



# Character Sets and Unicode in Firebird

**Stefan Heymann**

[www.destructor.de](http://www.destructor.de)  
[stefan@destructor.de](mailto:stefan@destructor.de)

This session will show you how non-ASCII text and text from different languages and scripts can be handled with Firebird. It will introduce you to the Firebird world of Character Sets, Collations, Unicode, case insensitive searching and Client character sets.

We will take a look at how international text influences the design of your database and the client programming. The session will also discuss the various new international features of Firebird 2.0.

---



# Topics

- Characters
- Character Sets
- Unicode
- Firebird
- Examples



# Characters



## Glyphs vs. Characters

A

A

A

A

A

A

A

Latin uppercase A

---



# Glyph, Character, Character Set

- A Glyph is something you can see with your eyes
  - A Character is an abstract concept
  - Rendering of characters as Glyphs is the job of the rendering machine (Postscript, GDI, TrueType, Web Browser, etc.)
  - We mostly care for processing the characters
  - A Character Set gives all characters a number:
    - Uppercase A = 65
    - Uppercase B = 66
    - etc.
-



## Glyphs

- Not all languages display glyphs as a string of left-to-right, contiguous rectangles
- Right-to-left (Arabic, Hebrew), top-to-bottom (Japanese, Chinese)
- Several characters can „melt“ into one glyph

Keystrokes	ل	ا	لا		غ	غ
Input characters	ل	ا	ل	ا	غ	غ
Encoded characters (byte values in hex)	0644	0627	0644	0627	0639	0639
Display	لااغغ					



# Character Sets



## ASCII: The Mother of Character Sets

- American Standard Code for Information Interchange: ASCII, ISO 646
  - 7 bits, characters ranging from 0 to 127 (00..7F)
  - 32 invisible control characters (NUL, TAB, CR, LF, FF, BEL, ESC, ...)
  - A..Z, a..z, Digits 0..9, Punctuation (;.-?)
  - Optimized for English
  - Only Latin characters, no accents, no umlauts
  - MIME code: US-ASCII
-





## ASCII for Europe

- Language specific characters get a new assignment
- [=Ä \=Ö ]=Ü
- Problem: Printer and screen must have same setting
- Impossible to mix French and German in one text:

„Amélie knackt gerne die Kruste von Crème Brulé  
mit dem Löffel“

- Died together with DOS
-



## Use the 8th Bit

- 128 new characters: 0 to 255 (00..FF)
  - A lot of 8-Bit character sets have evolved
    - ISO 8859-x = ASCII + 160..255
    - ISO 8859-1 = Latin-1  
(Western European languages)
    - Windows 1252 = ISO 8859-1 + 128..159
    - etc.
-



## ISO 8859

- Characters 0..127 identical to ASCII
- 128..159 undefined control characters
- 160..255 individually assigned



## ISO 8859-1 / Latin-1

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	í	φ	£	℥	¥	!	§		©	≡	«	¬	–	®	—
B0	B1	B2	B3	B4	B5	B6	B7	B8	B9	BA	BB	BC	BD	BE	BF
°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA	DB	DC	DD	DE	DF
Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF
ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

Afrikaans, Albanese, Basque, Danish, German, **English**, Faroese, Finnish, French, Icelandic, Italian, Catalan, Dutch, Norwegian, Portuguese, Rhaeto-Romance, Scottish Gaelic, Swedish, Spanish, Suaheli  
*Large parts of the world – wide-spread use*



## ISO 8859-2 / Latin-2

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	À		Ł	Ǻ	Ĺ	Š	Ś		Ŝ	Ş	Ť	Ž	-	Ž	Ž
B0	B1	B2	B3	B4	B5	B6	B7	B8	B9	BA	BB	BC	BD	BE	BF
	á		ł		ĺ	š			ŝ	ş	ť	ž		ž	ž
C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
	Ā	Ā	Ā	Ā	Ĺ	Č	Ç	Č	Ě	Ě	Ě	Ě	Ī	Ī	Ď
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA	DB	DC	DD	DE	DF
	Đ	Ñ	Ñ	Ō	Ō	Ō	×	Ř	Ů	Ú	Ů	Ü	Ý	Ť	ß
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
	ř	ā	ā	ā	ā	ī	č	ç	č	ě	ę	ě	ě	î	đ
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF
	đ	ñ	ñ	ō	ō	ō	÷	ř	ů	ú	ű	ü	ý	ţ	

Central and Eastern Europe (Czech, Polish, etc.)



## ISO 8859-4 / Latin-4

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	À	Ā	Ȧ	Ȧ	İ	Ł	Ś		Š	Ē	Ġ	Ʀ	–	Ž	–
B0	°	ā	˙	ı	ĩ	ł	˘	˘	š	ē	ġ	ʀ	ı	ž	ı
C0	Ā	Ā	Ā	Ā	Ā	Æ	ı	č	ě	ĕ	ĕ	ĕ	ĩ	î	ï
D0	Đ	Ń	Ō	Ȧ	Ō	Ö	×	Ø	Ů	Ú	Û	Ü	Ũ	Ū	ß
E0	ā	ā	ā	ā	ā	æ	ı	č	ě	ĕ	ĕ	ĕ	ĩ	î	ï
F0	đ	ń	ō	Ȧ	ō	ö	÷	ø	ů	ú	û	ü	ũ	ū	.

Northern Europe, Baltic, Greenlanic, Sami



## ISO 8859-5 / Cyrillic

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	Ё	Ђ	Ѓ	Є	Ѕ	І	Ї	Ј	Љ	Њ	Ћ	Ќ	–	Ў	Ч
B0	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О
C0	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю
D0	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о
E0	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю
F0	ѐ	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	ѕ	ў	ч

Cyrillic (Russia, Ukraine, etc.)

More important: KOI8-R (Russian), KOI8-U (Ukrainian)



## ISO 8859-9 / Latin-5

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	ı	¢	£	¤	¥	¦	§		©	ª	«	¬	­	®	¯
B0	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾
B1	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
C0	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î
C1	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î
D0	Ġ	Ñ	Ò	Ó	Ô	Õ	×	Ø	Ù	Ú	Û	Ü	Ý	Ş	ß
D1	Ñ	Ò	Ó	Ô	Õ	×	Ø	Ù	Ú	Û	Ü	Ý	Ş	ß	
E0	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î
E1	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î
F0	ğ	ñ	ò	ó	ô	õ	÷	ø	ù	ú	û	ü	ı	ş	ý
F1	ñ	ò	ó	ô	õ	÷	ø	ù	ú	û	ü	ı	ş	ý	

Turkish





## Windows Character Sets

- Congruent to ISO-8859
  - Additional assignment of characters **128..159** with visible (non-control) characters
    - Hyphens n length – and m length — (vs. Dash -)
    - Typographic „quotation marks“, etc.
  - Windows character sets officially registered at IANA
-



## windows-1252

- Congruent to ISO 8859-1 (Western European languages, including English)
- Additional characters in the 128..159 range:

€ , *f* „ ... † ‡ ^ %‰ Š < Œ Ž  
` / \" \" • – — ~ ™ š > œ ž Ÿ

Euro sign € since 2000 (Windows 1252-2000)

---



## Multi-Byte Character Sets MBCS

- Multiple Bytes per Character
  - Eastern Asian Languages (CJK)
  - String Length  $\neq$  Length of character chain
  - Making extraction of sub-strings more difficult
  - Firebird functions:
    - **BIT\_LENGTH**: length of a string in bits (!)
    - **CHAR\_LENGTH/CHARACTER\_LENGTH**: length of a string in characters
    - **OCTET\_LENGTH**: length of a string in bytes
-



# Unicode



## Why Unicode?

- One single Character Set for all languages/scripts
- No code overlaps
- Hardware and OS independant
- Standardisation ISO **10646** (vs. ASCII = ISO 646)



# Unicode

- Started with 16 Bits/Character, now 32 Bits/Char
- Ability to code 1,114,112 characters
- Currently only a fraction is used
- Current version: 5.0.0
- Defines Characters, *not* Glyphs
- Practically equal to ISO/IEC 10646



## Character Definition

- Unicode defines a numerical *Code Point* (scalar value) and an *Identifier* for every character

0041 LATIN CAPITAL LETTER A

00E4 LATIN SMALL LETTER A WITH DIAERESIS

0391 GREEK CAPITAL LETTER ALPHA

05D0 HEBREW LETTER ALEF

0950 DEVANAGARI OM

1D56C MATHEMATICAL BOLD FRAKTUR CAPITAL A



## Unicode Code Points

- Codespace: 0..10FFFF
  - Usual Notation is hexadecimal with preceding U+ and at least 4 digits
  - U+0020
  - U+0041
  - U+1D56C
-





## Unicode Character Names

- Consisting of the uppercase letters A..Z, digits 0..9, hyphen (-) and whitespace.
  - BYZANTINE MUSICAL SYMBOL LEIMMA ENOS  
CHRONOU
  - DESERET CAPITAL LETTER OW
  - BRAILLE PATTERN DOTS-1245
-



## Unicode Coding

- Storage of Code Points in memory
  - 32 Bits/Character easy but too fat
  - There are various codings around:
    - 8-Bit (UTF-8)
    - 16-Bit (UCS-2, UTF-16)
    - 32-Bit (UCS-4, UTF-32)
    - PunyCode for international Domain names
    - UTF-7, UTF-1, etc.
-



## UCS-2

- 16 Bits per Character
  - Code Area: 0000..FFFF  
= Basic Multilingual Plane (BMP)
  - Since Unicode 3.1 not all characters can be coded
  - Replaced by UTF-16
  - „Unicode“ often used as a synonym for UCS-2  
(which is wrong and can lead to misunderstanding)
-



## UTF-16

- 16 Bits per Character
- All characters  $> \text{FFFF}$  must be coded as a „Surrogate Pair“ and occupy 2 contiguous 16-Bit words
- Complete Code Space can be coded
- Difficult to calculate string length or substrings



## Endianness

- Problem with UCS-2, UTF-16: Low/High-Byte ordering
  - Differentiate in UTF-16BE and UTF-16LE in metadata
  - Byte Order Mark     BOM     U+FEFF
  - U+FEFF set at the very beginning of each text
  - U+FFFE is (and will be) an invalid code point
-



## UCS-2 vs. UTF-16

- UTF-16 ist backwards compatible
- UCS-2 often synonym for „Unicode“
- WideString, wchar\_t
- NT3, NT4: UCS-2
- Since Windows 2000: UTF-16



## UTF-8

- Coding as 8-Bit strings
  - US-ASCII characters untouched, all others occupy 2 to 4 contiguous bytes
  - Complete codespace can be coded
  - Advantage: „Latin“ texts quite compact and readable in unaware editors
  - Problem: string length, substrings, etc.
-



## UTF-8 coding

- US-ASCII characters untouched, Bit 7 reset 0xxxxxxx
  - All others are sequences of bytes with Bit 7 set 1xxxxxxx
  - First Byte: As many leading bits set as length of sequence: 110xxxxx
  - Following Bytes: Bit 7 set, Bit 6 reset: 10xxxxxx
  - Recognize the type of byte:
    - Complete character: 0xxxxxxx
    - Part of a sequence: 1xxxxxxx
    - Sequence head: 11xxxxxx
    - Sequence tail: 10xxxxxx
-





## UTF-8 coding

- Bits after the type bits remain for coding the Coding Point

ä – LATIN SMALL LETTER A WITH DIAERESIS

$00E4_{16} = 11100100_2$

110xxxxx 10xxxxxx  
    ooo11     100100

-----  
11000011 10100100 bin  
C3           A4       hex  
195          164       dez  
Ã            ä        Latin-1

00000	–	0007F	:	0xxxxxxx
00080	–	007FF	:	110xxxxx 10xxxxxx
00800	–	00FFF	:	1110xxxx 10xxxxxx 10xxxxxx
10000	–	1FFFF	:	11110xxx 10xxxxxx 10xxxxxx 10xxxxxx



# UTF-8 with ISO 8859-1 Rendering

JÃ¼rgen Klinsmann  
(German Soccer Legend)

---



## UCS-4, UTF-32

- 32 Bit per Character
  - Every code point as one word in memory
  - Fat but handy (1 word = 1 character)
  - Endianness issues (UTF-32BE, UTF-32LE, BOM)
  - Complete codespace can be coded
  - No practical differences between UCS-4 and UTF-32
-



## What is Plain Text?

- For every string, for every text (file, e-mail, attachment, download, etc) the encoding must be known.
- Plain text can begin with a BOM

There Ain't No Such Thing As Plain Text.  
-- Joel Spolsky



## Transcoding / Transliteration

- Transform from special character sets to Unicode and back
  - e.g. Windows-1252 -> Unicode -> ISO 8859-1
  - Translation tables: [www.unicode.org](http://www.unicode.org)
  - Characters can get lost (ك becomes ?)
  - Characters can get transformed (ç becomes c)
-



## Sorting

- Sorting rules also apply for searching
  - There are cultural differences
    - treat ä like a
    - treat ä like ae
    - treat ä as a seperate character after z
  - Unicode defines algorithms and delivers tables for sorting
-



## Case Mappings

- Not available in every language
- Not always reversible
- Turkish:  $\text{ı} \rightarrow \text{I}$ ,  $\text{i} \rightarrow \text{İ}$     English:  $\text{i} \rightarrow \text{I}$
- Not necessarily 1:1     $\beta \rightarrow \text{SS}$
- Case Mappings are language dependant



## Comparisons, Sorting

- Case Insensitive

Firebird = FIREBIRD ?

river = RiVeR ?

Fluß = FLUSS ?

a B c <--> B a c

- Accent Insensitive

Amélie = Amelie ?

a é i o ù <--> a i o é ù

---





# Firebird and Character Sets



## CHAR / VARCHAR Fields

- Every CHAR or VARCHAR column has a character set applied by Firebird  
(Remember? „There is no such thing as plain text“)
- This character set will be used for storage
- The character set can be defined when declaring the column:

```
create table persons (  
    pers_id integer not null primary key,  
    last_name  varchar (50) character set iso8859_1,  
    first_name varchar (50) character set iso8859_1  
);
```



## Default Character Set

- A default character set for all string columns can be defined together with CREATE DATABASE:

```
create database 'employee.gdb'  
  default character set ISO8859_1;
```

- You can override the default character set:

```
create table persons (  
  pers_id integer not null primary key,  
  last_name  varchar (50),  
  first_name varchar (50),  
  czech_name varchar (50) character set iso8859_2  
);
```

---



## Text Blobs

- Character Sets also apply to text blobs

```
create table persons (  
    pers_id integer not null primary key,  
    last_name  varchar (50),  
    first_name varchar (50),  
    resume blob sub_type text  
);
```



## Client Character Set

- The Client application defines a character set for its connection
  - All strings will be transliterated to/from this Client Character Set
  - Firebird 1.5: „Unable to transliterate between character sets“
  - Firebird 2.0: ServerCS <–> Unicode <–> ClientCS
-



## How to define the Client Character Set

- IB Objects:

```
IB_Connection1.CharSet := 'ISO8859_1';
```

- IBX:

```
IbDatabase1.Params.Add ('lc_ctype=ISO8859_1');
```

- PHP:

```
$db = ibase_connect ($Name, $Usr, $Pwd, "ISO8895_1");
```

---



## Collations

- Define sort ordering for ORDER BY
- Define casing rules for UPPER(), LOWER()

```
SELECT *  
FROM ...  
WHERE UPPER (NAME COLLATE DE_DE) = :SEARCHNAME  
ORDER BY LASTNAME COLLATE FR_CA
```

- Collations defined per Character Set
- Examples:

```
ISO8859_1: DE_DE, DU_NL, FR_FR, FR_CA, PT_PT, PT_BR  
WIN1250: WIN1250, BS_BA, WIN_CZ, WIN_CZ_CI_AI
```



## Defining the Collation for a Column

- There is no such thing as a default collation
- You can define the standard collation for a column:

```
create table persons (  
    pers_id integer not null primary key,  
    last_name  varchar (50) collate de_de,  
    first_name varchar (50) collate de_de,  
);
```





## UPPER(), LOWER()

- Firebird 1.0/1.5: UPPER() only works correctly if there is a collation defined for the parameter field.

Without collation no uppercasing of letters outside the a..z range.

- Firebird 2.0: UPPER() will return uppercased characters for all characters
  - Firebird 2.0: New LOWER() function
-



## Searching (Case insensitive)

- WHERE LIKE, WHERE STARTING WITH, WHERE =

```
select * from persons  
where upper (last_name) like '%TARK%'
```

```
select * from persons  
where upper (last_name) starting with 'STARK'
```

```
select * from persons  
where upper (last_name collate de_de) like '%Ä%'
```



## Indexed Case Insensitive Searches

- Firebird 1.5: Shadow Field and Trigger

```
create table persons (  
    name          varchar (50) collate de_de,  
    name_upper    varchar (50) collate de_de);
```

```
create index idx_person_name on persons (name_upper);
```

```
create or replace trigger biu_persons for persons  
before insert or update as  
begin  
    new.name_upper = upper (new.name);  
end;
```

```
select * from persons  
where name_upper = 'LÖW'
```

---



## Indexed Case Insensitive Searches

- Firebird 2.0: Expression Index

```
create table persons (  
  name varchar (50) collate de_de,  
  city varchar (50));
```

```
create index idx_person_name on persons  
  computed by ( upper (name) );
```

```
select * from persons  
where upper (name) = 'FIREBIRD'
```

```
create index idx_person_city on persons  
  computed by ( upper (city collate de_de) );
```

```
select * from persons  
where upper (city collate de_de) = 'MÜNCHEN'
```

---



## Firebird and Unicode

- Unicode internally stored as UTF-8
  - Character Set **UNICODE\_FSS**: Up until Firebird 1.5: an old version of UTF8 that accepts malformed strings and does not enforce correct maximum string length.
  - Character Set **UTF8**: Firebird 2.0: Replacement for UNICODE\_FSS
  - Unicode used to transliterate between character sets
  - Unicode collation implemented for comparisons and casings
-



## Unicode Collations for UTF8

- UCS\_BASIC: sorts in Unicode Code Point order
- UNICODE: uses the Unicode Collation Algorithm

```
select * from t order by c collate ucs_basic;  
A B a b á
```

```
select * from t order by c collate unicode;  
a A á b B
```



## Examples

- `LanguageMix.htm`
- `CharsetDemo.fdb`
- `CharsetDemo.dpr`
- `charsetdemo.php`



## Special Character Sets

- **NONE:** No character set applied. With this character set setting, Firebird is unable to perform conversion operations like UPPER() correctly on anything other than the standard 26 latin letters.
  - **OCTETS:** Same as NONE. Cannot be used as client connection character set. Space character is #x00
  - **ASCII:** US-ASCII
  - **UNICODE\_FSS:** an old version of UTF8 that accepts malformed strings and does not enforce correct maximum string length. Replaced by **UTF8** in FB2.0.
-





## Character Sets and Collations

- **ISO8859\_x** (e.g. ISO8895\_1, ISO8859\_2)  
Collations: DE\_DE, FR\_FR, CS\_CZ, etc.
  - **WIN125x** (e.g. WIN1252, WIN1250)  
Collations: WIN1252, WIN\_PTBR, PXW\_CSY, etc.
  - **DOSxxx** (e.g. DOS850, DOS852)  
Collations: DB\_DEU850, PDOX\_CSY
  - Complete List:  
[www.destructor.de/firebird/charsets.htm](http://www.destructor.de/firebird/charsets.htm)
-



## Other Character Sets

- BIG\_5: Chinese
  - CYRL, KOI8-R, KOI8-U: Russian, Ukrainian
  - KSC\_5601: Korean
  - SJIS\_0208: Japanese
  - EUCJ\_0208: Japanese
  - GB\_2312: Chinese
-



## Which one to select?

- DOS, dBASE and Paradox only for legacy support
  - WINxxx is extension of corresponding ISOxxx, but may lead to problems on non-Windows systems.
  - ISOxxx is missing a few characters of WINxxx (e.g. typographic dash signs) -> prepare to handle this
  - If you expect to store various languages, go for Unicode UTF-8 but use Firebird 2.0 then.
-



## Links

- Unicode Consortium, ISO 10646  
Book: The Unicode Standard, Version 5.0  
ISBN 0-321-48091-0  
[www.unicode.org](http://www.unicode.org)
  - Troy Wolbrink: Tnt Delphi Free Unicode Controls  
[www.tntware.com](http://www.tntware.com)
  - My Firebird Website and Conference Blog/Gallery  
[www.destructor.de/firebird](http://www.destructor.de/firebird)
  - Papers of this talk  
[www.destructor.de/talks](http://www.destructor.de/talks)
-



# Questions?

**Stefan Heymann**

[www.destructor.de](http://www.destructor.de)  
[stefan@destructor.de](mailto:stefan@destructor.de)

---



**Thank You!**  
**Danke! Merci! Grazie!**  
**Gracias! Obrigado!**  
**Děkuji! Paldies! Hvala!**  
**СПАСИБО**

**Stefan Heymann**

[www.destructor.de](http://www.destructor.de)  
[stefan@destructor.de](mailto:stefan@destructor.de)

---