

better  
office

# Using Unicode and Other Encodings in your Programs

Jeroen Pluimers  
better office benelux  
[jpluimers@better-office.com](mailto:jpluimers@better-office.com)

Microsoft  
.net

CODE  
GEAR

Microsoft  
.net

## Note

- This is an overview session
  - It explains about Unicode and Encodings
  - It points to lots of references for further Delphi and Unicode reading
  - It has very few demos
  - But it does some show pitfalls

better  
office

CODE  
GEAR

Microsoft  
.net

## ASCII - Encoding - Display

Computer codes	Encoding table	Characters
41	41	A
42	42	B
43	43	C
44	44	D
45	45	E
46	46	F
47	47	G
48	48	H
49	49	I
50	50	J
51	51	K
52	52	L
53	53	M
54	54	N
55	55	O
56	56	P
57	57	Q
58	58	R
59	59	S
60	60	T
61	61	U
62	62	V
63	63	W
64	64	X
65	65	Y
66	66	Z
67	67	[
68	68	\
69	69	]
70	70	^
71	71	_
72	72	`
73	73	a
74	74	b
75	75	c
76	76	d
77	77	e
78	78	f
79	79	g
80	80	h
81	81	i
82	82	j
83	83	k
84	84	l
85	85	m
86	86	n
87	87	o
88	88	p
89	89	q
90	90	r
91	91	s
92	92	t
93	93	u
94	94	v
95	95	w
96	96	x
97	97	y
98	98	z
99	99	{
100	100	
101	101	}
102	102	~
103	103	
104	104	
105	105	
106	106	
107	107	
108	108	
109	109	
110	110	
111	111	
112	112	
113	113	
114	114	
115	115	
116	116	
117	117	
118	118	
119	119	
120	120	
121	121	
122	122	
123	123	
124	124	
125	125	
126	126	
127	127	

- This is what we have been doing forever
- Bytes were equal to characters:
  - C := Char(B);


better  
office

CODE  
GEAR


Picture courtesy of <http://www.sketchpad.net/macwinfontintro.1.htm>

Microsoft  
.net

## Data – Encoding – Display



ASCII  
 CP 437  
 Windows 1252  
 UTF-8  
 ...



better  
office

Pictures courtesy of  
<http://datatlib.ed.ac.uk/> and <http://www.ascenderfont.com/info/simplified-chinese-fonts.aspx>

CODE  
GEAR

Microsoft  
.net

## Some encodings

- Single byte (8-bit)
  - ASCII aka ISO 646 (actually 7-bit; only 0-127 are defined)
  - EBCDIC aka CP390 – if you do AS/400, iSeries, OS390/Mainframes
  - ISO 8859 split into 'regions' as ...1 through ...15, except ...12
  - DOS Code Pages most popular: CP437, CP850, CP852, CP858
  - Windows Code Pages most popular: Windows-1250, Windows-1252
- Two byte (16-bits)
  - ISO-2022 for ...JP, ...CN, ...KR (kinds of 7-bit multi byte encoding)
  - EUC for ...JP, ...CN, ...KR (kinds of 8-bit multi byte encoding)
- Multi byte (8, 16 or 32-bits)
  - Unicode UTF-8, UTF-16, UTF-32
  - GB 18030 Chinese; mapable to Unicode
- Wikipedia has tables for most 8-bit encodings
- UTF8 encodings tables: <http://www.utf8-chartable.de>
- See: [http://en.wikipedia.org/wiki/Character\\_encoding](http://en.wikipedia.org/wiki/Character_encoding)  
[http://en.wikipedia.org/wiki/Code\\_page](http://en.wikipedia.org/wiki/Code_page)  
[http://en.wikipedia.org/wiki/Multi-byte\\_character\\_set](http://en.wikipedia.org/wiki/Multi-byte_character_set)  
[http://en.wikipedia.org/wiki/CJK\\_characters](http://en.wikipedia.org/wiki/CJK_characters)  
[http://en.wikipedia.org/wiki/Comparison\\_of\\_Unicode\\_encodings](http://en.wikipedia.org/wiki/Comparison_of_Unicode_encodings)

better  
office

CODE  
GEAR

Microsoft  
.net

## Unicode fonts

- No font can have all Unicode code points
  - Unicode defines 100k+ code points
  - Fonts have a maximum of 65k code points
- Some fonts having many Unicode code points
  - GNU Unifont <http://unifoundry.com/unifont.html>
  - Code2000 [http://www.code2000.net/code2000\\_page.htm](http://www.code2000.net/code2000_page.htm)
  - New Gulim / GulimChe [http://r.office.microsoft.com/r/rlidDownloadDetail?p1=2002&p2=ie\\_ko](http://r.office.microsoft.com/r/rlidDownloadDetail?p1=2002&p2=ie_ko)
  - Arial Unicode MS included since MS Office 2000
  - Everson Mono <http://www.evertype.com/emono/>
- Notable font
  - Lucida Sans Unicode – first ever font supporting Unicode
- See: [http://en.wikipedia.org/wiki/Unicode\\_typefaces](http://en.wikipedia.org/wiki/Unicode_typefaces)

better  
office

CODE  
GEAR

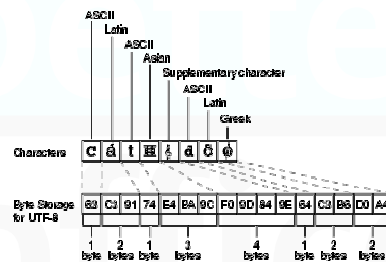
## Explanation of encodings

- Encoding methods as falling into two main categories:
  - The Old Way:  
Each character has a specific and direct representation on a computer.
  - The New Way:  
Each character is a concept, and can be represented in multiple ways.
- The new way is what Unicode does:
  - it maps every letter of every language to a unique number (called a code point),
  - so an upper-case "D" in Unicode maps to U+0044 (aka U-00000044)
  - for UTF-8, encoding of these code points is done like this:

U-00000000 - U-0000007F	10xxxxxx
U-00000080 - U-000007FF	110xxxxx 10xxxxxx
U-00000800 - U-0000FFFF	1110xxxx 10xxxxxx 10xxxxxx
U-00010000 - U-0001FFFF	11110xxx 10xxxxxx 10xxxxxx 10xxxxxx
U-00200000 - U-03FFFFFF	111110xx 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx
U-04000000 - U-7FFFFFFF	1111110x 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx 10xxxxxx

See: <http://danielmiessler.com/study/encoding/>

## Mostly used: UTF-8 encoding



Picture courtesy of <http://stanford.edu/dept/itss/docs/oracle/10g/server.101/b10749/ch2charset.htm>

## Read these

- Marco Cantu
  - Whitepaper "Delphi and Unicode" New  
[http://etnaweb04.embarcadero.com/resources/technical\\_papers/Delphi-and-Unicode\\_Marco-Cantu.pdf](http://etnaweb04.embarcadero.com/resources/technical_papers/Delphi-and-Unicode_Marco-Cantu.pdf)
  - Old  
[http://www.codegear.com/article/38548/images/38548/Delphi\\_and\\_Unicode.pdf](http://www.codegear.com/article/38548/images/38548/Delphi_and_Unicode.pdf)
  - Presentation "Delphi and Unicode"  
[http://it-republik.de/konferenzen/delphi\\_live/material/DelphiLive09\\_Cantu\\_Unicode.pdf](http://it-republik.de/konferenzen/delphi_live/material/DelphiLive09_Cantu_Unicode.pdf)
- Nick Hodges
  - Whitepaper "Delphi in a Unicode World"  
[http://etnaweb04.embarcadero.com/resources/technical\\_papers/Delphi%20in%20a%20Unicode%20World\\_updated.pdf](http://etnaweb04.embarcadero.com/resources/technical_papers/Delphi%20in%20a%20Unicode%20World_updated.pdf)

## On-line help: topics

- What's New in Delphi and C++Builder 2009
  - [ms-help://embarcadero.rs2009/devcommon/whatsnewliburon.xml.html](http://ms-help://embarcadero.rs2009/devcommon/whatsnewliburon.xml.html)
  - <http://cdn.embarcadero.com/article/38869>
- String Types
  - Base: ShortString, AnsiString, WideString, UnicodeString
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/devcommon/stringtypes.xml.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/devcommon/stringtypes.xml.html)
    - (including internal memory layouts of the string types)
  - Enabling for Unicode
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/devcommon/enablingunicode.xml.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/devcommon/enablingunicode.xml.html)
  - RawByteString (AnsiString with CodePage \$ffff)
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/delphiwin32/System\\_RawByteString.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/delphiwin32/System_RawByteString.html)
  - UTF8String (AnsiString with CodePage 65001)
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/delphiwin32/System\\_UTF8String.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/delphiwin32/System_UTF8String.html)
- String <-> PChar conversions
  - PChar fundamentals
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/devwin32/stringdependencies.xml.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/devwin32/stringdependencies.xml.html)
  - Returning a PChar Local Variable
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/devwin32/passinglocalvariableasapchar.xml.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/devwin32/passinglocalvariableasapchar.xml.html)
  - Passing a Local Variable as a PChar
    - [http://docs.embarcadero.com/products/rad\\_studio/delphiAndcpp2009/Help/Update2/EN/html/devwin32/returningapcharlocalvariable.xml.html](http://docs.embarcadero.com/products/rad_studio/delphiAndcpp2009/Help/Update2/EN/html/devwin32/returningapcharlocalvariable.xml.html)

## Unchanged since D2009

- AnsiChar tkChar      PAnsiChar none      AnsiString tkLString
- WideChar tkWChar      PWideChar none      WideString tkWString
- ShortString tkString
  - ShortString and WideString are NOT reference counted
  - AnsiString (and UnicodeString) ARE reference counted
  - There is no UnicodeShortString

## Changed since D2009

- Pre D2009
  - Char = AnsiChar
  - PChar = PAnsiChar
  - String = AnsiString
  - tkString (TypeInfo)
- D2009+
  - Char = UnicodeChar
  - PChar = PUnicodeChar
  - String = UnicodeString
  - tkUString (TypeInfo)

-8	RefCount
-4	Length
0	String content

-12	Code Page
-10	Element Size
-8	RefCount
-4	Length
0	String content

## New since D2009

- AnsiString**
  - can have CodePage, i.e. AnsiString(1252) for CodePage 1252
  - CodePage is same as CodePage of the running system (NOT the compiling system)
  - RawByteString = AnsiString(\$ffff)
  - UTF8String = AnsiString(65001)
- TypeInfo** for tkLString now has CodePage field
- UnicodeString**
  - Internally uses UTF-16 (like WideString)
- Implicit conversion** from
  - AnsiString(CodePage1) to AnsiString(CodePage2)
  - AnsiString to UnicodeString
- Classes**
  - TCharacter (Characters)
  - TEncoding (SysUtils)
  - TStringBuilder (SysUtils)
  - TTextReader (Classes)
    - TStreamReader
    - TStringReader
  - TTextWriter (Classes)
    - TStreamWriter
    - TStringWriter
- Units**
  - Character
  - AnsiStrings
- CodePage list:**
  - [http://msdn.microsoft.com/en-us/library/dd317756\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/dd317756(VS.85).aspx)

## CodePages used internally

- "Code Page Identifiers"
  - [http://msdn.microsoft.com/en-us/library/dd317756\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/dd317756(VS.85).aspx)

```
// some consts that are internal to the system unit
const
  CP_UTF8 = 65001; // UTF8String CodePage
  CP_UTF16 = 1200; // string CodePage

// some literals that are internal to the system unit
const
  CP_ACP = 0; // default AnsiString CodePage
  CP_RawByteString = $ffff; // RawByteString CodePage
```

## Some things to watch for

## String constants don't always have the right CodePage

- <http://qc.embarcadero.com/wc/qcmain.aspx?d=69517>

```
const
  CP_VietnameseString = 1258;
type
  VietnameseString = type
    AnsiString(CP_VietnameseString);
// these literal string consts
// do not have the right CodePage
// they all get 1252 assigned:
const
  LF_UTF8String = UTF8String(#10);
  LF_RawByteString = RawByteString(#10);
  LF_VietnameseString = VietnameseString(#10);
```

## When converting by hand:

```
function ReadMessageBufferLen(IOHandler: TIDIOHandler;
  ABytes: Integer):
  RawByteString;
{$IF Declared(RTLVersion) and (RTLVersion >= 20.0)}
var
  Bytes: TBytes;
{$IFDEF}
begin
  {$IF Declared(RTLVersion) and (RTLVersion >= 20.0)}
  IOHandler.ReadBytes(Bytes, ABytes);
  SetLength(Result, Length(Bytes));
  Move(Bytes[0], Result[1], ABytes);
{$ELSE}
  Result := IOHandler.ReadString(ABytes, Indy8BitEncoding);
{$IFDEF}
end;
```

## Detecting wrong encodings

- Bad**

This looks OK, but is in fact not UTF-8, but an ISO-8859-1 character
- Good**

This looks bad, but is in fact the correct UTF-8 encoding for "à"

Microsoft  
.net

## Reencoding files

- Two examples:
  - Utf82Ascii
    - Convert from UTF-8 to ASCII
  - Iso8859\_12Utf8
    - Convert from ISO 8859-1 to UTF-8
- Shows:
  - TEncoding (default and custom)
  - TStreamReader/TStreamWriter
  - THandleStream

better  
office

CODE  
GEAR

better  
office

## Q & A

Jeroen Pluimers  
better office benelux  
[jpluimers@better-office.com](mailto:jpluimers@better-office.com)

If you have questions after the session, please mail me

Downloads will be on my blog: <http://wiert.wordpress.com>

Microsoft  
.net

CODE  
GEAR