



What Developers Should Know about Character Sets, Unicode etc.

Stefan Heymann

www.destructor.de
stefan@destructor.de



Topics

- Why?
- Fonts, Glyphs
- ASCII & Co.
- Unicode
- Application
- Firebird and Character Sets



Why?



Did you ever have to specify ...

- Content-Type in HTML?
- encoding in XML?
- CHARACTER SET in Firebird?
- NLS_LANG in Oracle?
- ...



Growing Requirements

- ASCII is enough ...
- Internationalization of User Interfaces
- Internationalization of Websites
- Handling of Texts from various languages
- Development tools built (mainly) by native English speakers



Joel on Software

So I have an announcement to make: if you are a programmer working in 2003 and you don't know the basics of characters, character sets, encodings, and Unicode, and I catch you, I'm going to punish you by making you peel onions for 6 months in a submarine. I swear I will.

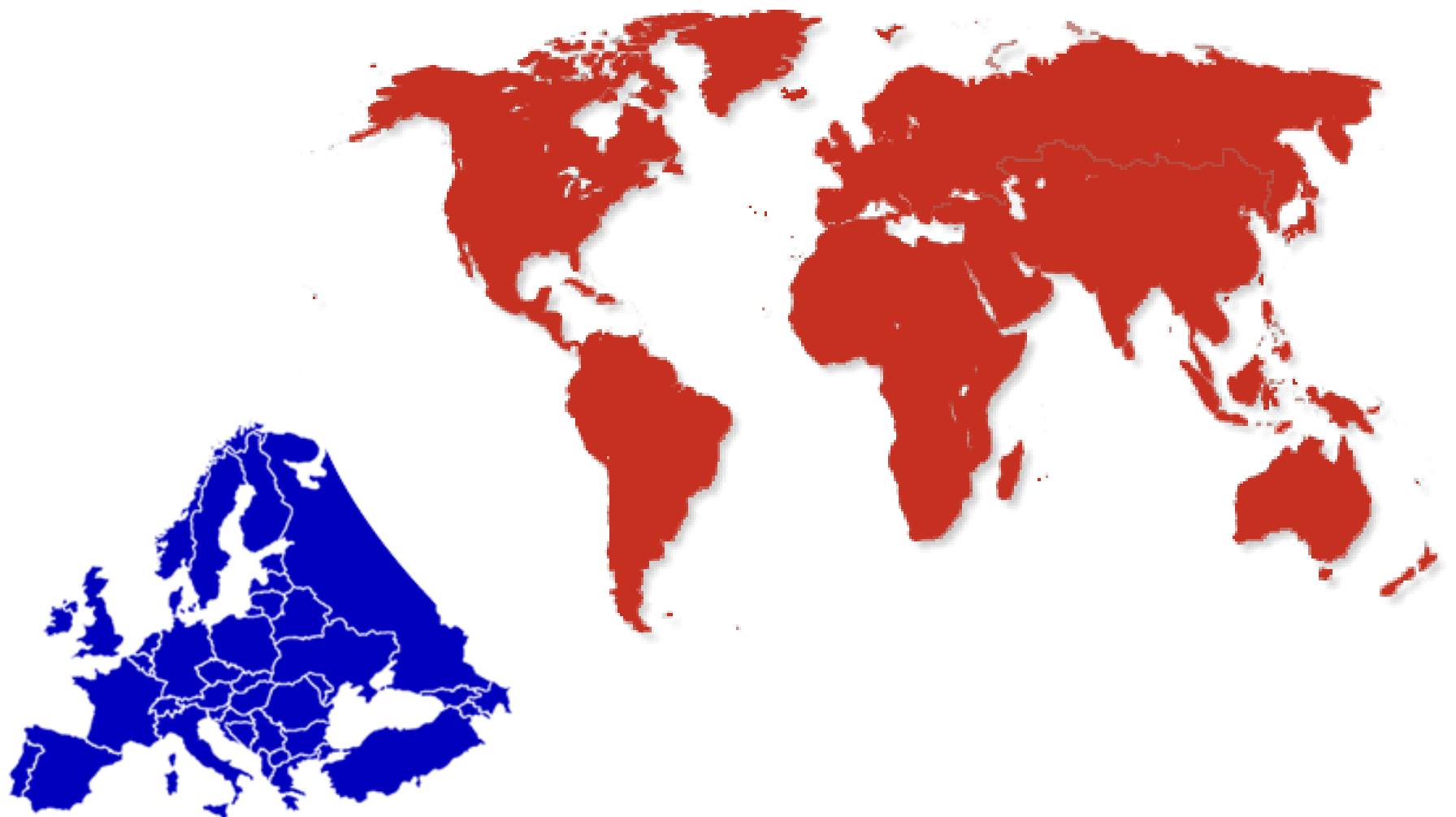
-- Joel Spolsky, www.joelonsoftware.com



Fonts, Glyphs



Characters of the World





Glyphs vs. Characters

A

A

A

A

A

A

A



Glyph, Character, Character Set

- Rendering of characters as glyphs is the job of the rendering machine (Postscript, GDI, TrueType, etc.)
- As developers we mostly care for processing the characters
- A character set encodes characters as numerical values
A = 65



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	لا لا غ غ					



ASCII & Co.



The Mother of Character Sets

- American Standard Code for Information Interchange: ASCII, ISO 646
- 7 bits wide, characters ranging from 0 to 127
- 32 invisible control characters (NUL, CR, LF, FF, BEL, ESC, ...)
- Alphabet, Digits, Punctuation (;.-?)
- Optimized for English
- No accents and umlauts



ASCII for Europe

- Seldom used characters get new assignment
- [=Ä \=Ö]=Ü
- Problem: Printer and screen must have same setting
- No mixing in one text possible: „Amélie knackt gerne die Kruste von Crème Brulé mit dem Löffel“



Use the 8th Bit

- 128 new characters
- A lot of character sets have evolved
- $\text{ISO 8859-}x = \text{ASCII} + 160..255$
- $\text{ISO 8859-1} = \text{Latin-1}$
(Western European languages)
- $\text{Windows 1252} = \text{ISO 8859-1} + 128..159$



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
Ř	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	Ń	Ń	Ō	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	ĺ	č	ç	č	ě	ě	ě	ě	ī	ī	ď
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF
đ	ñ	ň	ō	ô	õ	ö	÷	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	Š		Š	Ē	Ġ	Ʀ	–	Ž	–
B0	°	ā	˙	˙	ĩ	ł	˘		š	ē	ġ	ʀ	ŀ	ž	ŋ
C0	Ā	Ā	Ā	Ā	Ā	Æ	ı	č	ē	ē	Ē	Ē	ĩ	ĩ	ĩ
D0	Đ	Ń	Ō	Ȧ	Ō	Ö	×	Ø	Ů	Ú	Û	Ü	Ũ	Ū	ß
E0	ā	ā	ā	ā	ā	æ	ı	č	ē	ē	ē	ē	ĩ	ĩ	ĩ
F0	đ	ñ	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
B1	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
C0	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю
C1	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
C2	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я	
D0	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о
D1	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
D2	в	г	д	е	ж	з	и	й	к	л	м	н	о	п	
E0	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю
E1	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
E2	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я	
F0	ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч
F1	ё	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч
F2	ђ	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч	
F3	ѓ	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч		
F4	є	ѕ	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч			
F5	ѕ	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч				
F6	і	ї	ј	љ	њ	ћ	ќ	–	ў	ч					
F7	ї	ј	љ	њ	ћ	ќ	–	ў	ч						
F8	ј	љ	њ	ћ	ќ	–	ў	ч							
F9	љ	њ	ћ	ќ	–	ў	ч								
FA	њ	ћ	ќ	–	ў	ч									
FB	ћ	ќ	–	ў	ч										
FC	ќ	–	ў	ч											
FD	–	ў	ч												
FE	ў	ч													
FF	ч														

Cyrillic (Russia, Ukraine, etc.)

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



ISO 8859-6 / Arabic

A0				A4هـ								ACء	ADـ		
											BBة				BF؟
	C1ع	C2آ	C3أ	C4ؤ	C5إ	C6ي	C7ا	C8ب	C9ة	CAت	CBث	CCج	CDح	CEخ	CFد
D0ذ	D1ر	D2ز	D3س	D4ش	D5ص	D6ض	D7ط	D8ظ	D9ع	DAغ					
E0ـ	E1ف	E2ق	E3ك	E4ل	E5م	E6ن	E7و	E8و	E9ى	EAي	EB#	EC،	ED#	EEر	EF،
F0ر	F1و	F2و													

Arabic

Does not contain all characters, so it's not used often)



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
C0	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
E0	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

Turkish



Windows Character Sets

- Mostly congruent to ISO-8859
- Additional assignment of characters 128..159 with visible (non-control) characters
- Hyphen, typographic „quotation marks“, etc.
- Windows character sets officially registered at IANA



windows-1252

- Congruent to ISO 8859-1
- Additional characters in the 128..159 range:

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

€ sign since 2000



Multi-Byte Character Sets MBCS

- Multiple Bytes per Character
- Eastern Asian Languages (CJK)
- String Length \neq Length of character chain
- Making sub-strings more difficult



Unicode



Why Unicode?

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



Unicode

- Started with 16 Bits/Character, now 32 Bits/Char
- Possibility to code 1,114,112 characters
- Only a fraction is used
- Current version: 4.1.0
- 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: hexadecimal, with preceding „U+“, 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 too fat
- There are various codings around:
 - 8-Bit (UTF-8)
 - 16-Bit (UCS-2, UTF-16)
 - 32-Bit (UCS-4, UTF-32)
 - UTF-7, UTF-1, PunyCode, ...



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 to 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 need 2 to 4 contiguous bytes
- Complete codespace can be coded
- Advantage: „Latin“ texts quite compact
- 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`



UTF-8 coding

- You can determine the type of each byte:
 - Complete character: 0xxxxxxx
 - Part of a sequence: 1xxxxxxx
 - Sequence Start: 11xxxxxx
 - Sequence body: 10xxxxxx



UTF-8 coding

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

ä – LATIN SMALL LETTER A WITH DIAERESIS

$00E4_{16} = 11100100_2$

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

110xxxxx	10xxxxxx
ooo11	100100

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



UTF-8 with ISO rendering



JÃ¼rgen Klinsmann
(German Soccer Coach)



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)
- 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



Encodings Overview

- ISO 8859-x
- Windows-12xx
- KOI8-R, KOI8-U
- Shift-JIS
- etc.
- Unicode: UTF-8, (UCS-2), UTF-16, UTF-32
- There is *no* „Unicode“ encoding!



Unicode Specials

Skip



Character Ranges

- General Scripts (Latin, Greek, Cyrillic, etc.)
- Symbols (Symbol, Dingbat, Punctuation, Math, etc.)
- CJK Phonetics and Symbols
- CJK Ideographs
- Yi Syllables, Hangul Syllables
- Surrogate Area, Private Use Area
- Compatibility and Specials



Control Characters

- 0000..001F: ASCII control characters
- 007F: ASCII Delete/Rubout
- 0080..009F: Control characters



Combined Characters

- There are 2 ways to code an ä:
- 00E4: LATIN SMALL LETTER A WITH DIAERESIS
- 0061: LATIN SMALL LETTER A +
0308: COMBINING DIAERESIS
- You can „compose“ characters/glyphs using
Combining Characters



Normalization

- Problem: You cannot compare text byte-for-byte
- Normalisation/Canonicalisation required
- Translation tables
- Composition, Decomposition
- Technical Report #15: Unicode Normalization Forms



Bi-directionality

- Unicode defines how characters are arranged in memory
- Some languages go from right-to-left
- E.g. Arabic, Hebrew
- Rules and control characters for left-to-right and right-to-left text



Transcoding

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



Surrogate Pairs for UTF-16

- For coding characters $> \text{U+FFFF}$
- Code Points D800..DFFF reserved
- Two contiguous 16-Bit words used: High and Low Surrogate = Surrogate Pair
- High Surrogate: U+D800..U+DBFF
- Low Surrogate: U+DC00..U+DFFF



Surrogate Pairs for UTF-16

- You can see if a 16-Bit word ...
 - is a non-surrogate character
 - belongs to a surrogate pair
 - is the high or low surrogate



Sorting

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



Case Mappings

- Not available in every language
- Not always reversible
- Turkish: $\text{ı} \rightarrow \text{I}$, $\text{i} \rightarrow \text{İ}$
- No necessarily 1:1 $\beta \rightarrow \text{SS}$
- Case Mappings are language depending
- Unicode Technical Report #21
„Case Mappings“



The Unicode Standard

- Unicode Consortium www.unicode.org
- ISO 10646
- Book: The Unicode Standard, Version 4.0
ISBN 0-321-18578-1
- All chapters and files, complete character database
available for free at www.unicode.org
- Technical Reports



Application

Skip



Windows Notepad

- Save as ... > Encoding
 - ANSI (= Windows 1252)
 - Unicode (= UCS-2 or UTF-16)
 - Unicode Big Endian
 - UTF-8



XML

- Every XML entity is Unicode by definition. Default is UTF-16 or UTF-8.
- `<?xml version="1.0" encoding="utf-8" ?>`
- With *every* encoding, every Unicode character can be inserted as a *Character Reference*
 - `€` €
 - `ॐ` ॐ „DEVANAGARI OM“



HTML

- Unicode aware since HTML 4.0
- Character References like XML
- <head>

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



Windows API

- Win32 API is Unicode aware from the beginning
- All string function also available as „W“ function
- UCS-2 up to Windows NT 4.0
- UTF-16 since Windows 2000



Borland Delphi

- VCL controls are not Unicode aware (local character sets usable)
- Troy Wolbrink: Tnt Delphi Unicode Controls
<http://tnt.ccci.org/>
- Mike Lischke (Virtual Treeview, etc.)
<http://www.delphi-gems.com/>
- JCL: JclUnicode.pas
- Delphi Roadmap: VCL controls Unicode aware some time in the future



Oracle

- Character set defined when database is created
- Client can have other characters set, automatic transliteration
- Environment variable NLS_LANG
- WE8MSWIN1252, WE8ISO8859P1, UTF8, etc.



Firebird and Character Sets



Character Sets

- Default Character Set defined at database create time:

```
CREATE DATABASE 'employee.gdb'  
DEFAULT CHARACTER SET ISO8859_1;
```

- Different Character Set for every table column

```
CREATE TABLE users (  
  czech_name VARCHAR(50) CHARACTER SET ISO8859_2
```

- Default Collation (sort order) for every table column

```
CREATE TABLE test (  
  name VARCHAR(50) COLLATE DE_DE );
```

- Available collations depending on character set used



Collations

- Define sort ordering
- There is no default collation for the database
- `COLLATE DE_DE`: German sorting, applying the rules for Germany
- `COLLATE FR_CA`: French sorting, applying the rules for Canada

```
ORDER BY LASTNAME COLLATE FR_CA, FIRSTNAME COLLATE FR_CA  
WHERE LASTNAME COLLATE FR_CA = :lastnametosearch
```



UPPER()

- The UPPER() function only works correctly if there is a collation defined for the parameter field
- Situation will improve with Firebird 2.0 (which will also bring a new LOWER() function)



Local Character Set

- Client side character set can be different from server side character set. Automatically transliterated (sort of ;-)
- Conversions between character sets are always done
CHARSET1 -> UNICODE -> CHARSET2
- With NONE or OCTETS bytes are just copied:
NONE/OCTETS -> CHARSET2 and CHARSET1 -> NONE/OCTETS.
- IB Objects
Ib_Connection1.CharSet := 'ISO8859_1';
- IBX:
IbDatabase1.Params.Add ('lc_ctype=ISO8859_1');



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



„Normal“ Character Sets

- **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



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



Unicode „Character Sets“

- **UNICODE_FSS:** Unicode UTF-8. Old implementation that accepts malformed strings and does not enforce correct max. string length. *DEPRECATED*.
- **UTF8:** New in Firebird 2.0
 - Collation **UCS_BASIC:** sorts in Unicode code-point order
 - Collation **UNICODE:** Sorts using the Unicode Collation Algorithm (UCA)



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
- You don't care much about character set issues? And you work for one platform anyway? Select the appropriate ISOxxx character set. It will store everything correctly (you don't care about typographic dashes, do you ;-) -- and you have the best selection of collations.



Links

- Unicode specification, tables, articles:
www.unicode.org
- My Firebird Website (with an article about and a complete listing of Firebird character sets):
www.destructor.de/firebird



Thank You!

Stefan Heymann

www.destructor.de
stefan@destructor.de