# Character Sets and Unicode in Firebird

### **Stefan Heymann**

www.consic.de heymann@consic.de

After a short introduction to the world of Character Sets and Unicode, this session will show you how to bring it all to work in Firebird. You will learn what all those character sets and collations are and how you can properly use them to get the right characters into the database and onto your screen.

# **Topics**

- Characters
- Character Sets
- Unicode
- Firebird
- Examples

Session: Character Sets and Firebird

Speaker: Stefan Heymann Page: 3

## **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 <u>Character</u> is an abstract concept
- Rendering of characters as <u>Glyphs</u> is the job of the rendering machine (Postscript, GDI, TrueType, Web Browser, etc.)
- We mostly care for processing the characters
- A <u>Character Set</u> assigns a number to a character:

Uppercase A = 65

Uppercase B = 66

etc.

# **Glyphs**

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

Keystrokes	j	I	7		غ	غ		
Input characters	ن	I	ئ	I	غ	غ		
Encoded characters (byte values in hex)	0644	0627	0644	0627	0639	0639		
Display	لالاغغ							

### **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 <u>English</u>
- 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: 128 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**

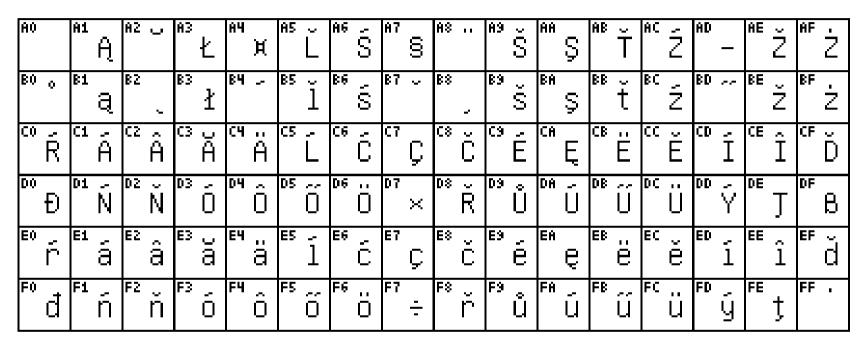
- Characters 0..127 identical to ASCII
- 128..159 undefined control characters
- 160..255 individually assigned

## ISO 8859-1 / Latin-1

AO		A1 j	A	ф	яз £	A4 其	as ¥	A6	A7 S	A8	A9 (C)	AA a	AB {{	AC ¬	AD —	AE ®	AF _
BO	۰	B1 <u>+</u>	: B:	2 2	83	84	85 µ	B6 ¶	B7 •	B8 _	B9 1	BA Q	BB }}	BC 1/4	BD 1/2	BE ¾	BF ن
CO	Ã	cı A	Ci	Â	ЗÃ	сч А	cs Å	ce Æ	°7 Ç	œÈ	ΘÉ	ca Ê	св <u>:</u>	αĨ	ΩÍ	Î	CF ::
DO	Đ	D1 ~	i I	Õ	D3 Ó	Ô	os õ	De	D7 ×	os Ø	D9 Ù	DA Ú	DB Û	Ü	PΡ	DE Þ	В
ΕO	à	ы á	E	â	¤ ã	ëч "	å	æ €6	e7 Ç	ĕ	é	ê	ë	EC Ĩ	ED 1	î	EF 1
F0	ð	F1 ñ	F	Ò	F3 Ó	F٩ Ô	FS Õ	F6 Ö	F7 ÷	F8 Ø	F9 Ü	FA Ú	fb Û	FC Ü	FD Ý	FE Þ	FF

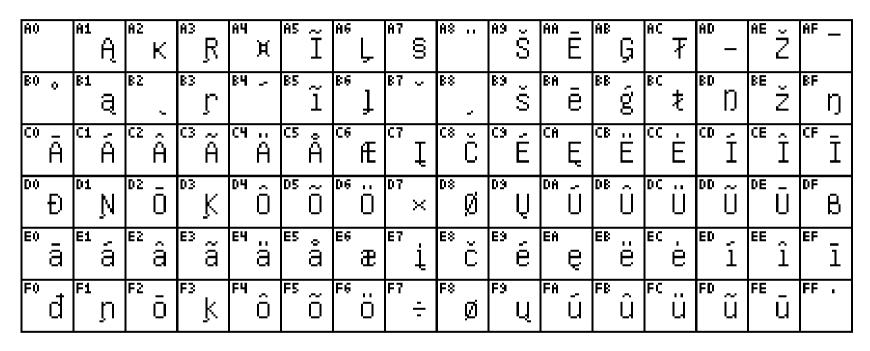
Afrikaans, Albanese, Basque, Danish, German, **English**, Faroese, Finnish, French, Icelandic, Italian, Catalan, Dutch, Norwegian, Portuguese, Rhaeto-Romance, Scottish Gaelic, Schwedish, Spanish, Suaheli *Large parts of the world – wide-spread use* 

## ISO 8859-2 / Latin-2



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

## ISO 8859-4 / Latin-4



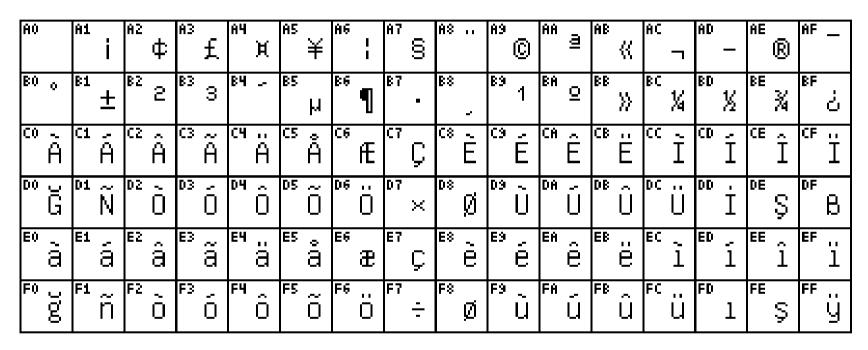
Northern Europe, Baltic, Greenlanic, Sami

## **ISO 8859-5 / Cyrillic**

A0	A1	Ë	A2 .	Ъ	ÅЗ	<b>вч</b> Е	A5 S		I	<sup>A7</sup> Ï	A8	J	<sup>н9</sup> Љ	<sup>ва</sup> Њ	AB	Ћ	ac É	• [	_	РΕ	AF ↓
ВО	B1	Б	B2 	В	ВЗ	вч Д	B5 E	- B	Ж	3	B8	И	вэ И	BA K	BB	Л	PC P	BD	Н	BE ()	BF ∏
°Р	C1	C	CZ	Т	СЗ	СЧ	C5 }		ЕП	ч	C\$	Ш	СЭЩ	ся Ъ	СВ	Ы	űΕ	CD	Э	СЕЮ	Я
ю а	D1	б	D2	В	D3 	оч Д	D5 E		Ж	07 3	D8	и	<sup>09</sup> О	DA K	DB	Л	pc DC	DD I	Н	DE O	DF 
	D1 E1	б		В	Γ	оч Д Еч Ф		<b>?</b>	Ж	3	E\$	И	Й		DB EB	Л	_	  ED	Н	DE O EE HO	рғ П Еғ Я

Cyrillic (Russia, Ukraine, etc.)
More important: KOI8-R (Russian), KOI8-U (Ukrainian)

## ISO 8859-9 / Latin-5



**Turkish** 

# Windows Character Sets ("Codepages")

- Partly 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

## **Microsoft Windows Codepages**

- 874 Thai
- 932 Japanese
- 936 Simplified Chinese
- 949 Korean
- 950 Traditional Chinese
- 1250 Central European
- 1251 Cyrillic
- 1252 Western European
- 1253 Greek
- 1254 Turkish
- 1255 Hebrew
- 1256 Arabic
- 1257 Baltic
- 1258 Vietnamese

#### windows-1252

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

Euro sign € since 2000 (Windows 1252-2000)

## **Multi-Byte Character Sets MBCS**

- Multiple Bytes per Character
- Eastern Asian Languages (CJK)
- String Length <> 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 **10**646 (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
- "Basic Multilingual Plane" (BMP): 0..U+FFFF
   Can be encoded in 16 bits
- Current version: 6.0.0 (February 2011)
- 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 4 or 5 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 several codings around:
  - 8-Bit (UTF-8, formerly called "FSS")
  - 16-Bit (UCS-2, UTF-16)
  - 32-Bit (UCS-4, UTF-32)
  - PunyCode for international Domain names
  - "Exotic": UTF-7, UTF-1, etc.

#### UCS-2

- 16 Bits per Character
- Code Area: 0000..FFFF
   Basic Multilingual Plane (BMP)
- Characters beyond the BMP (defined since Unicode 3.1) can not be encoded
- 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 > FFFF must be coded as a "Surrogate Pair" and occupy 2 consecutive 16-Bit words
- Complete Code Space can be encoded
- Difficult to calculate string length or substrings
- Used by Windows 2000 and later
- Used by Delphi 2009/2010/XE (UnicodeString type)

#### **Endianness**

- Problem with UCS-2, UTF-16: Low/High-Byte ordering
- Differentiate in UTF-16BE and UTF-16LE in metadata
- Byte Order Mark <u>BOM</u> 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
- "Unicode" is used as a (bad) synonym for UCS-2 or UTF-16
- WideString, wchar\_t
- Windows NT3, NT4: UCS-2
- Since Windows 2000: UTF-16
- Delphi 2009/etc. UnicodeString: UTF-16

#### UTF-8

- Coding as 8-Bit strings
- Called "File System Safe" (FSS) in its early days
- 7-bit US-ASCII characters untouched, all others occupy 2 to 4 consecutive bytes
- Complete codespace can be encoded
- Advantage: "Latin" texts quite compact and readable in unaware editors
- Problem: string length, substrings, etc.
- No BOM necessary, but sometimes used as an indicator for UTF-8 text

## **UTF-8** coding

- US-ASCII characters untouched, Bit 7 reset 0xxxxxxxx
- All others are sequences of bytes with Bit 7 set 1xxxxxxx
- Head Byte: As many leading bits set as length of sequence: 110xxxxx
- Tail Bytes: Bit 7 set, Bit 6 reset: 10xxxxxxx
- 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 encoding the Code Point

 $\ddot{a}$  - LATIN SMALL LETTER A WITH DIAERESIS  $00E4_{16} = 228_{10} = 11100100_2$ 

110xxxxx 10xxxxxx 00011 100100 00000 - 00007F: 0xxxxxx

00080 - 0007FF: 110xxxxx 10xxxxxx

00800 - 00FFFF: 1110xxxx 10xxxxxx 10xxxxxx

10000 - 1FFFFF: 11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

11000011 10100100 bin

C3 A4 hex

195 164 dez

## UTF-8 with ISO 8859-1 Rendering

Jürgen Klinsmann

## **UCS-4, UTF-32**

- 32 Bits per Character
- Every code point as one 32-bit 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 (ISO) and UTF-32 (Unicode)

### 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**

- Transformation from special character sets to Unicode and back
- e.g. Windows-1252 -> Unicode -> ISO 8859-1
- Translation tables: www.unicode.org
- Characters can get lost (خا 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: ı -> I, i -> İ English: i -> I
- Not necessarily 1:1 ß → SS
- Case Mappings are language dependant

# **Comparisons, Sorting**

Case InsensitiveFirebird = FIREBIRD ?

river = RiVeR ?

Fluß = FLUSS?

a B c <--> B a c

Accent Insensitive
 Amélie = Amelie ?

aéioù <--> aioéù

# **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.x: ServerCS <-> Unicode <-> ClientCS

# **How to define the Client Character Set**

```
IBObjects:

IB_Connection1.CharSet := 'ISO8859_1';
IBX:

IbDatabase1.Params.Add ('lc_ctype=ISO8859_1');
IBDAC:

Connection.Options.CharSet := 'ISO8859_1';

Connection.Options.UseUnicode := false;

Connection.Options.EnableMemos := false;
PHP:

$db = ibase_connect ($Name, $Usr, $Pwd, "ISO8859_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
```

- Define comparison rules (A = B, A <> B)
   WHERE NAME COLLATE UNICODE\_CI =
   : PERSNAME COLLATE UNICODE\_CI
- 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 UTF8: UCS_BASIC, UNICODE, UNICODE_CI, UNICODE_CI_AI
```

# **Defining the Collation for a Column**

- There is no such thing as a default collation
- Oh, wait: there IS since Firebird 2.5
- 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, ...);
```

Or you can define it with every string usage:
 where name collate unicode\_ci = myname

# **Default Collations**

- Feature introduced in Firebird 2.5 / ODS 11.2
- Define the default collation for the default character set in CREATE DATABASE:

```
create database <file name>
    [ page_size <page size> ]
    [ length = <length> ]
    [ user <user name> ]
    [ password <user password> ]
    [ set names <charset name> ]
    [ default character set <charset name>
        [ [ collation <collation name> ] ]

create database elvis:presley
    page_size 16384
    user presley
    password guitar
    default character set iso8859_1
        collation de_de;
```

# **Default Collations (continued)**

 New syntax to change the default collation for existing databases:

```
ALTER CHARACTER SET <charset_name>
SET DEFAULT COLLATION <collation_name>
```

### Example:

```
ALTER CHARACTER SET ISO8859_1
SET DEFAULT COLLATION DE_DE;
```

# **User-Defined Collations**

- CREATE COLLATION
- Poorly documented
- Examples don't work
- Dead? Badly documented? More investigation necessary

# 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.x: UPPER() will return uppercased characters for all characters, no collation required
- Firebird 2.x: New LOWER() function

# **Case insensitive Searching**

WHERE LIKE, WHERE STARTING WITH, WHERE =

# **Indexed Case Insensitive Searches**

# Firebird 1.5: Shadow Field and Trigger

# **Indexed Case Insensitive Searches**

# Firebird 2.x: 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: an old version of UTF8 that accepts malformed strings and does not enforce correct maximum string length. Fixed in Firebird 2.5.
- Character Set **UTF8** since Firebird 2.0: Replacement for UNICODE\_FSS
- Unicode used for transliteration between character sets: CS <-- Unicode --> CS
- Unicode collation implemented for comparisons and casings (ICU DLLs)

# **Unicode Collations for UTF8**

- UCS\_BASIC: sorts in Unicode Code Point order
- UNICODE: uses the Unicode Collation Algorithm
- UNICODE\_CI: Case-insensitive [FB2.1]
- UNICODE\_CI\_AI: Case-/Accent-insensitive [FB2.5]

```
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
```

# **Special Character Sets**

- **NONE:** Plain octets, 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 (English only)
- **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

# **Other Character Sets**

- BIG\_5: Chinese
- CYRL, KOI8-R, KOI8-U: Cyrillic, 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 several languages, use Unicode UTF8

# Links

- The Unicode Consortium. The Unicode Standard www.unicode.org
- My Firebird Website and Conference Blog/Gallery www.destructor.de/firebird

# **Questions?**

### **Stefan Heymann**

www.consic.de heymann@consic.de

www.destructor.de stefan@destructor.de

# Thank You! Danke! Merci! Grazie! iGracias! Obrigado! Děkuji! Paldies! Hvala! CΠΑCИБО