

# Unicode character properties

## Unicode character properties: what are they?

- Attributes of Unicode characters specified by Unicode standard
- Determine behavior of each character (e.g. joining)
- Must be considered during orthography development
  - (Is my proposed use of the character compatible with its properties?)
- When proposing a new character, all relevant property values must be included in the proposal.

This will be a brief review of some of the character properties that are important for Arabic characters and which will need to be defined for any proposed new characters.

## Unicode character properties: Where to find them

- The Unicode Character Database (UCD) data files
  - Plain-text files (but with strict formatting rules)
    - Line-oriented and structured XML variants are available
  - <http://www.unicode.org/Public/UCD/latest/>
- Additional ways to view properties
  - Richard Ishida's UniView webapp: <https://r12a.github.io/uniview/>
  - Unicode properties Excel workbook: <http://scripts.sil.org/excelunicodedata>
    - Currently at v 7.0
  - Unicode's Unibook windows app: <http://unicode.org/unibook/>
  - 
  - 
  -

## Unicode properties: Joining type and group (ArabicShaping.txt)

- Definitions
  - **R**=Right Joining, **L**=Left Joining, **D**=Dual Joining, **C** Join\_Causing, **U**=Non Joining
  - Joining Groups: AIN, **ALEF**, **BEH**, DAL, FARSI YEH, FEH, GAF, HAH, HEH, HEH GOAL, KAF, KNOTTED HEH, LAM, **MEEM**, **No\_Joining\_Group**, NOON, NYA, QAF, REH, SAD, **SEEN**, SWASH KAF, TAH, TEH MARBUTA, TEH MARBUTA GOAL, WAW, YEH, YEH BARREE, YEH WITH TAIL
- Examples:
  - 0622; ALEF WITH MADDA ABOVE; **R**; **ALEF**
  - 0634; SEEN WITH 3 DOTS ABOVE; **D**; **SEEN**
  - 06DD; ARABIC END OF AYAH; **U**; **No\_Joining\_Group**
  - 08A1; BEH WITH HAMZA ABOVE; **D**; **BEH**
  - 08A7; MEEM WITH 3 DOTS ABOVE; **D**; **MEEM**
  - 200D; ZERO WIDTH JOINER; **C**; **No\_Joining\_Group**

2017 ©SIL International

Unicode prescribes every character's "Joining type" and "Joining group". These are all defined in a file called "ArabicShaping.txt". It's important to understand this so that you will know the expected behavior for a characters.

The Joining type describes how a character joins. An example of Right joining is an Alef. It only joins on the right. And example of a Dual joining character is a seen or a beh. As you know, these join on the left or the right. The only Join Causing character that we are likely to need is U+200D. A number of Arabic characters are in the Non Joining type. An example of this would be the End of Ayah. This doesn't mean it doesn't require special rendering, it just means it doesn't actually link with another character.

Next we have Joining Groups. This defines what kind of joining behavior to expect. If a character is in the BEH group, then we would expect it to behave as the standard "Beh" behaves.

## Unicode properties: Combining class

27: fathatan, open fathatan (◌َ◌َ)

32: kasra, small kasra (◌ِ◌ِ)

28: dammatan, open dammatan (◌ُ◌ُ)

33: shadda (◌ّ)

29: kasratan, open kasratan (◌ِ◌ِ)

34: sukun (◌◌)

30: fatha, small fatha (◌◌)

35: superscript alef (◌◌)

31: damma, small damma (◌◌)

220: all other below combining marks

230: all other above combining marks

2017 © SIL International

The combining class value determines how combining marks are reordered in your text during Unicode normalization (most often used for string storage and comparison).

Can anyone suggest why this question is important during orthography design?

Ans: Suppose you have diacritics from two different classes, say a shadda (cc=33) and damma (cc=31), on the same base letter. Text processes are free to change the relative order of the diacritics. In particular, Unicode normalization will always put the damma before the shadda no matter what order they were in the original text. This means – from an orthography design perspective – that one *must not make meaning distinctions* between, damma followed by shadda and shadda followed by damma -- searching, sorting, comparison should all treat either order as meaning the same thing. (Also, font designers should render damma+shadda identically to shadda+damma, but not all fonts do this correctly).

Does anyone see anything illogical in these combining class values?

Ans: Do you think it is logical for a fatha or damma to precede shadda?

## Unicode properties: Combining class (UnicodeData.txt)

- Vowels

- 064E;ARABIC FATHA;Mn;30;NSM;;;;;N;ARABIC FATHAH;;;;
- 064F;ARABIC DAMMA;Mn;31;NSM;;;;;N;ARABIC DAMMAH;;;;
- 0650;ARABIC KASRA;Mn;32;NSM;;;;;N;ARABIC KASRAH;;;;
- 08F5;ARABIC FATHA WITH DOT ABOVE;Mn;230;NSM;;;;;N;;;;;
- 08F6;ARABIC KASRA WITH DOT BELOW;Mn;220;NSM;;;;;N;;;;;

- Shadda, sukun, hamza

- 0651;ARABIC SHADDA;Mn;33;NSM;;;;;N;ARABIC SHADDAH;;;;
- 0652;ARABIC SUKUN;Mn;34;NSM;;;;;N;;;;;
- 0654;ARABIC HAMZA ABOVE;Mn;230;NSM;;;;;N;;;;;

- Koranic marks

- 06DC;ARABIC SMALL HIGH SEEN;Mn;230;NSM;;;;;N;;;;;
- 06EA;ARABIC EMPTY CENTRE LOW STOP;Mn;220;NSM;;;;;N;;;;;

Extending this to more marks: if you have a shadda (cc=33), fatha (cc=30) and a koranic mark above (cc=230), from a technical point of view, these mean the same thing no matter which of the 6 possible orders they are in your text!

If there are multiple combining marks *from the same class*, say from cc=220, Unicode normalization will not change the relative order of these marks in your text. While this potentially means that you can treat different orders as having different meanings, extreme caution is needed as there are other rules at play (such as for the koranic and tone marks)



## Unicode properties: Decomposition

- Unicode character properties (UnicodeData.txt)
  - 0623;ARABIC LETTER ALEF WITH HAMZA ABOVE;Lo;0;AL;0627 0654;;;;;N;ARABIC LETTER HAMZAH ON ALEF;;;;
  - 08A1;ARABIC LETTER BEH WITH HAMZA ABOVE;Lo;0;AL;;;;;N;;;;;
- Normalization issues (never decompose)
- Collation (where does the character sort?)
  - U+08A1 could sort at end of beh-like characters
- Confusability
  - Can this character be confused with other characters?
  - U+08A1 could be confused with U+0628 + U+0654. This has implications for internet domain names and spoofing
  - Document that there should not be a decomposition to U+0628 U+0654

2017 © SIL International

Another important Unicode character property is Decomposition. An example is 0623 which is an alef with a Hamza above. Unicode character properties declare that this character *is the same thing as* the two character sequence alef plus combining Hamza. From the perspective of what does the text mean, there is no difference between 0623 and the sequence 0627+0654.

This is common in Latin script, right? Someone give me an example?

Recalling the earlier example (in the previous presentation) of the new character beh+Hamza – note that the decomposition field is empty indicating there is no decomposition.

Ok, does decomposition need to be taken into account when making orthography decisions?

[discuss]

In the Unicode proposal a number of questions have to be answered, including:

- impact on normalization (we already noted there is no decomposition so it has no impact on normalization)
- Collation: we have to specify a default sort order for this character relevant to

- existing characters
- Confusibility



## Unicode properties: Decompositions

Does this have a bearing on orthography decisions?

## Q & A

2017 ©SIL International