

Unicode in Education



Adil Allawi
Technical Director
adil@diwan.com

Why Education and Unicode:

- Cost of printing and distribution is high - paper is becoming a scarce resource
- Many classes only need part of a book - less waste
- More than any other kind of book - educational books need to be reused for different media - print/web/ebook etc
- Encoding must be right or text is lost for the future
- What else do you put on “one laptop per child” ?

So why not use PDF's?

- The good
 - a lot of the documents already exist in electronic form that prints to Postscript or can be scanned and saved as PDF
- The bad
 - limited to the feature set of the browser/library
 - royalties drive up the end user costs - and for education they have to be low
 - it sidesteps the real issue
 - Without properly handling the encoding and structure you end up with just an image of the document
 - text flow is lost, links between items are lost.



Encoding and rendering

- Take away everything and we have a story of encodings
- Books contain handwritten examples.
- Fonts created were built just to render the book in print
- Need to convert simple and practical encodings
- Trouble with Bidi Algorithms

Simple encodings

- Mostly easy to encode into Unicode
- No new ones being created
- For Arabic there are a few:
 - 2 Mac Arabic ones (System 6 and 7)
 - Windows Arabic
 - ISO-8859-6
 - DOS Arabic
 - 3B2 (an odd 16-bit encoding)

Practical encodings

- They work around system/application restrictions
- Font and Keyboard based
- Difficult to re-encode
- New ones are being created all the time
 - Unicode has only made this problem more complex

Practical Problems

- Based on a well-known DTP software

ض ش أ آ

- From a well-known institution

FB53	FB54	FB55	FB56	FB57	FB58	FB59	FB5A	FB5B	FB5C	FB5D	FB5E	FB5F	FB60	FB61	FB62	FB63	FB64	FB65
وَلَمْ	نَكَ	نُطْعِمُ	مِسْكِي	٤٤	نُحُوضُكَ	مَعَ	فَايَضِي	٤٥	وَكُنَّا	يَوْمَ	نُكَذِّبُ	الَّذِينَ	٤٦	حَتَّى	أَتَيْنَا	الْيَقِينَ	٤٧	

- Arabic 8-bit $\text{ى} = \text{U}+06\text{CC}$ or $\text{U}+0649$
 - In Mac Arabic it is the same for Alef Maksura or Egyptian-Arabic Yeh or Farsi Yeh.

Unicode Problems

- Many, many encoders use Arabic Presentation Forms A and B for encodings and placing characters
- is ١١/١٢ a date or a fraction?
- also ك ك are essentially the same letter
- so are ي and ي
- But the link is not made

Example - A Math Font

$$\frac{\text{قام}}{\text{ك}} = \text{لام}$$

$$\frac{\text{ع}}{\text{نو}} = \text{لام}$$

0021	0022	0023	0024	0025	0026	0027	0028	0029	002A	002B
%	∴	')	(×	+	∅	≠	≥	+
002C	002D	002E	002F	0030	0031	0032	0033	0034	0035	0036
∏	-	√	/	•	١	٢	٣	٤	٥	٦
0037	0038	0039	003A	003B	003C	003D	003E	003F	0040	0041
√	∧	٩	<	٥	√	=	آ	↔	*	A
0058	0059	005A	005B	005C	005D	005E	005F	0060	0061	0062
×	∩	∅	≠	•	-	↔	√	∃	↔	↔
0626	0627	0628	0629	062A	062B	062C	062D	062E	062F	0630
٧	٢	٣	٤	٥	٦	٧	٨	٩	١٠	١١
0631	0632	0633	0634	0635	0636	00D7	0637	0638	0639	063A
⊥	ر	س	○	ص	ص	×	ط		ع	ذ
0640	0641	0642	0643	00E0	0644	00E2	0645	0646	0647	0648
÷	و	ق	ك	à	ل	â	م	د	ه	و



Encoding Math Fonts

- Keep the practical encoding for data entry
 - simple - works with existing programs
- re-encode in the file format
 - some glyphs will encode as 2 characters
e.g. initial Ha $\curvearrowright = \text{ح} + \text{zwj}$
 - some need separate markup
e.g. subscript 2 `_٢`
- Try to use existing encoding
Curly Arabic - just matches to normal Arabic

■ But what about






Encoding Quran in Unicode

- Announced at the 17th IUC
- Extended to teaching Quran

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

أَقْرَأْ بِاسْمِ رَبِّكَ الَّذِي خَلَقَ ﴿١﴾ خَلَقَ الْإِنْسَانَ مِنْ عَلَقٍ ﴿٢﴾ أَقْرَأْ وَرَبُّكَ

الْأَكْرَمُ ﴿٣﴾ الَّذِي عَلَّمَ بِالْقَلَمِ ﴿٤﴾ عَلَّمَ الْإِنْسَانَ مَا لَمْ يَعْلَمْ ﴿٥﴾ كَلَّا إِنَّ



Encoding Method

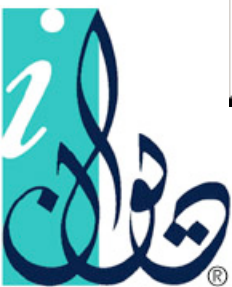
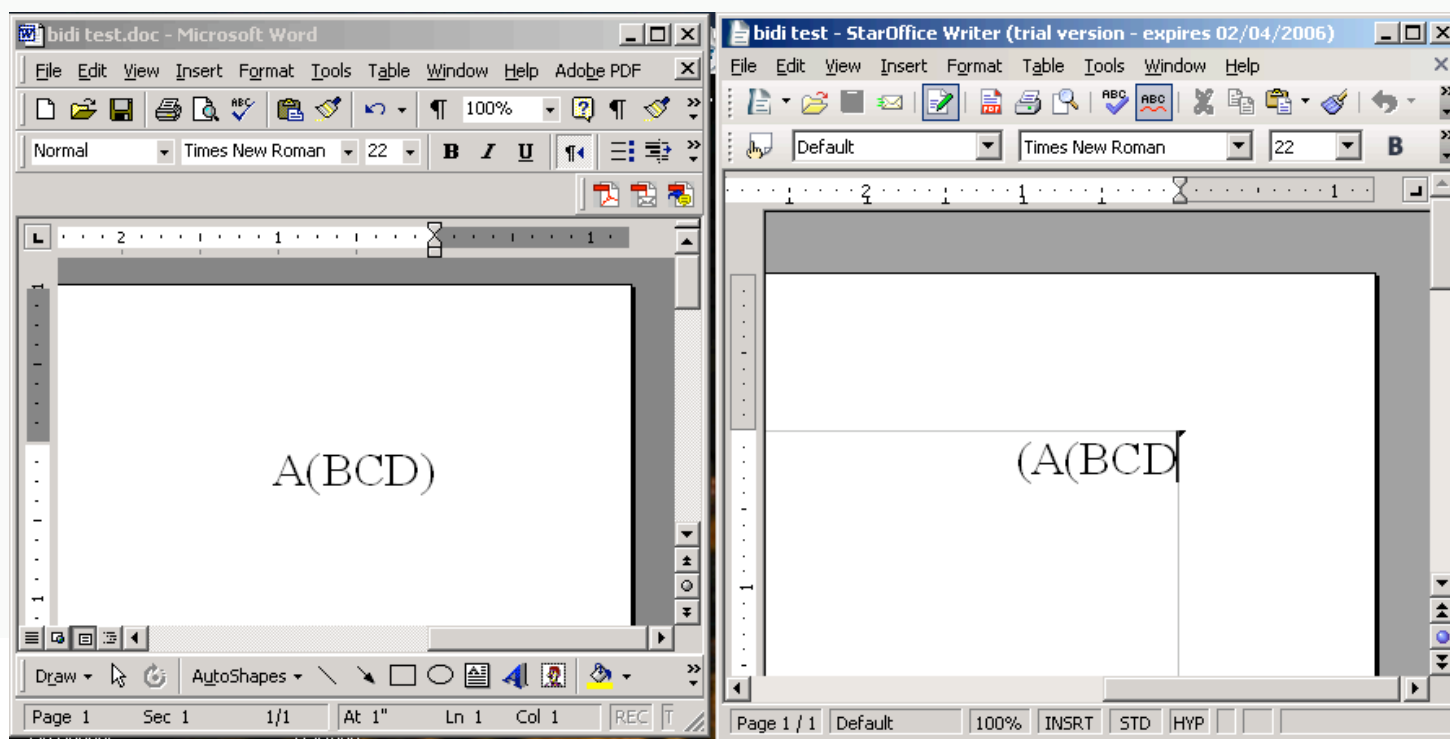
- Fitted encoding within the Unicode standard
- Example: Tanween $\overset{=}{\text{---}}\overset{2\text{e}}{\text{---}}$ is one character but 3 display forms:

عَظِيمٌ عَظِيمٌ عَظِيمٌ

عَظِيمٌ عَظِيمٌ عَظِيمٌ

De-bidi the Bidi Algorithm

- Typical problem
Word uses its own bidi algorithm



De-bidi the bidi algorithm

- There are 4 bi-di algorithms
 - Apple Mac OS 9
 - Windows 95
 - Unicode
 - Word for Windows
- And then there is text after printing
- Need to take apart bi-di'd text and reencode/reorder to work with Unicode Bidi.

Same issue different language

- Iñupiaq

From Alaska Native Education Program:

- 004C + 0323 ᐅ
 - 006C + 0200A +0323 ᐅ
 - 0141 + 0323 ᐅ
 - 0141 + 0200A + 0323 ᐅ
- Thin space is inserted because of rendering issues on web browsers

Solutions

- If the rendering engine does not fit your requirements - create your own.
- Use local knowledge to define the encoding and get it right first time. Don't let the programmers create workarounds.
- If you really must bypass Opentype or Uniscribe or browser text renderers - create something to convert back to Unicode or all the effort will be lost in time.

Looking to the Future

- Unicode is important for education and for preserving knowledge. Yet very little effort is put into encoding for education outside the major languages.
- We are creating many more encoding issues for the future as people are 'fixing' problems with 'clever' workarounds. Unicode is not stopping this process.
- How will all this be indexed and searched?
- One day the memory will be lost for a particular workaround and the data and research will be damaged for the future.