60 bytes on the disk, even if the value is a single character. The data size for an nvarchar(30) type column is not fixed; it varies row by row or value by value. A single character value takes 2 bytes, and a value with 30 characters long takes 60 bytes to be stored. Another difference between nvarchar and nchar types is performance. Nchar type is always stored in fixed location in every row, which can be retrieved faster than nvarchar type which is stored in different locations for different rows. However, the benefit of less stored space for nvarchar types usually overcomes the cost of locating the value in a row.

If possible, avoid using a column collation which is different from the database's collation. You will have less collation conflicts. If you want a query to use a special collation's sorting rule, use explicitly the collate clause in that query.

Use Windows collations instead of SQL Server collations. The only exception here is the default collation for en\_US locale which is sql\_latin1\_general\_cp1\_ci\_as.

Never store UTF-8 string in varchar types, you will get data corruption.

String comparison always ignores trailing spaces (Unicode U+0020), no matter what collation you are using. LEN function always returns the number of characters excluding trailing spaces. DataLength function returns the storage size in terms of bytes.

Use \_BIN2 collation, instead of \_BIN collation if you want binary, code page based string comparison behavior.

# Globalization standards

---

## The code page issue

In a computer, characters are represented by different patterns of bits being either ON or OFF. There are 8 bits in a byte, and the 8 bits can be turned ON and OFF in 256 different patterns. A program that uses 1 byte to store each character can therefore represent up to 256 different characters by assigning a character to each of the bit patterns.

There are 16 bits in 2 bytes, and 16 bits can be turned ON and OFF in 65,536 unique patterns. A program that uses 2 bytes to represent each character can represent up to 65,536 characters.

Single-byte code pages are definitions of the characters mapped to each of the 256 bit patterns possible in a byte.

**Code pages define bit patterns for uppercase and lowercase characters, digits, symbols, and special characters such as !, @, #, or %.**

Each European language, such as German or Spanish, has its own single-byte code page. Although the bit patterns used to represent the Latin alphabet characters A through Z are the same for all the code pages, the bit patterns used to represent accented characters such as 'é' and 'á' vary from one code page to the next.

If data is exchanged between computers running different code pages, all character data must be converted from the code page of the sending computer to the code page of the receiving computer. If the source data has extended characters that are not defined in the code page of the receiving computer, data is lost.

When a database serves clients from many different countries, it is difficult to pick a code page for the database that contains all the extended characters required by all the client computers. Also, there is a lot of processing time spent doing the constant conversions from one code page to another.

Single-byte character sets are also inadequate to store all the characters used by many languages. For example, some Asian languages have thousands of characters, so must use two bytes per character.

Double-byte character sets have been defined for these languages. Still, each of these languages have their own code page, and there are difficulties in transferring data from a computer running one double-byte code page to a computer running another.

## Standard Organizations

### **ANSI**

The **American National Standards Institute** ([www.ansi.org](http://www.ansi.org)) is an organization that standardizes various areas, both public and private. It is an acronym for American National Standards Institute. ANSI has standardized many areas of computing.

One of the standards in computers was the character set (letters, numbers, and symbols) that a computer can use. This is a standardized encoding table (a code page) covering all uppercase and lowercase English letters, digits, punctuation characters, as well as some special and control characters.

At the very beginning computers could work with 128 different combinations (one bit was reserved for other purposes).

The ASCII code page covered 128 characters, but newer computer systems were able to work with 256 codes and engineers soon noticed that 128 codes are not sufficient for all characters. Even the diacritical marks of the Western European languages could not be covered.

So the standard committees (ANSI, ISO) and computer companies (IBM, Apple, Microsoft) started extending the ASCII code page with various character sets.

The complementary 128 codes had been filled with graphical symbols, mathematical signs, Western European diacritical marks etc. Each organization elaborated its own "standards".

ANSI and Microsoft invented the code page 1252 (ANSI Latin-1), the International Standards Organization (ISO) established the ISO 8859-1 (ISO Latin-1), IBM developed the code page 850 (IBM Latin-1), Apple created the Macintosh Roman character set, etc.

Each Windows system has a default ANSI code page according to system regional settings (932 commonly known as Shift JIS for Japan computers, 1252 commonly known as ANSI Latin-1 for Western European computer).

### **Unicode Consortium**

The Unicode Consortium is a non-profit organization founded to develop, extend and promote use of the Unicode Standard, which specifies the representation of text in modern software products and standards.

The membership of the consortium represents a broad spectrum of corporations and organizations in the computer and information processing industry.

The consortium is supported financially solely through membership dues. Membership in the Unicode Consortium is open to organizations and individuals anywhere in the world who support the Unicode Standard and wish to assist in its extension and implementation.

The Unicode standard is a character coding system designed to support written texts of diverse modern, classical and historical languages. It's based on double-byte character encoding, so it can enumerate 65,536 characters. It's hopefully the "one and only" future standard, and it may solve the "code page soup" problem.

Unicode is compatible with the ISO 10646 standard. The current version 5.2 includes 107,361 coded characters used in written languages of the Americas, Europe, the Middle East, Africa, India, Asia, and Pacifica.

The industry is converging on Unicode for all internationalization. For example: Microsoft Windows Operating Systems have been built on a base of Unicode since Windows NT; AIX, Sun and HP/UX offer Unicode support.

All web standards: HTML, XML, WML etc. are supporting or requiring Unicode. Most versions of Web Browser support Unicode. Sybase, Oracle, SQL, DB2 all offer or are developing Unicode support.

*Unicode is the only way of globalization.*

**UTF-7 encoding:** (UTF-7 stands for UCS Transformation Format, 7-bit form). This encoding supports all Unicode character values, and can also be accessed as code page 65000.

**UTF-8 encoding:** (UTF-8 stands for UCS Transformation Format, 8-bit form). This encoding supports all Unicode character values, and can also be accessed as code page 65001.

**UCS-2 encoding:** Some applications (especially those that are Web based) must deal with Unicode data that is encoded with the UTF-8 encoding method. SQL Server uses a different Unicode encoding (UCS-2) and does not recognize UTF-8 as valid character data.

**UTF-16 encoding:** This encoding stores the basic Unicode characters using single 16 bit units and others characters using two 16 bit units. UTF-16 is the primary encoding mechanism used by Microsoft Windows Client and Server OS.

**Note:** *The UCS-2 encoding scheme is actually a subset of the UTF-16 scheme. Every UCS-2 encoded code point is identical to the encoding of the same code point in UTF-16. Also, most new implementations using the Unicode standard now employ UTF-16, UTF-8 or UTF-32 instead of UCS-2.*

# Encodings in SQL Server

## Collation

The physical storage of character strings in Microsoft SQL Server 2008 is controlled by collations.

A collation specifies the bit patterns that represent each character and the rules by which characters are sorted and compared.

Each SQL Server collation specifies three properties:

- The sort order to use for Unicode data types (nchar and nvarchar). A sort order defines the sequence in which characters are sorted, and the way characters are evaluated in comparison operations.
- The sort order to use for non-Unicode character data types (char, varchar, and text).
- The code page used to store non-Unicode character data.

SQL Server 2008 supports objects that have different collations being stored in a single database. Separate SQL Server 2008 collations can be specified down to the level of columns. Each column in a table can be assigned different collations. Versions of SQL Server before SQL Server 2000 support only one collation for each instance of SQL Server. All databases and database objects created in an instance of SQL Server 7.0 or earlier have the same collation.

SQL Server 2008 supports the following code pages.

Code page	Description
<b>437</b>	MS-DOS U.S. English

<b>850</b>	Multilingual (MS-DOS Latin1)
<b>874</b>	Thai
<b>932</b>	Japanese
<b>936</b>	Chinese (Simplified)
<b>949</b>	Korean
<b>950</b>	Chinese (Traditional)
<b>1250</b>	Central European
<b>1251</b>	Cyrillic
<b>1252</b>	Latin1 (ANSI)
<b>1253</b>	Greek
<b>1254</b>	Turkish
<b>1255</b>	Hebrew
<b>1256</b>	Arabic
<b>1257</b>	Baltic
<b>1258</b>	Vietnamese

Microsoft SQL Server collations can be categorized in two groups: Windows collations and SQL Server collations.

## Windows Collations

Windows collations are collations defined for SQL Server to support the Windows system locales available for the operating system on which SQL Server instances are installed. For information on new Windows collations support (collations based on Windows system locales) added in SQL Server 2008 and all other Windows collations, see [Windows Collation Name \(Transact-SQL\)](#).

By specifying a Windows collation for SQL Server, the instance of SQL Server uses the same code pages and sorting and comparison rules as an application that is running on a computer for which you have specified the associated Windows locale. For example, the French Windows collation for SQL Server matches the collation attributes of the French locale for Windows.

There are more Windows locales than there are SQL Server Windows collations. The names of Windows locales are based on language and territory; for example, French (Canada). However, several languages share common alphabets and rules for sorting and comparing characters. For example, several Windows locales, including all the Portuguese and English Windows locales, use the Latin1 code page (1252) and follow a common set of rules for sorting and comparing characters. Latin1\_General, the SQL Server-supported Windows collation based on the 1252 code page and sorting rules, supports all of these Windows locales.

Also, Windows locales specify attributes that are not covered by SQL Server supported Windows collations such as currency, date, and time formats. Because countries and regions such as Great Britain and the United States have different currency, date, and time formats, they require different Windows locales. They do not require different SQL Server collations, because they have the same alphabet and rules for sorting and comparing characters. In SQL Server, Windows collations are combined with suffixes that define sorting and comparison rules based on case, accent, kana, and width sensitivity. For more information about these suffixes, see [Windows Collation Sorting Styles](#).

These are some examples of Windows collation names:

- Latin1\_General\_CI\_AS

Collation uses the Latin1 General dictionary sorting rules, code page 1252, is case-insensitive and accent-sensitive.

- Estonian\_CS\_AS

Collation uses the Estonian dictionary sorting rules, code page 1257, is case-sensitive and accent-sensitive.

Please see the appendix for a list of Microsoft SQL Server 2008 Windows collations.

### **SQL Server collations**

SQL Server collations are a compatibility option to match the attributes of common combinations of code-page number and sort orders that have been specified in earlier versions of SQL Server. Many of these collations support suffixes for case, accent, kana, and width sensitivity, but not always. For more information, see [Using SQL Server Collations](#).

SQL Server collations apply non-Unicode sorting rules to non-Unicode data, and Unicode sorting rules to Unicode data, by using a corresponding Windows collation for the Unicode data. This difference can cause inconsistent results for comparisons of the same characters. Therefore, if you have a mix of Unicode and non-Unicode columns in your database, they should all be defined by using Windows collations so that the same sorting rules are used across Unicode and non-Unicode data.

To maintain compatibility with earlier versions of SQL Server, or applications that were developed with SQL Server collations in earlier versions of SQL Server, SQL Server offers the SQL\_Latin1\_General\_CP1\_CI\_AS collation as the default collation for server installations on computers that use the English (United States) Windows system locale.

There can be differences in performance between Windows collations and SQL Server collations. For more information, see [Storage and Performance Effects of Unicode](#).

## **Character data types**

SQL Server 2008 supports two categories of character data types:

- The Unicode data types `nchar` and `nvarchar`. These data types use the Unicode character representation. Code pages do not apply to these data types.
- The non-Unicode character data types `char`, `varchar`, and `text`. These data types use the character representation scheme defined in a single or double-byte code page.

Collations do not control the code page used for Unicode columns, only attributes such as comparison rules and case sensitivity.

## International Data and Unicode

The easiest way to manage character data in international databases is to always use the Unicode nchar and nvarchar data types in place of their non-Unicode equivalents (char, varchar, and text).

If all the applications that work with international databases also use Unicode variables instead of non-Unicode variables, character translations do not have to be performed anywhere in the system. All clients will see exactly the same characters in data as all other clients.

For systems that could use single-byte code pages, the fact that Unicode data needs twice as much storage space as non-Unicode character data is at least partially offset by eliminating the need to convert extended characters between code pages. Systems using double-byte code pages do not have this issue.

For instance, **SQL Server stores all textual system catalog data in columns having Unicode data types**. The names of database objects such as tables, views, and stored procedures are stored in Unicode columns. This allows applications to be developed using only Unicode, which avoids all issues with code page conversions.

## Sort order

A sort order specifies the rules used by SQL Server to interpret, collate, compare, and present character data. For example, a sort order defines whether 'a' is less than, equal to, or greater than 'b'. A sort order defines whether the collation is case-sensitive, for example whether 'm' is equal or not equal to 'M'. It also defines if the collation is accent-sensitive, for example whether 'á' is equal or not equal to 'ä'.

SQL Server uses two sort orders with each collation, one for Unicode data and another one for the character code page.

Many SQL Server collations use the same code page, but have a different sort order for the code page.

This allows applications to choose:

- Whether characters will simply be sorted based on the numeric value represented by their bit patterns. Binary sorting is fastest because SQL Server does not have to make any adjustments and can use fast, simple sorting algorithms. Binary sort orders are always case-sensitive. Because the bit patterns in a code page may not be arranged in the same sequence as defined by the dictionary rules for a specific language, binary sorting sometimes does not sort characters in a sequence users who speak that language might expect.
- Between case-sensitive or case-insensitive behavior.
- Between accent-sensitive or accent-insensitive behavior.

# Encodings in the .Net Framework

The .NET Framework is a platform for building, deploying, and running Web services and applications that provide a highly productive, standards-based, multi-language environment for integrating existing or legacy investments with next-generation applications and services.

The .NET Framework uses Unicode UTF-16 to represent characters, although in some cases it uses UTF-8 internally.

The System.Text namespace provides classes that allow you to encode and decode characters, with support that includes the following encodings:

- Unicode UTF-16 encoding. Use the `UnicodeEncoding` class to convert characters to and from UTF-16 encoding.
- Unicode UTF-8 encoding. Use the `UTF8Encoding` class to convert characters to and from UTF-8 encoding.
- ASCII encoding. ASCII encodes the Latin alphabet as single 7-bit characters. Because this encoding only supports character values from U+0000 through U+007F, in most cases it is inadequate for internationalized applications. You can use the `ASCIIEncoding` class to convert characters to and from ASCII encoding whenever you need to interoperate with legacy encodings and systems.
- Windows/ISO Encodings. The [System.Text.Encoding](#) class provides support for a wide range of Windows/ISO encodings.

The .NET Framework provides support for data encoded using code pages. You can use the `Encoding.GetEncoding Method (Int32)` to create a target encoding object for a specified code page. Specify a code page number as the `Int32` parameter.

The one additional type of support introduced to ASP.NET is the ability to clearly distinguish between file, request, and response encodings. To set the encoding in ASP for code, page directives, and configuration files, you'll need to do the following.

In code:

```
Response.ContentEncoding= <value>
```

```
Request.ContentEncoding= <value>
```

```
File.ContentEncoding= <value>
```

### In Page directive:

Several new directives have been added to ASP.NET. The Language attribute must now be placed within a Page directive, as shown in the following example:

```
<%@Page Language="VB" Codepage="932"%>  
<%@OutputCache Duration="10 VaryByParam="location"%>
```

But for migration purposes, the shorter ASP-style syntax is also supported for the Page directive only.

```
<%@ Language="VB" Codepage="932"%>  
<%@OutputCache Duration="10 VaryByParam="location"%>
```

Directives can be located anywhere in an .aspx file, but standard practice is to place them at the beginning of the file. Case is not important in ASP.NET directive statements, and quotes are not required around the attribute values.

### In a configuration file:

You can set the globalization settings in the application web.config file. If not present, the settings are the ones settled in the machine.config file, and if none, the system's default ANSI code page.

Here is the syntax in the configuration file:

```
<configuration>
  <globalization
    enableClientBasedCulture="true|false"
    requestEncoding="any valid encoding string"
    responseEncoding="any valid encoding string"
    fileEncoding="any valid encoding string"

    responseHeaderEncoding = "any valid encoding string"
    resourceProviderFactoryType = string
    enableBestFitResponseEncoding = "true|false"

    culture="any valid culture string"
    uiCulture="any valid culture string"/>
</configuration>
```

---

For more information about setting globalization in configuration files, please go to <http://msdn.microsoft.com/en-us/library/hy4kkhe0.aspx>.

# Encodings in Web Pages

However, Web pages currently consist of content that can be in Windows or other character-encoding schemes besides Unicode. Therefore, when form or query-string values come in from the browser in an HTTP request, they must be converted from the character set used by the browser into Unicode for processing by the .Net Framework. Similarly, when output is sent back to the browser, any strings returned by scripts must be converted from Unicode back to the code page used by the client.

Generally speaking, there are four different ways of setting the character set or the encoding of a Web page:

- **Windows code pages or ISO character encodings:** With this approach, you can select from the list of supported code pages to create your Web content. The downside of this approach is that you are limited to languages that are included in the selected character set, making true multilingual Web content impossible. This limits you to a single-script Web page.
- **Number entities:** Number entities can be used to represent a few symbols out of the currently selected code page or encoding. Let's say, for example, you have decided to create a Web page using the previous approach with the Latin ISO charset 8859-1. Now you also want to display some Greek characters in a mathematical equation; Greek characters, however, are not part of the Latin code page. Unfortunately, this approach makes it impossible to compose large amounts of text and makes editing your Web content very hard.
- **UTF-16:** Unlike Win32 applications where UTF-16 is by far the best approach, for Web content UTF-16 can be used safely only on Windows NT networks that have full Unicode support. Therefore, this is not a suggested encoding for Internet sites where the capabilities of the client Web browser as well the network Unicode support are not known.
- **UTF-8: This Unicode encoding is the best and safest approach for multilingual Web pages.** It allows you to encode the whole repertoire of Unicode characters. Also, all versions of Internet Explorer 4 and later support this encoding, which is not restricted to network or wire capabilities. The UTF-8 encoding allows you to create multilingual Web content without having to change the encoding based on the target language.

## Setting and encoding Web Pages

Since Web content is currently based on Windows or other encoding schemes, you'll need to know how to set and manipulate encodings. The following describes how to do this for HTML pages, Active Server Pages (ASP), ASP.Net, and XML pages.

Internet Explorer uses the character set specified for a document to determine how to translate the bytes in the document into characters on the screen or on paper.

By default, Internet Explorer uses the character set specified in the HTTP content type returned by the server to determine this translation. If this parameter is not given, Internet Explorer uses the character set specified by the meta-element in the document, taking into account the user's preferences if no meta-element is specified.

To apply a character set to an entire document, you must insert the meta-element before the body element. For clarity, it should appear as the first element after the head, so that all browsers can translate the meta-element before the document is parsed. The meta-element applies to the document containing it. This means, for example, that a compound document (a document consisting of two or more documents in a set of frames) can use different character sets in different frames. Here is how it works:

```
<META HTTP-EQUIV="Content-Type" CONTENT="text/html; charset=<value>">
```

You substitute with any supported character-set-friendly name (for example, UTF-8) or any code-page name (for example, windows 1251). See the appendix for more code page name.

# Working with Unicode Data in Microsoft Platform

## Unicode Basics

Storing data in multiple languages within one database is difficult to manage when you use only character data and code pages. It is also difficult to find one code page for the database that can store all the required language-specific characters. Additionally, it is difficult to guarantee the correct translation of special characters when being read or updated by different clients running various code pages. **Databases that support international clients should always use Unicode data types instead of non-Unicode data types.**

For example, consider a database of customers in North America that must handle three major languages:

- Spanish names and addresses for Mexico
- French names and addresses for Quebec
- English names and addresses for the rest of Canada and the United States

When you use only character columns and code pages, you must take care to make sure the database is installed with a code page that will handle the characters of all three languages. You must also take care to guarantee the correct translation of characters from one of the languages when read by clients running a code page for another language.

With the growth of the Internet, it is even more important to support many client computers that are running different locales. Selecting a code page for character data types that will support all the characters required by a worldwide audience would be difficult.

The easiest way to manage character data in international databases is to always use the Unicode nchar, nvarchar, and nvarchar(max) data types, instead of their non-Unicode equivalents, char, varchar, and text.

Unicode is a standard for mapping code points to characters. Because it is designed to cover all the characters of all the languages of the world, there is no need for different code pages to handle different sets of characters. SQL Server supports the Unicode Standard, Version 3.2.

If all the applications that work with international databases also use Unicode variables instead of non-Unicode variables, character translations do not have to be performed anywhere in the system. Clients will see the same characters in the data as all other clients.

SQL Server stores all textual system catalog data in columns having Unicode data types. The names of database objects, such as tables, views, and stored procedures, are stored in Unicode columns. This enables applications to be developed by using only Unicode, and helps avoid all issues with code page conversions.

## Storage and Performance Effects of Unicode

SQL Server stores Unicode data by using the UCS-2 encoding scheme. Under this mechanism, all Unicode characters are stored by using 2 bytes.

The difference in storing character data between Unicode and non-Unicode depends on whether non-Unicode data is stored by using double-byte character sets. All non-East Asian languages and the Thai language store non-Unicode characters in single bytes. Therefore, storing these languages as Unicode uses two times the space that is used specifying a non-Unicode code page. On the other hand, the non-Unicode code pages of many other Asian languages specify character storage in double-byte character sets (DBCS). Therefore, for these languages, there is almost no difference in storage between non-Unicode and Unicode.

The following table shows the non-Unicode code pages that specify character data storage in double-byte character sets.

Language	Code page
<b>Simplified Chinese</b>	936
<b>Traditional Chinese</b>	950
<b>Japanese</b>	932
<b>Korean</b>	949

The effect of Unicode data on performance is complicated by a variety of factors that include the following:

- The difference between Unicode sorting rules and non-Unicode sorting rules
- The difference between sorting double-byte and single-byte characters
- Code page conversion between client and server

SQL Server performs string comparisons of non-Unicode data defined with a Windows collation by using Unicode sorting rules. Because these rules are much more complex than non-Unicode sorting rules, they are more resource-intensive. So, although Unicode sorting rules are frequently more expensive, there is generally little difference in performance between Unicode data and non-Unicode data defined with a Windows collation.

The only case when SQL Server uses non-Unicode sorting rules is on non-Unicode data that is defined by using SQL Server collation. Sorts and scans in this instance are generally faster than when Unicode sorting rules apply. Unicode sorting rules apply to all Unicode data, defined by using either a Windows collation or SQL Server collation.

Of secondary importance, sorting lots of Unicode data can be slower than non-Unicode data, because the data is stored in double bytes. On the other hand, sorting Asian characters in Unicode is faster than sorting Asian DBCS data in a specific code page, because DBCS data is actually a mixture of single-byte and double-byte widths, while Unicode characters are fixed-width.

Other performance issues are primarily determined by the issue of converting the encoding mechanism between an instance of SQL Server and the client. Generally, the effects on performance of client/server code-page conversion are negligible. Nevertheless, you should understand what is occurring at this layer.

The **ODBC API, version 3.6 or earlier, and the DB-Library API do not recognize Unicode**. For clients that use data access methods defined by these APIs, resources are used to implicitly convert Unicode data to the client code page. Also, **there is a risk of data corruption on the client side** when the client code page does not recognize certain Unicode characters.

**Later versions of ODBC**, starting with Microsoft Data Access Components version 2.7 that was included with SQL Server version 7.0, and OLE DB and ADO **are Unicode aware** and assume a UCS-2 encoding mechanism. Therefore, if the application is Unicode enabled, there are no conversion issues when you work strictly with Unicode data from an instance of SQL Server. If a client is using a Unicode-enabled API but the data storage mechanism in the instance of SQL Server is not Unicode, there are no conversion issues. However, there is a risk that any data insert or update operations will be corrupted if the code points for any character cannot be mapped to the SQL Server code page.

## Unicode Best Practices

Deciding whether to store non-DBCS data as Unicode is generally determined by an awareness of the effects on storage, and about how much sorting, conversion, and possible data corruption might happen during client interactions with the data. Sorting and conversion may affect performance, depending on where it occurs. However, for most applications the effect is negligible. Databases with well-designed indexes are especially unlikely to be affected. However, data corruption will affect not only the integrity of an application and database, but also the business as a whole.

Considering this trade-off, storing character data in a specific code page may make sense if both of the following are true:

- Conserving storage space is an issue, because of hardware limitations. Or, you are performing frequent sorts of lots of data, and testing indicates that a Unicode storage mechanism severely affects performance.
- You are sure the code pages of all clients accessing this data match yours, and that this situation will not unexpectedly change.

Most of the time, the decision to store character data, even non-DBCS data, in Unicode should be based more on business needs instead of performance. In a global economy that is encouraged by rapid growth in Internet traffic, it is becoming more important than ever to support client computers that are running different locales. Additionally, it is becoming increasingly difficult to pick a single code page that supports all the characters required by a worldwide audience.

## Server-side Programming with Unicode

To make a database Unicode-aware involves defining Unicode-aware client interactions in addition to using the **nchar**, **nvarchar**, and **nvarchar(max)** data types to define Unicode storage. You can define Unicode-aware client interactions by performing the following on the database server side:

- Switch from non-Unicode data types to Unicode data types in table columns and in CONVERT() and CAST() operations.
- Substitute using ASCII() and CHAR() functions with their Unicode equivalents, UNICODE() and NCHAR().
- Define variables and parameters of stored procedures and triggers in Unicode.
- Prefix Unicode character string constants with the letter N.

### Using UNICODE(), NCHAR(), and other Functions

The ASCII() function returns the non-Unicode character code of the character passed in. Therefore, use the counterpart **UNICODE()** function for Unicode strings where you would use the ASCII function on non-Unicode strings. The same is true of the CHAR function; NCHAR is its Unicode counterpart.

Because the **SOUNDEX()** function is defined based on English phonetic rules, it is not meaningful on Unicode strings unless the string contains only the Latin characters A through Z and a through z.

ASCII, CHAR, and SOUNDEX can be passed Unicode parameters, but these arguments are implicitly converted to non-Unicode strings. This could cause the possible loss of Unicode characters before processing, because these functions operate on non-Unicode strings by definition.

Besides the UNICODE() and NCHAR() functions, **the following string manipulation functions support Unicode wherever possible:**

- CHARINDEX()
- LEFT()
- LEN()
- UPPER()
- LOWER()
- LTRIM()
- RTRIM()
- PATINDEX()
- REPLACE()
- QUOTENAME()
- REPLICATE()
- REVERSE()
- STUFF()
- SUBSTRING()
- UNICODE()

These functions accept Unicode arguments, respect the 2-byte character boundaries of Unicode strings, and use Unicode sorting rules for string comparisons when the input parameters are Unicode.

### **Defining Parameters in stored Procedures**

**Defining parameters with a Unicode data type guarantees that client requests or input are implicitly converted to Unicode on the server and not corrupted in the process.** If the parameter is specified as an OUTPUT parameter, a Unicode type also minimizes the chance of corruption on its way back to the client.

In the following stored procedure, the variable is declared as a Unicode data type.

```
CREATE PROCEDURE Product_Info
    @name nvarchar(40)
AS
SELECT p.ListPrice, v.Name
FROM Production.Product p
INNER JOIN Purchasing.ProductVendor pv
ON p.ProductID = pv.ProductID
INNER JOIN Purchasing.Vendor v
ON pv.VendorID = v.VendorID
WHERE p.Name = @name;
```

### ***Using the N Prefix***

**Unicode string constants that appear in code executed on the server, such as in stored procedures and triggers, must be preceded by the capital letter N.** This is true even if the column being referenced is already defined as Unicode. Without the N prefix, the string is converted to the default code page of the database. This may not recognize certain characters.

For example, the stored procedure created in the previous example can be executed on the server in the following way:

```
EXECUTE Product_Info @name = N'Chain'
```

The requirement to use the N prefix applies to both string constants that originate on the server and those sent from the client.

## Client-side Programming with Unicode

The topics in this section explain how to preserve the integrity of character data when you program client-side database applications.

### **Managing Data Conversion between Client/Server Code Pages**

This topic describes how to preserve the integrity of character data when the database does not store the character data by using Unicode data types and when your client-side applications that interact with the data are also not Unicode-aware. In this situation, the code page of your data storage and the code page of the client-side application should be the same. If these code pages differ, the conversion that occurs between the client and the server might cause the loss of some characters.

Disabling the AutoTranslate feature of the SQL Server ODBC driver to insert data defined by a different code page from the server is not supported. Also, even if AutoTranslate is disabled, it does not prevent code page translation for SQL language events. The result is that if the client and database code pages do not match, code page translation will generally be applied to any non-Unicode character string that is sent to or from the server.

If you can, you should avoid this situation. **The best choice for a code page-specific server is to communicate only with the clients using the same code page.** The second-best choice is to use another code page that has almost the same character set. For example, code page 1252 (Latin1) and code page 850 (Multilingual Latin1) can store almost the same set of characters, so that most characters in these two code pages can be converted from one code page to another without data loss.

If you must communicate with clients using different code pages, the supported solution is to store your data in Unicode columns. If any one of these options is not feasible, the other alternative is to store the data in binary columns using the binary, varbinary, or varbinary(max) data types. However, binary data can only be sorted and compared in binary order. This makes it less flexible than character data.

## **Managing Data Conversion between a Unicode Server and a Non-Unicode Client**

This topic describes how to preserve the integrity of character data when the server-side data storage is in Unicode, but the client-side application that interacts with the data uses a specific code page.

### **Data Input**

When non-Unicode data is sent from the client to be stored on the server in Unicode, data from any client with any code page can be stored correctly if one of the following conditions is true:

- Character strings are sent to the server as parameters of a remote procedure call (RPC).
- String constants are preceded with the capital letter N. This is required regardless of whether your client-side application is Unicode-aware. Without the N prefix, SQL Server will convert the string to the code page that corresponds to the default collation of the database. Any characters not found in this code page will be lost.

### **Data Retrieval**

If the client application is not Unicode-enabled and retrieves the data into non-Unicode buffers, a client will only be able to retrieve or modify data that can be represented by the client machine's code page. This means that ASCII characters can always be retrieved, because the representation of ASCII characters is the same in all code pages, while any non-ASCII data depends on code-page-to-code-page conversion.

For example, suppose you have an application that is currently running only in the United States (U.S.), but is deployed to Japan. Because the SQL Server database is Unicode-aware, both the English and Japanese text can be stored in the same tables, even though the application has not yet been modified to deal with text as Unicode. As long as the application complies with one of the two previous options, Japanese users can use the non-Unicode application to input and retrieve Japanese data, and U.S. users can input and retrieve English data. All data from both sets of users is stored intact in the same column of the database and represented as Unicode. In this situation, a Unicode-enabled reporting application that generates reports that span the complete data set can be

deployed. However, English users cannot view the Japanese rows, because the application cannot display any characters that do not exist in their code page (1252).

This situation might be acceptable if the two groups of users do not have to view each other's records. If an application user must be able to view or modify records with text that cannot be represented by a single code page, there is no alternative but to modify the application so that it can use Unicode.

### **Web-based Applications**

If the client-side program is Web-based or connects to an Active Server Pages (ASP) page, there are metadata specifications on both the client-side HTML page and the server-side ASP page. These specifications must be made to specify how character strings should be converted between the server, the ASP engine, and the client browser.

On the client side HTML page, the META attribute must specify that the character set data should be converted to the encoding scheme of the client by specifying a CHARSET code. For example, the following HTML page instructs the client to convert character data to the 950 (Chinese Traditional) code page by specifying big5 as the CHARSET code.

```
<HTML>
<HEAD>
<META HTTP-EQUIV="Content-Type" CONTENT="text/html; CHARSET=big5">
<!--
-->
</HEAD>
```

```
<BODY>

<!--
body
-->

</BODY>

</HTML>
```

On the server-side ASP page, you must instruct the ASP Web application what code page the client browser is using. You can specify the `Session.CodePage` property, or the `@CodePage` directive. These methods will handle the conversion of data from server to client and also both GET and POST client requests. In the following examples, both methods are used to specify conversion to and from the code page of the client, which is 950 (Chinese Traditional).

```
<%@ Language=VBScript codepage=950 %>

<% Session.CodePage=950 %>
```

And finally, you must remember to prefix any string literals with the letter N.

### **Managing Data Conversion between Unicode Encoding Schemes**

This topic describes how to preserve the integrity of character data when both server-side data storage and the client application that interacts with the data are Unicode-enabled, but use different Unicode encoding schemes. SQL Server stores Unicode in the UCS-2 encoding scheme. However, many clients process Unicode in another encoding scheme, generally UTF-8. This scenario frequently occurs for Web-based applications.

Because you are essentially still converting from one encoding scheme to another; many of the same solutions discussed in the previous topics also apply. Unicode character string constants sent to the server must be preceded with a capital N. For Web-based applications, you specify the CHARSET code under the META attribute of the client-side HTML page. For example, specify CHARSET = utf-8 if the Unicode encoding scheme is UTF-8. On the server side, specify the encoding scheme of the client by using the Session.CodePage property or the @Codepage directive. For example, codepage=65001 specifies a UTF-8 encoding scheme. If you follow these directions, Internet Information Services (IIS) 5.0 or later versions will seamlessly handle the conversion from UTF-8 to UCS-2 and back without additional effort on your part.

In Visual Basic applications, character strings are processed in the UCS-2 encoding scheme. Therefore, you do not have to specify encoding scheme conversion explicitly between these applications and an instance of SQL Server.

## Using Unicode with bcp and OPENROWSET

To prevent character loss when you use the bcp utility to copy data between servers that have different code pages, you can specify that the data be copied in Unicode format. For more information, see [Copying Data Between Different Collations](#).

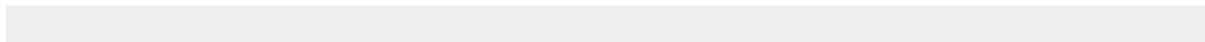
Unicode format can also be specified when you use OPENROWSET to access external data. For more information, see [OPENROWSET \(Transact-SQL\)](#).

## Using Unicode with XML Data

SQL Server stores XML data using the UTF-16 encoding scheme. Because UTF-16 data is variable-width, it is processed according to a byte-oriented protocol. This means that UTF-16 data can be treated in a way that is independent of the byte ordering on different computers (little endian versus big endian). Therefore, UTF-16 is well-suited for traversing different computers that use different encodings and byte-ordering systems. Because XML data is typically shared widely across networks, it makes sense to maintain the default UTF-16 storage of XML data in your database, and when you export XML data to clients.

If you must specify a different encoding, you can use FOR XML requests and specify the following:

- The Output Encoding property of an XML-formatted data stream Response object in Active Server Pages (ASP).
- For example, the following ASP code tells the browser to display an incoming XML data stream in UCS-2:



```
<% cmdXML.Properties("Output Encoding") = "UCS-2" %>
```

- An output encoding in a URL when you make an HTTP request.
- The following example specifies UCS-2 as the output encoding of the XML document returned by this request:

```
http://IISServer/nwind?sql=SELECT+*+FROM+Customers+FOR+XML+AUTO&outputencoding=UCS-2
```

- An output encoding in an XML template or style sheet.
- The following example specifies UCS-2 as the output encoding in the header of this XML template document:

```
<?xml version='1.0' encoding='UCS-2'?>  
<root xmlns:sql='urn:schemas-microsoft-com:xml-sql'  
sql:xsl='MyXSL.xsl'>  
<sql:query>
```

```
SELECT FirstName, LastName FROM Employees FOR XML AUTO
```

```
</sql:query>
```

```
</root>
```

# SQL Server Globalization Development Recommendations

## Use nchar/nvarchar instead of char/vartimetype

Nchar/nvarchar type is the Unicode data type used by SQL Server, which can store any characters defined by Unicode.org. Char/varchar type is always associated with a code page, so the number of supported characters is limited. For example:

- Starting from Windows 2000, Windows has full Unicode support implemented by using UTF-16 Encoding. Windows APIs take WCHAR\* as input which represent a Unicode String.
- Virtual C++ has WCHAR, which is Unicode char type.
- .Net String is Unicode String only encoded in UTF-16.
- Java String is Unicode String only encoded in UTF-16.
- Char type might have a data corruption issue if your client locale is different from the server locale.

## Deprecate old code page technique

Use Unicode in your application and in your SQL Server database. The only benefit of using char type is the space saving for single byte code page. You can always use the Data Compression feature in SQL Server 2008 if you care about disk space, which can save more space than using char type.

## Use N for your string literal in T-SQL

Remember to put N' for your string literal in T-SQL. String literal N'abc' is a Unicode nvarcharstring literal. Literal 'ab©' is a varchar literal, it always associates with a code page (string literal always use current database's collation). Suppose char © is not in the code page, you will get a question mark (?) when inserting the literal into a table, even the column is nvarchar type.

### **Use nvarchar instead of nchar**

Use nvarchar instead of nchar. The difference between nvarchar and nchar is the storage. A column with nchar(30) type always take 60 bytes to store on the disk, even the value is a single character. The data size for an nvarchar(30) type column is not fixed; it varies row by row or value by value. A single character value takes 2 bytes to store, and a value with 30 characters long takes 60bytes to be stored. Another different between nvarchar and nchar type is the performance. Nchar types always stored in fixed location in every row, which can be retrieved faster than nvarchar type which is stored in different location for different row. However, I believe the benefit of less stored space for nvarchar types usually overcomes the cost of locating the value in a row.

### **Avoid column collation that is different from the dataset collation**

If possible, avoid using a column collation which is different with the database's collation. You will have less collation conflict. If you want a query use a special collation's sorting rule, use explicit collate clause in that query.

### **Use Windows Collation instead of SQL Server Collation**

Use Windows collation instead of SQL Server collation. The only exception here is the default collation for en\_US locale which is sql\_latin1\_general\_cp1\_ci\_as.

### **Never store UTF-8 string in varchar types**

Never store UTF-8 string in varchar types, you *will* get data corruption.

### **String comparison and trailing spaces**

String comparison always ignores trailing spaces (Unicode U+0020), no matter what collation you are using.

### **LEN Function reminder**

LEN function always returns the number of characters excluding the trailing spaces. DataLength function returns the storage size in term of bytes.

### **\_BIN2 Collation**

Use \_BIN2 collation, instead of \_BIN collation if you want binary, code page base string comparison behavior.

# Working with other data sources in Microsoft Platform

The following section deals with inserting data from source systems into Microsoft environment. Each part describes the possibilities and constraints of the specified method.

## Working with UTF-16 source data in SQL Server (Unicode)

SQL Server offers the following possibilities when working with UTF-16 data sets:

- SQL Server can store any Unicode characters in UTF-16 encoding. The reason is that the storage format for UCS-2 and UTF-16 are the same.
- SQL Server can display any Unicode characters in UTF-16 encoding. The reason is that we internally call Windows functions to display characters, the Windows functions and fonts can take care of the supplementary character (a character take 4 bytes in UTF-16) correctly.
- SQL Server can input any Unicode characters in UTF-16 encoding. The reason is that we internally call Windows IMEs (Input Method Editors) to input, the Windows IMEs can take care of the supplementary character (a character take 4 bytes in UTF-16) correctly.
- SQL Server can sort/compare any defined Unicode characters in UTF-16 encoding. Note, not all code points are mapped to valid Unicode characters. For example, The Unicode Standard, Version 5.1 defines code points for around 10,000 characters. All these characters can be compared/sorted in SQL Server latest version: SQL Server 2008.

SQL Server has the following constraints when working with UTF-16 data sets:

SQL Server **cannot detect invalid UTF-16 sequences**. Unpaired surrogate character is not valid in UTF-16 encoding, but SQL Server accept it as valid. Note, in reality, it is unlikely end-user will input invalid UTF-16 sequence since they are not supported in any language or by any IMEs.

- SQL Server **treats a UTF-16 supplementary character as two characters**. The Len function return 2 instead of 1 for such input.
- SQL Server **has a potential risk of breaking a UTF-16 supplementary character into an un-paired surrogate character**, as in calling the substring function. Note, in the real scenario, the chance of this to happen is much lower, because supplementary characters are rare and string function will only break when it happens at the boundary. For example, calling substring(s,5,1) will break if and only if the character at index 5 is a supplementary character.

## Working with UTF-8 source data in SQL Server (Unicode)

SQL Server doesn't support UTF-8 encoding for Unicode data, it supports the UTF-16 encoding. Conversion from UTF-8 to UTF-16 must be addressed through ETL.

With SQL Server Integration Services 2008, a simple Data Conversion component inserted in a package can convert data from UTF-8 to UTF-16.

## Working with Non-Unicode source data in SQL Server (Unicode)

### **Expected inputs**

Code page of the source must be compatible with the code pages supported by ETL tools.

**Illustration: Informatica mechanism**

Code page compatibility matrix

With the "Codepage Relaxation" feature, Src - Informatica Server - Tgt do not need to be in compatible codepages.

However there are some limitations:

- Informatica Server should run in UNICODE mode.
- User has to make sure the code page of Src should be compatible with that of Tgt.
- Both client and server have to set code page relaxation on.

**Illustration: SQL Server Integration Services mechanism**

With SQL Server 2008 Integration Services, a simple Data Conversion component inserted in a package can convert data from a source code page to Unicode. A prerequisite to this step is to make sure that SQL Server 2008 Integration Services is compatible with the source code page.

For a list of compatible code pages, please refer to Appendix 3.

## Working with source data in SQL Server (Non-Unicode)

Code page of the source must be compatible with the code page of SQL Server if you want to store data in Non-Unicode data types.

SQL Server 2008 supports the following code pages.

Code page	Description
<b>1258</b>	Vietnamese
<b>1257</b>	Baltic
<b>1256</b>	Arabic
<b>1255</b>	Hebrew
<b>1254</b>	Turkish
<b>1253</b>	Greek
<b>1252</b>	Latin1 (ANSI)
<b>1251</b>	Cyrillic
<b>1250</b>	Central European
<b>950</b>	Chinese (Traditional)
<b>949</b>	Korean
<b>936</b>	Chinese (Simplified)
<b>932</b>	Japanese
<b>874</b>	Thai
<b>850</b>	Multilingual (MS-DOS Latin1)
<b>437</b>	MS-DOS U.S. English



# Appendix

## Common code page list

Here is the list of the most popular code pages:

Code Page	Name	DisplayName
<b>37</b>	IBM037	IBM EBCDIC (US-Canada)

<b>437</b>	IBM437	OEM United States
<b>500</b>	IBM500	IBM EBCDIC (International)
<b>708</b>	ASMO-708	Arabic (ASMO 708)
<b>720</b>	DOS-720	Arabic (DOS)
<b>737</b>	ibm737	Greek (DOS)
<b>775</b>	ibm775	Baltic (DOS)
<b>850</b>	ibm850	Western European (DOS)
<b>852</b>	ibm852	Central European (DOS)
<b>855</b>	IBM855	OEM Cyrillic
<b>857</b>	ibm857	Turkish (DOS)
<b>858</b>	IBM00858	OEM Multilingual Latin I
<b>860</b>	IBM860	Portuguese (DOS)
<b>861</b>	ibm861	Icelandic (DOS)
<b>862</b>	DOS-862	Hebrew (DOS)
<b>863</b>	IBM863	French Canadian (DOS)
<b>864</b>	IBM864	Arabic (864)
<b>865</b>	IBM865	Nordic (DOS)

<b>866</b>	cp866	Cyrillic (DOS)
<b>869</b>	ibm869	Greek, Modern (DOS)
<b>870</b>	IBM870	IBM EBCDIC (Multilingual Latin-2)
<b>874</b>	windows-874	Thai (Windows)
<b>875</b>	cp875	IBM EBCDIC (Greek Modern)
<b>932</b>	shift_jis	Japanese (Shift-JIS)
<b>936</b>	gb2312	Chinese Simplified (GB2312)
<b>949</b>	ks_c_5601-1987	Korean
<b>950</b>	big5	Chinese Traditional (Big5)
<b>1026</b>	IBM1026	IBM EBCDIC (Turkish Latin-5)
<b>1047</b>	IBM01047	IBM Latin-1
<b>1140</b>	IBM01140	IBM EBCDIC (US-Canada-Euro)
<b>1141</b>	IBM01141	IBM EBCDIC (Germany-Euro)
<b>1142</b>	IBM01142	IBM EBCDIC (Denmark-Norway-Euro)
<b>1143</b>	IBM01143	IBM EBCDIC (Finland-Sweden-Euro)
<b>1144</b>	IBM01144	IBM EBCDIC (Italy-Euro)
<b>1145</b>	IBM01145	IBM EBCDIC (Spain-Euro)

<b>1146</b>	IBM01146	IBM EBCDIC (UK-Euro)
<b>1147</b>	IBM01147	IBM EBCDIC (France-Euro)
<b>1148</b>	IBM01148	IBM EBCDIC (International-Euro)
<b>1149</b>	IBM01149	IBM EBCDIC (Icelandic-Euro)
<b>1200</b>	utf-16	Unicode
<b>1201</b>	unicodeFFFE	Unicode (Big-Endian)
<b>1250</b>	windows-1250	Central European (Windows)
<b>1251</b>	windows-1251	Cyrillic (Windows)
<b>1252</b>	Windows-1252	Western European (Windows)
<b>1253</b>	windows-1253	Greek (Windows)
<b>1254</b>	windows-1254	Turkish (Windows)
<b>1255</b>	windows-1255	Hebrew (Windows)
<b>1256</b>	windows-1256	Arabic (Windows)
<b>1257</b>	windows-1257	Baltic (Windows)
<b>1258</b>	windows-1258	Vietnamese (Windows)
<b>1361</b>	Johab	Korean (Johab)
<b>10000</b>	macintosh	Western European (Mac)

<b>10001</b>	x-mac-japanese	Japanese (Mac)
<b>10002</b>	x-mac-chinesetrad	Chinese Traditional (Mac)
<b>10003</b>	x-mac-korean	Korean (Mac)
<b>10004</b>	x-mac-arabic	Arabic (Mac)
<b>10005</b>	x-mac-hebrew	Hebrew (Mac)
<b>10006</b>	x-mac-greek	Greek (Mac)
<b>10007</b>	x-mac-cyrillic	Cyrillic (Mac)
<b>10008</b>	x-mac-chinesesimp	Chinese Simplified (Mac)
<b>10010</b>	x-mac-romanian	Romanian (Mac)
<b>10017</b>	x-mac-ukrainian	Ukrainian (Mac)
<b>10021</b>	x-mac-thai	Thai (Mac)
<b>10029</b>	x-mac-ce	Central European (Mac)
<b>10079</b>	x-mac-icelandic	Icelandic (Mac)
<b>10081</b>	x-mac-turkish	Turkish (Mac)
<b>10082</b>	x-mac-croatian	Croatian (Mac)
<b>20000</b>	x-Chinese-CNS	Chinese Traditional (CNS)
<b>20001</b>	x-cp20001	TCA Taiwan

<b>20002</b>	x-Chinese-Eten	Chinese Traditional (Eten)
<b>20003</b>	x-cp20003	IBM5550 Taiwan
<b>20004</b>	x-cp20004	TeleText Taiwan
<b>20005</b>	x-cp20005	Wang Taiwan
<b>20105</b>	x-IA5	Western European (IA5)
<b>20106</b>	x-IA5-German	German (IA5)
<b>20107</b>	x-IA5-Swedish	Swedish (IA5)
<b>20108</b>	x-IA5-Norwegian	Norwegian (IA5)
<b>20127</b>	us-ascii	US-ASCII
<b>20261</b>	x-cp20261	T.61
<b>20269</b>	x-cp20269	ISO-6937
<b>20273</b>	IBM273	IBM EBCDIC (Germany)
<b>20277</b>	IBM277	IBM EBCDIC (Denmark-Norway)
<b>20278</b>	IBM278	IBM EBCDIC (Finland-Sweden)
<b>20280</b>	IBM280	IBM EBCDIC (Italy)
<b>20284</b>	IBM284	IBM EBCDIC (Spain)
<b>20285</b>	IBM285	IBM EBCDIC (UK)

<b>20290</b>	IBM290	IBM EBCDIC (Japanese katakana)
<b>20297</b>	IBM297	IBM EBCDIC (France)
<b>20420</b>	IBM420	IBM EBCDIC (Arabic)
<b>20423</b>	IBM423	IBM EBCDIC (Greek)
<b>20424</b>	IBM424	IBM EBCDIC (Hebrew)
<b>20833</b>	x-EBCDIC- KoreanExtended	IBM EBCDIC (Korean Extended)
<b>20838</b>	IBM-Thai	IBM EBCDIC (Thai)
<b>20866</b>	koi8-r	Cyrillic (KOI8-R)
<b>20871</b>	IBM871	IBM EBCDIC (Icelandic)
<b>20880</b>	IBM880	IBM EBCDIC (Cyrillic Russian)
<b>20905</b>	IBM905	IBM EBCDIC (Turkish)
<b>20924</b>	IBM00924	IBM Latin-1
<b>20932</b>	EUC-JP	Japanese (JIS 0208-1990 and 0212-1990)
<b>20936</b>	x-cp20936	Chinese Simplified (GB2312-80)
<b>20949</b>	x-cp20949	Korean Wansung
<b>21025</b>	cp1025	IBM EBCDIC (Cyrillic Serbian-Bulgarian)
<b>21866</b>	koi8-u	Cyrillic (KOI8-U)

<b>28591</b>	iso-8859-1	Western European (ISO)
<b>28592</b>	iso-8859-2	Central European (ISO)
<b>28593</b>	iso-8859-3	Latin 3 (ISO)
<b>28594</b>	iso-8859-4	Baltic (ISO)
<b>28595</b>	iso-8859-5	Cyrillic (ISO)
<b>28596</b>	iso-8859-6	Arabic (ISO)
<b>28597</b>	iso-8859-7	Greek (ISO)
<b>28598</b>	iso-8859-8	Hebrew (ISO-Visual)
<b>28599</b>	iso-8859-9	Turkish (ISO)
<b>28603</b>	iso-8859-13	Estonian (ISO)
<b>28605</b>	iso-8859-15	Latin 9 (ISO)
<b>29001</b>	x-Europa	Europa
<b>38598</b>	iso-8859-8-i	Hebrew (ISO-Logical)
<b>50220</b>	iso-2022-jp	Japanese (JIS)
<b>50221</b>	csISO2022JP	Japanese (JIS-Allow 1 byte Kana)
<b>50222</b>	iso-2022-jp	Japanese (JIS-Allow 1 byte Kana - SO/SI)
<b>50225</b>	iso-2022-kr	Korean (ISO)

<b>50227</b>	x-cp50227	Chinese Simplified (ISO-2022)
<b>51932</b>	euc-jp	Japanese (EUC)
<b>51936</b>	EUC-CN	Chinese Simplified (EUC)
<b>51949</b>	euc-kr	Korean (EUC)
<b>52936</b>	hz-gb-2312	Chinese Simplified (HZ)
<b>54936</b>	GB18030	Chinese Simplified (GB18030)
<b>57002</b>	x-iscii-de	ISCII Devanagari
<b>57003</b>	x-iscii-be	ISCII Bengali
<b>57004</b>	x-iscii-ta	ISCII Tamil
<b>57005</b>	x-iscii-te	ISCII Telugu
<b>57006</b>	x-iscii-as	ISCII Assamese
<b>57007</b>	x-iscii-or	ISCII Oriya
<b>57008</b>	x-iscii-ka	ISCII Kannada
<b>57009</b>	x-iscii-ma	ISCII Malayalam
<b>57010</b>	x-iscii-gu	ISCII Gujarati
<b>57011</b>	x-iscii-pa	ISCII Punjabi
<b>65000</b>	utf-7	Unicode (UTF-7)

<b>65001</b>	utf-8	Unicode (UTF-8)
<b>65005</b>	utf-32	Unicode (UTF-32)
<b>65006</b>	utf-32BE	Unicode (UTF-32 Big-Endian)

## Windows Collations

The following is the list of supported Windows collations in SQL Server 2008.

<b>Windows locale</b>	<b>SQL Server 2008 (100)</b>	<b>SQL Server 2005 (90) or SQL Server 2000</b>
<b>Afrikaans (South Africa)</b>	Latin1_General_100_	Latin1_General_

<b>Albanian (Albania)</b>	Albanian_100_	Albanian_
<b>Alsatian (France)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Amharic (Ethiopia)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Arabic (Algeria)</b>	Arabic_100_	Arabic_
<b>Arabic (Bahrain)</b>	Arabic_100_	Arabic_
<b>Arabic (Egypt)</b>	Arabic_100_	Arabic_
<b>Arabic (Iraq)</b>	Arabic_100_	Arabic_
<b>Arabic (Jordan)</b>	Arabic_100_	Arabic_
<b>Arabic (Kuwait)</b>	Arabic_100_	Arabic_
<b>Arabic (Lebanon)</b>	Arabic_100_	Arabic_
<b>Arabic (Libya)</b>	Arabic_100_	Arabic_
<b>Arabic (Morocco)</b>	Arabic_100_	Arabic_
<b>Arabic (Oman)</b>	Arabic_100_	Arabic_
<b>Arabic (Qatar)</b>	Arabic_100_	Arabic_
<b>Arabic (Saudi Arabia)</b>	Arabic_100_	Arabic_
<b>Arabic (Syria)</b>	Arabic_100_	Arabic_
<b>Arabic (Tunisia)</b>	Arabic_100_	Arabic_

<b>Arabic (U.A.E.)</b>	Arabic_100_	Arabic_
<b>Arabic (Yemen)</b>	Arabic_100_	Arabic_
<b>Armenian (Armenia)</b>	Cyrillic_General_100_	New to SQL Server 2008
<b>Assamese (India)</b>	Assamese_100_ <sup>1</sup>	New to SQL Server 2008
<b>Azeri (Azerbaijan, Cyrillic)</b>	Azeri_Cyrillic_100_	Azeri_Cyrillic_90_
<b>Azeri (Azerbaijan, Latin)</b>	Azeri_Latin_100_	Azeri_Latin_90_
<b>Bashkir (Russia)</b>	Bashkir_100_	New to SQL Server 2008
<b>Basque (Basque)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Belarusian (Belarus)</b>	Cyrillic_General_100_	Cyrillic_General_
<b>Bengali (Bangladesh)</b>	Bengali_100_ <sup>1</sup>	New to SQL Server 2008
<b>Bengali (India)</b>	Bengali_100_ <sup>1</sup>	New to SQL Server 2008
<b>Bosnian (Bosnia and Herzegovina, Cyrillic)</b>	Bosnian_Cyrillic_100_	New to SQL Server 2008
<b>Bosnian (Bosnia and Herzegovina, Latin)</b>	Bosnian_Latin_100_	New to SQL Server 2008
<b>Breton (France)</b>	Breton_100_	New to SQL Server 2008
<b>Bulgarian (Bulgaria)</b>	Cyrillic_General_100_	Cyrillic_General_
<b>Catalan (Catalan)</b>	Latin1_General_100_	Latin1_General_

<b>Chinese (Hong Kong SAR, PRC)</b>	Chinese_Traditional_Stroke_Count_100_	Chinese_Hong_Kong_Stroke_90_
<b>Chinese (Macao SAR)</b>	Chinese_Traditional_Pinyin_100_	New to SQL Server 2008
<b>Chinese (Macau)</b>	Chinese_Traditional_Stroke_Order_100_	New to SQL Server 2008
<b>Chinese (PRC)</b>	Chinese_Simplified_Pinyin_100_	Chinese_PRC_90_, Chinese_PRC_
<b>Chinese (PRC)</b>	Chinese_Simplified_Stroke_Order_100_	Chinese_PRC_Stroke_90_, Chinese_PRC_Stroke_
<b>Chinese (Singapore)</b>	Chinese_Simplified_Pinyin_100_	Chinese_PRC_90_, Chinese_PRC_
<b>Chinese (Singapore)</b>	Chinese_Simplified_Stroke_Order_100_	New to SQL Server 2008
<b>Chinese (Taiwan)</b>	Chinese_Traditional_Bopomofo_100_	Chinese_Taiwan_Bopomofo_90_, Chinese_Taiwan_Bopomofo_
<b>Chinese (Taiwan)</b>	Chinese_Traditional_Stroke_Count_100_	Chinese_Taiwan_Stroke_90_, Chinese_Taiwan_Stroke_
<b>Corsican (France)</b>	Corsican_100_	New to SQL Server 2008
<b>Croatian (Bosnia and Herzegovina, Latin)</b>	Croatian_100_	New to SQL Server 2008
<b>Croatian (Croatia)</b>	Croatian_100_	Croatian_

<b>Czech (Czech Republic)</b>	Czech_100_	Czech_
<b>Danish (Denmark)</b>	Danish_Greenlandic_100_	Danish_Norwegian_
<b>Dari (Afghanistan)</b>	Dari_100_	New to SQL Server 2008
<b>Divehi (Maldives)</b>	Divehi_100_ <sup>1</sup>	Divehi_90_
<b>Dutch (Belgium)</b>	Latin1_General_100_	Latin1_General_
<b>Dutch (Netherlands)</b>	Latin1_General_100_	Latin1_General_
<b>English (Australia)</b>	Latin1_General_100_	Latin1_General_
<b>English (Belize)</b>	Latin1_General_100_	Latin1_General_
<b>English (Canada)</b>	Latin1_General_100_	Latin1_General_
<b>English (Caribbean)</b>	Latin1_General_100_	Latin1_General_
<b>English (India)</b>	Latin1_General_100_	New to SQL Server 2008
<b>English (Ireland)</b>	Latin1_General_100_	Latin1_General_
<b>English (Jamaica)</b>	Latin1_General_100_	Latin1_General_
<b>English (Malaysia)</b>	Latin1_General_100_	New to SQL Server 2008
<b>English (New Zealand)</b>	Latin1_General_100_	Latin1_General_
<b>English (Philippines)</b>	Latin1_General_100_	Latin1_General_
<b>English (Singapore)</b>	Latin1_General_100_	New to SQL Server 2008

<b>English (South Africa)</b>	Latin1_General_100_	Latin1_General_
<b>English (Trinidad and Tobago)</b>	Latin1_General_100_	Latin1_General_
<b>English (United Kingdom)</b>	Latin1_General_100_	Latin1_General_
<b>English (United States)</b>	SQL_Latin1_General_CP1_	SQL_Latin1_General_CP1_
<b>English (Zimbabwe)</b>	Latin1_General_100_	Latin1_General_
<b>Estonian (Estonia)</b>	Estonian_100_	Estonian_
<b>Faroese (Faroe Islands)</b>	Latin1_General_100_	Latin1_General_
<b>Filipino (Philippines)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Finnish (Finland)</b>	Finnish_Swedish_100_	Finnish_Swedish_
<b>French (Belgium)</b>	French_100_	French_
<b>French (Canada)</b>	French_100_	French_
<b>French (France)</b>	French_100_	French_
<b>French (Luxembourg)</b>	French_100_	French_
<b>French (Monaco)</b>	French_100_	French_
<b>French (Switzerland)</b>	French_100_	French_
<b>Frisian (Netherlands)</b>	Frisian_100_	New to SQL Server 2008

<b>Galician (Spain)</b>	Latin1_General_100_	Latin1_General_
<b>Georgian (Georgia)</b>	Georgian_Modern_Sort_100_	Georgian_Modern_Sort_
<b>Georgian (Georgia)</b>	Cyrillic_General_100_	New to SQL Server 2008
<b>German - Phone Book Sort (DIN)</b>	German_PhoneBook_100_	German_PhoneBook_
<b>German (Austria)</b>	Latin1_General_100_	Latin1_General_
<b>German (Germany)</b>	Latin1_General_100_	Latin1_General_
<b>German (Liechtenstein)</b>	Latin1_General_100_	Latin1_General_
<b>German (Luxembourg)</b>	Latin1_General_100_	Latin1_General_
<b>German (Switzerland)</b>	Latin1_General_100_	Latin1_General_
<b>Greek (Greece)</b>	Greek_100_	Greek_
<b>Greenlandic (Greenland)</b>	Danish_Greenlandic_100_	New to SQL Server 2008
<b>Gujarati (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Hausa (Nigeria, Latin)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Hebrew (Israel)</b>	Hebrew_100_	Hebrew_
<b>Hindi (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Hungarian (Hungary)</b>	Hungarian_100_	Hungarian_

<b>Hungarian Technical Sort</b>	Hungarian_Technical_100_	Hungarian_Technical_
<b>Icelandic (Iceland)</b>	Icelandic_100_	Icelandic_
<b>Igbo (Nigeria)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Indonesian (Indonesia)</b>	Latin1_General_100_	Latin1_General_
<b>Inuktitut (Canada, Latin)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Inuktitut (Syllabics) Canada</b>	Latin1_General_100_	New to SQL Server 2008
<b>Irish (Ireland)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Italian (Italy)</b>	Latin1_General_100_	Latin1_General_
<b>Italian (Switzerland)</b>	Latin1_General_100_	Latin1_General_
<b>Japanese (Japan XJIS)</b>	Japanese_XJIS_100_	Japanese_90_, Japanese_
<b>Japanese (Japan)</b>	Japanese_Bushu_Kakusu_100_	New to SQL Server 2008
<b>Kannada (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Kazakh (Kazakhstan)</b>	Kazakh_100_	Kazakh_90_
<b>Khmer (Cambodia)</b>	Khmer_100_ <sup>1</sup>	New to SQL Server 2008
<b>K'iche (Guatemala)</b>	Modern_Spanish_100_	New to SQL Server 2008
<b>Kinyarwanda (Rwanda)</b>	Latin1_General_100_	New to SQL Server 2008

<b>Konkani (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Korean (Korea Dictionary Sort)</b>	Korean_100_	Korean_90_ Korean_Wansung_
<b>Kyrgyz (Kyrgyzstan)</b>	Cyrillic_General_100_	Cyrillic_General_
<b>Lao (Lao PDR)</b>	Lao_100_ <sup>1</sup>	New to SQL Server 2008
<b>Latvian (Latvia)</b>	Latvian_100_	Latvian_
<b>Lithuanian (Lithuania)</b>	Lithuanian_100_	Lithuanian_
<b>Lower Sorbian (Germany)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Luxembourgish (Luxembourg)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Macedonian (Macedonia, FYROM)</b>	Macedonian_FYROM_100_	Macedonian_FYROM_90_
<b>Malay (Brunei Darussalam)</b>	Latin1_General_100_	Latin1_General_
<b>Malay (Malaysia)</b>	Latin1_General_100_	Latin1_General_
<b>Malayalam (India)</b>	Indic_General_100_ <sup>1</sup>	New to SQL Server 2008
<b>Maltese (Malta)</b>	Maltese_100_	New to SQL Server 2008
<b>Maori (New Zealand)</b>	Maori_100_	New to SQL Server 2008

<b>Mapudungun (Chile)</b>	Mapudungan_100_	New to SQL Server 2008
<b>Marathi (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Mohawk (Canada)</b>	Mohawk_100_	New to SQL Server 2008
<b>Mongolian (Mongolia)</b>	Cyrillic_General_100_	Cyrillic_General_
<b>Mongolian (PRC)</b>	Cyrillic_General_100_	New to SQL Server 2008
<b>Nepali (Nepal)</b>	Nepali_100_ <sup>1</sup>	New to SQL Server 2008
<b>Norwegian (Bokmål, Norway)</b>	Norwegian_100_	New to SQL Server 2008
<b>Norwegian (Nynorsk, Norway)</b>	Norwegian_100_	New to SQL Server 2008
<b>Occitan (France)</b>	French_100_	New to SQL Server 2008
<b>Oriya (India)</b>	Indic_General_100_ <sup>1</sup>	New to SQL Server 2008
<b>Pashto (Afghanistan)</b>	Pashto_100_ <sup>1</sup>	New to SQL Server 2008
<b>Persian (Iran)</b>	Persian_100_	New to SQL Server 2008
<b>Polish (Poland)</b>	Polish_100_	Polish_
<b>Portuguese (Brazil)</b>	Latin1_General_100_	Latin1_General_
<b>Portuguese (Portugal)</b>	Latin1_General_100_	Latin1_General_
<b>Punjabi (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_

<b>Quechua (Bolivia)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Quechua (Ecuador)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Quechua (Peru)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Romanian (Romania)</b>	Romanian_100_	Romanian_
<b>Romansh (Switzerland)</b>	Romansh_100_	New to SQL Server 2008
<b>Russian (Russia)</b>	Cyrillic_General_100_	Cyrillic_General_
<b>Sami (Inari, Finland)</b>	Sami_Sweden_Finland_100_	New to SQL Server 2008
<b>Sami (Lule,Norway)</b>	Sami_Norway_100_	New to SQL Server 2008
<b>Sami (Lule, Sweden)</b>	Sami_Sweden_Finland_100_	New to SQL Server 2008
<b>Sami (Northern, Finland)</b>	Sami_Sweden_Finland_100_	New to SQL Server 2008
<b>Sami (Northern,Norway)</b>	Sami_Norway_100_	New to SQL Server 2008
<b>Sami (Northern, Sweden)</b>	Sami_Sweden_Finland_100_	New to SQL Server 2008
<b>Sami (Skolt, Finland)</b>	Sami_Sweden_Finland_100_	New to SQL Server 2008
<b>Sami (Southern, Norway)</b>	Sami_Norway_100_	New to SQL Server 2008
<b>Sami (Southern, Sweden)</b>	Sami_Sweden_Finland_100_	New to SQL Server 2008
<b>Sanskrit (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Serbian (Bosnia and</b>	Serbian_Cyrillic_100_	New to SQL Server 2008

<b>Herzegovina, Cyrillic)</b>		
<b>Serbian (Bosnia and Herzegovina, Latin)</b>	Serbian_Latin_100_	New to SQL Server 2008
<b>Serbian (Serbia, Cyrillic)</b>	Serbian_Cyrillic_100_	New to SQL Server 2008
<b>Serbian (Serbia, Latin)</b>	Serbian_Latin_100_	New to SQL Server 2008
<b>Sesotho sa Leboa/Northern Sotho (South Africa)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Setswana/Tswana (South Africa)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Sinhala (Sri Lanka)</b>	Indic_General_100_ <sup>1</sup>	New to SQL Server 2008
<b>Slovak (Slovakia)</b>	Slovak_100_	Slovak_
<b>Slovenian (Slovenia)</b>	Slovenian_100_	Slovenian_
<b>Spanish (Argentina)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Bolivia)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Chile)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Colombia)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Costa Rica)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Dominican Republic)</b>	Modern_Spanish_100_	Modern_Spanish_

<b>Spanish (Ecuador)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (El Salvador)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Guatemala)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Honduras)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Mexico)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Nicaragua)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Panama)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Paraguay)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Peru)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Puerto Rico)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Spain)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Spain, Traditional Sort)</b>	Traditional_Spanish_100_	Traditional_Spanish_
<b>Spanish (United States)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Spanish (Uruguay)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Spanish (Venezuela)</b>	Modern_Spanish_100_	Modern_Spanish_
<b>Swahili (Kenya)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Swedish (Finland)</b>	Finnish_Swedish_100_	Finnish_Swedish_

<b>Swedish (Sweden)</b>	Finnish_Swedish_100_	Finnish_Swedish_
<b>Syriac (Syria)</b>	Syriac_100_ <sup>1</sup>	Syriac_90_
<b>Tajik (Tajikistan)</b>	Cyrillic_General_100_	New to SQL Server 2008
<b>Tamazight (Algeria, Latin)</b>	Tamazight_100_	New to SQL Server 2008
<b>Tamil (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Tatar (Russia)</b>	Tatar_100_	Tatar_90_, Cyrillic_General_
<b>Telugu (India)</b>	Indic_General_100_ <sup>1</sup>	Indic_General_90_
<b>Thai (Thailand)</b>	Thai_100_	Thai_
<b>Tibetan (PRC)</b>	Tibetan_100_ <sup>1</sup>	New to SQL Server 2008
<b>Turkish (Turkey)</b>	Turkish_100_	Turkish_
<b>Turkmen (Turkmenistan)</b>	Turkmen_100_	New to SQL Server 2008
<b>Uighur (PRC)</b>	Uighur_100_	New to SQL Server 2008
<b>Ukrainian (Ukraine)</b>	Ukrainian_100_	Ukrainian_
<b>Upper Sorbian (Germany)</b>	Upper_Sorbian_100_	New to SQL Server 2008
<b>Urdu (Pakistan)</b>	Urdu_100_	New to SQL Server 2008
<b>Uzbek (Uzbekistan, Cyrillic)</b>	Cyrillic_General_100_	Cyrillic_General_

<b>Uzbek (Uzbekistan, Latin)</b>	Uzbek_Latin_100_	Uzbek_Latin_90_
<b>Vietnamese (Vietnam)</b>	Vietnamese_100_	Vietnamese_
<b>Welsh (United Kingdom)</b>	Welsh_100_	New to SQL Server 2008
<b>Wolof (Senegal)</b>	French_100_	New to SQL Server 2008
<b>Xhosa/isiXhosa (South Africa)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Yakut (Russia)</b>	Yakut_100_	New to SQL Server 2008
<b>Yi (PRC)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Yoruba (Nigeria)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Zulu/isiZulu (South Africa)</b>	Latin1_General_100_	New to SQL Server 2008
<b>Deprecated, not available at server level in SQL Server 2008</b>	Hindi	Hindi
<b>Deprecated, not available at server level in SQL Server 2008</b>	Korean_Wansung_Unicode	Korean_Wansung_Unicode
<b>Deprecated, not available at server level in SQL Server 2008</b>	Lithuanian_Classic	Lithuanian_Classic

**Deprecated, not  
available at server level  
in SQL Server 2008**

Macedonian

Macedonian

## SQL Server 2008 Integration Services locales and code pages

### Supported locales for flat files

Afrikaans	Azeri (Cyrillic, Azerbaijan)	Czech
Afrikaans (South Africa)	Azeri (Latin, Azerbaijan)	Czech (Czech Republic)
Albanian	Bashkir (Russia)	Danish
Albanian (Albania)	Basque	Danish (Denmark)
Alsatian (France)	Basque (Basque)	Dari (Afghanistan)
Amharic (Ethiopia)	Belarusian	Divehi
Arabic	Belarusian (Belarus)	Divehi (Maldives)
Arabic (Algeria)	Bengali (Bangladesh)	Dutch
Arabic (Bahrain)	Bengali (India)	Dutch (Belgium)
Arabic (Egypt)	Bosnian (Cyrillic, Bosnia and Herzegovina)	Dutch (Netherlands)
Arabic (Iraq)	Bosnian (Latin, Bosnia and Herzegovina)	English
Arabic (Jordan)	Breton (France)	English (Australia)
Arabic (Kuwait)	Bulgarian	English (Belize)
Arabic (Lebanon)	Bulgarian (Bulgaria)	English (Canada)
Arabic (Libya)	Catalan	English (Caribbean)
Arabic (Morocco)	Catalan (Catalan)	English (India)
Arabic (Oman)	Chinese (Hong Kong S.A.R.)	English (Ireland)
Arabic (Qatar)	Chinese (Macao S.A.R.)	English (Jamaica)
Arabic (Saudi Arabia)	Chinese (People's Republic of China)	English (Malaysia)
Arabic (Syria)	Chinese (Simplified)	English (New Zealand)
Arabic (Tunisia)	Chinese (Singapore)	English (Republic of the Philippines)
Arabic (U.A.E.)	Chinese (Taiwan)	English (Singapore)
Arabic (Yemen)	Chinese (Traditional)	English (South Africa)
Armenian	Corsican (France)	English (Trinidad and Tobago)
Armenian (Armenia)	Croatian	English (United Kingdom)
Assamese (India)	Croatian (Croatia)	English (United States)
Azeri	Croatian (Latin, Bosnia and Herzegovina)	English (Zimbabwe)

Estonian	Greenlandic (Greenland)
Estonian (Estonia)	Gujarati
Faroese	Gujarati (India)
Faroese (Faroe Islands)	Hausa (Latin, Nigeria)
Filipino (Philippines)	Hebrew
Finnish	Hebrew (Israel)
Finnish (Finland)	Hindi
French	Hindi (India)
French (Belgium)	Hungarian
French (Canada)	Hungarian (Hungary)
French (France)	Icelandic
French (Luxembourg)	Icelandic (Iceland)
French (Principality of Monaco)	Igbo (Nigeria)
French (Switzerland)	Indonesian
Frisian (Netherlands)	Indonesian (Indonesia)
Galician	Inuktitut (Latin, Canada)
Galician (Galician)	Inuktitut (Syllabics, Canada)
Georgian	Invariant Language (Invariant Country)
Georgian (Georgia)	Irish (Ireland)
German	isiXhosa (South Africa)
German (Austria)	isiZulu (South Africa)
German (Germany)	Italian
German (Liechtenstein)	Italian (Italy)
German (Luxembourg)	Italian (Switzerland)
German (Switzerland)	Japanese
Greek	Japanese (Japan)
Greek (Greece)	Kannada

Kannada (India)	Maltese (Malta)
Kazakh	Maori (New Zealand)
Kazakh (Kazakhstan)	Mapudungun (Chile)
Khmer (Cambodia)	Marathi
K'iche (Guatemala)	Marathi (India)
Kinyarwanda (Rwanda)	Mohawk (Mohawk)
Kiswahili	Mongolian
Kiswahili (Kenya)	Mongolian (Cyrillic, Mongolia)
Konkani	Mongolian (Traditional Mongolian, PRC)
Konkani (India)	Nepali (Nepal)
Korean	Norwegian
Korean (Korea)	Norwegian, Bokmål (Norway)
Kyrgyz	Norwegian, Nynorsk (Norway)
Kyrgyz (Kyrgyzstan)	Occitan (France)
Lao (Lao P.D.R.)	Oriya (India)
Latvian	Pashto (Afghanistan)
Latvian (Latvia)	Persian
Lithuanian	Persian (Iran)
Lithuanian (Lithuania)	Polish
Lower Sorbian (Germany)	Polish (Poland)
Luxembourgish (Luxembourg)	Portuguese
Macedonian	Portuguese (Brazil)
Macedonian (Former Yugoslav Republic of Macedonia)	Portuguese (Portugal)
Malay	Punjabi
Malay (Brunei Darussalam)	Punjabi (India)
Malay (Malaysia)	Quechua (Bolivia)
Malayalam (India)	Quechua (Ecuador)

Quechua (Peru)	Sesotho sa Leboa (South Africa)
Romanian	Setswana (South Africa)
Romanian (Romania)	Sinhala (Sri Lanka)
Romansh (Switzerland)	Slovak
Russian	Slovak (Slovakia)
Russian (Russia)	Slovenian
Sami, Inari (Finland)	Slovenian (Slovenia)
Sami, Lule (Norway)	Spanish
Sami, Lule (Sweden)	Spanish (Argentina)
Sami, Northern (Finland)	Spanish (Bolivia)
Sami, Northern (Norway)	Spanish (Chile)
Sami, Northern (Sweden)	Spanish (Colombia)
Sami, Skolt (Finland)	Spanish (Costa Rica)
Sami, Southern (Norway)	Spanish (Dominican Republic)
Sami, Southern (Sweden)	Spanish (Ecuador)
Sanskrit	Spanish (El Salvador)
Sanskrit (India)	Spanish (Guatemala)
Scottish Gaelic (United Kingdom)	Spanish (Honduras)
Serbian	Spanish (Mexico)
Serbian (Cyrillic, Bosnia and Herzegovina)	Spanish (Nicaragua)
Serbian (Cyrillic, Montenegro)	Spanish (Panama)
Serbian (Cyrillic, Serbia and Montenegro (Former))	Spanish (Paraguay)
Serbian (Cyrillic, Serbia)	Spanish (Peru)
Serbian (Latin, Bosnia and Herzegovina)	Spanish (Puerto Rico)
Serbian (Latin, Montenegro)	Spanish (Spain)
Serbian (Latin, Serbia and Montenegro (Former))	Spanish (United States)
Serbian (Latin, Serbia)	Spanish (Uruguay)

Spanish (Venezuela)  
Swedish  
Swedish (Finland)  
Swedish (Sweden)  
Syriac  
Syriac (Syria)  
Tajik (Cyrillic, Tajikistan)  
Tamazight (Latin, Algeria)  
Tamil  
Tamil (India)  
Tatar  
Tatar (Russia)  
Telugu  
Telugu (India)  
Thai  
Thai (Thailand)  
Tibetan (PRC)  
Turkish  
Turkish (Turkey)  
Turkmen (Turkmenistan)  
Ukrainian  
Ukrainian (Ukraine)  
Upper Sorbian (Germany)  
Urdu  
Urdu (Islamic Republic of Pakistan)  
Uyghur (PRC)  
Uzbek

Uzbek (Cyrillic, Uzbekistan)  
Uzbek (Latin, Uzbekistan)  
Vietnamese  
Vietnamese (Vietnam)  
Welsh (United Kingdom)  
Wolof (Senegal)  
Yakut (Russia)  
Yi (PRC)  
Yoruba (Nigeria)

## Supported code pages for flat files

10000 (MAC - Roman)	1253 (ANSI - Greek)
10001 (MAC - Japanese)	1254 (ANSI - Turkish)
10002 (MAC - Traditional Chinese Big5)	1255 (ANSI - Hebrew)
10003 (MAC - Korean)	1256 (ANSI - Arabic)
10004 (MAC - Arabic)	1257 (ANSI - Baltic)
10005 (MAC - Hebrew)	1258 (ANSI/OEM - Viet Nam)
10006 (MAC - Greek I)	1361 (Korean - Johab)
10007 (MAC - Cyrillic)	20000 (CNS - Taiwan)
10008 (MAC - Simplified Chinese GB 2312)	20001 (TCA - Taiwan)
10010 (MAC - Romania)	20002 (Eten - Taiwan)
10017 (MAC - Ukraine)	20003 (IBM5550 - Taiwan)
10021 (MAC - Thai)	20004 (TeleText - Taiwan)
10029 (MAC - Latin II)	20005 (Wang - Taiwan)
10079 (MAC - Icelandic)	20105 (IA5 IRV International Alphabet No.5)
10081 (MAC - Turkish)	20106 (IA5 German)
10082 (MAC - Croatia)	20107 (IA5 Swedish)
1026 (IBM EBCDIC - Turkish (Latin-5))	20108 (IA5 Norwegian)
1047 (IBM EBCDIC - Latin-1/Open System)	20127 (US-ASCII)
1140 (IBM EBCDIC - U.S./Canada (37 + Euro))	20261 (T.61)
1141 (IBM EBCDIC - Germany (20273 + Euro))	20269 (ISO 6937 Non-Spacing Accent)
1142 (IBM EBCDIC - Denmark/Norway (20277 + Euro))	20273 (IBM EBCDIC - Germany)
1143 (IBM EBCDIC - Finland/Sweden (20278 + Euro))	20277 (IBM EBCDIC - Denmark/Norway)
1144 (IBM EBCDIC - Italy (20280 + Euro))	20278 (IBM EBCDIC - Finland/Sweden)
1145 (IBM EBCDIC - Latin America/Spain (20284 + Euro))	20280 (IBM EBCDIC - Italy)
1146 (IBM EBCDIC - United Kingdom (20285 + Euro))	20284 (IBM EBCDIC - Latin America/Spain)
1148 (IBM EBCDIC - International (500 + Euro))	20285 (IBM EBCDIC - United Kingdom)
1149 (IBM EBCDIC - Icelandic (20871 + Euro))	20290 (IBM EBCDIC - Japanese Katakana Extended)
1250 (ANSI - Central Europe)	20297 (IBM EBCDIC - France)
1251 (ANSI - Cyrillic)	20420 (IBM EBCDIC - Arabic)
1252 (ANSI - Latin I)	20423 (IBM EBCDIC - Greek)

20424 (IBM EBCDIC - Hebrew)	50222 (ISO-2022 Japanese JIS X 0201-1989)
20833 (IBM EBCDIC - Korean Extended)	50225 (ISO-2022 Korean)
20838 (IBM EBCDIC - Thai)	50227 (ISO-2022 Simplified Chinese)
20866 (Russian - KOI8)	50229 (ISO-2022 Traditional Chinese)
20871 (IBM EBCDIC - Icelandic)	51949 (EUC-Korean)
20880 (IBM EBCDIC - Cyrillic (Russian))	52936 (HZ-GB2312 Simplified Chinese)
20905 (IBM EBCDIC - Turkish)	54936 (GB18030 Simplified Chinese)
20924 (IBM EBCDIC - Latin-1/Open System (1047 + Euro))	57002 (ISCII - Devanagari)
20932 (JIS X 0208-1990 & 0212-1990)	57003 (ISCII - Bengali)
20936 (Simplified Chinese GB2312)	57004 (ISCII - Tamil)
21025 (IBM EBCDIC - Cyrillic (Serbian, Bulgarian))	57005 (ISCII - Telugu)
21027 (Ext Alpha Lowercase)	57006 (ISCII - Assamese)
21866 (Ukrainian - KOI8-U)	57007 (ISCII - Oriya)
28591 (ISO 8859-1 Latin I)	57008 (ISCII - Kannada)
28592 (ISO 8859-2 Central Europe)	57009 (ISCII - Malayalam)
28593 (ISO 8859-3 Latin 3)	57010 (ISCII - Gujarati)
28594 (ISO 8859-4 Baltic)	57011 (ISCII - Punjabi (Gurmukhi))
28595 (ISO 8859-5 Cyrillic)	65000 (UTF-7)
28596 (ISO 8859-6 Arabic)	65001 (UTF-8)
28597 (ISO 8859-7 Greek)	708 (Arabic - ASMO)
28598 (ISO 8859-8 Hebrew: Visual Ordering)	720 (Arabic - Transparent ASMO)
28599 (ISO 8859-9 Latin 5)	737 (OEM - Greek 437G)
28603 (ISO 8859-13 Latin 7)	775 (OEM - Baltic)
28605 (ISO 8859-15 Latin 9)	850 (OEM - Multilingual Latin I)
37 (IBM EBCDIC - U.S./Canada)	852 (OEM - Latin II)
38598 (ISO 8859-8 Hebrew: Logical Ordering)	855 (OEM - Cyrillic)
437 (OEM - United States)	857 (OEM - Turkish)
500 (IBM EBCDIC - International)	858 (OEM - Multilingual Latin I + Euro)
50220 (ISO-2022 Japanese with no halfwidth Katakana)	860 (OEM - Portuguese)
50221 (ISO-2022 Japanese with halfwidth Katakana)	861 (OEM - Icelandic)

862 (OEM - Hebrew)  
863 (OEM - Canadian French)  
864 (OEM - Arabic)  
865 (OEM - Nordic)  
866 (OEM - Russian)  
869 (OEM - Modern Greek)  
870 (IBM EBCDIC - Multilingual/ROECE (Latin-2))  
874 (ANSI/OEM - Thai)  
875 (IBM EBCDIC - Modern Greek)  
932 (ANSI/OEM - Japanese Shift-JIS)  
936 (ANSI/OEM - Simplified Chinese GBK)  
949 (ANSI/OEM - Korean)  
950 (ANSI/OEM - Traditional Chinese Big5)

## Windows/SAP code pages mappings

<b>2-char Lang ID</b>	<b>1-char Lang ID</b>	<b>Language</b>	<b>Win32 code page</b>	<b>SAP code page</b>	<b>ISO character set</b>
ar	a	Arabic	1256	1100	Windows-1256

<b>he</b>	b	Hebrew	1255	1800	Windows-1255
<b>cs</b>	c	Czech	1250	1404	Windows-1250
<b>de</b>	d	German	1252	1100	us-ascii
<b>en</b>	e	English	1252	1100	us-ascii
<b>fr</b>	f	French	1252	1100	us-ascii
<b>el</b>	g	Greek	1253	1700	Windows-1253
<b>hu</b>	h	Hungarian	1250	1401	Windows-1250
<b>it</b>	i	Italian	1252	1100	us-ascii
<b>ja</b>	j	Japanese	932	8000	x-sjis
<b>da</b>	k	Danish	1252	1100	Windows-1250
<b>pl</b>	l	Polish	1250	1401	Windows-1250
<b>zf</b>	m	Traditional Chinese	950	8300	big5
<b>nl</b>	n	Dutch	1252	1100	us-ascii
<b>no</b>	o	Norwegian	1252	1100	us-ascii
<b>pt</b>	p	Portuguese	1252	1100	us-ascii
<b>sk</b>	q	Slovakian	1250	1404	Windows-1250
<b>ru</b>	r	Russian	1251	1500	Windows-1251
<b>es</b>	s	Spanish	1252	1100	us-ascii
<b>tr</b>	t	Turkish	1254	1610	Windows-1254
<b>fi</b>	u	Finnish	1252	1100	us-ascii
<b>sv</b>	v	Swedish	1252	1100	us-ascii
<b>bg</b>	w	Bulgarian	1251	1500	Windows-1251
<b>zh</b>	1	Simplified Chinese	936	8400	GB_2312_80
<b>th</b>	2	Thai	874	8600	x-TIS-620
<b>ko</b>	3	Korean	949	8500	ks_c_5601-1987
<b>ro</b>	4	Romanian	1250	1401	Windows-1250
<b>sl</b>	5	Slovenian	1250	1401	Windows-1250
<b>hr</b>	6	Croatian	1250	1401	Windows-1250
<b>zz</b>	z	Language-independent	1252	1100	Z