

**Subject:** Dependency Graph for Segmentation  
**From:** Mark Davis (Google, Inc.), Laurențiu Iancu (Microsoft Corporation)  
**Date:** 2016-11-01

We had an action of identify those other Unicode Properties that the segmentation algorithms explicitly depend on. That is, a change to those other properties can cause the results of segmentation to change. In some cases, the results only depend on a few particular values of the property, not all or most of the values.

148	A001	Mark Davis, Laurentiu Iancu	Create a dependency graph for derived properties where the derivation is either explicit or implicit, for the November 2016 UTC meeting. The initial focus is on the four segmentation properties.
-----	------	--------------------------------	--

There are other (implicit) dependencies, but the explicit ones are the focus of this doc. (An example of an implicit dependency is the GCB=SM. The excluded marks are gc=Mc from SE Asian scripts (Myanmar, Tai Tham etc.) which are post-base vowels or tone marks (approx. InSC=Vowel\_Dependent, InPC=Right), but are ultimately hand-tuned based on feedback from experts.)

[Grapheme Cluster Break](#)

[Word Break](#)

[Sentence Break](#)

[Line Break](#)

## Grapheme Cluster Break

1. General\_Category=Unassigned, Control, Format, Surrogate, Line\_Separator, Paragraph\_Separator, Spacing\_Mark, Nonspacing\_Mark, Enclosing\_Mark
2. Default\_Ignorable\_Code\_Point
3. Grapheme\_Extend
4. Regional\_Indicator\*
5. Hangul\_Syllable\_Type
6. Indic\_Syllabic\_Category = Consonant\_Preceding\_Repha, Consonant\_Prefixed
7. Prependen\_Concatenation\_Mark
8. Emoji\_Modifier\_Base
9. Emoji\_Modifer
10. Glue\_After\_Zwj\*
11. E\_Base\_GAZ\*

\* The items marked with a star don't currently have UCD properties. There are, however, proposals to have UCD properties for them.

## Word Break

1. General\_Category=Other\_Letter, Spacing\_Mark, Format, Connector\_Punctuation
2. Script=Hebrew, Katakana, Hiragana
3. Grapheme\_Extend
4. Regional\_Indicator\*
5. Line\_Break=Complex\_Context, Infix\_Numeric, Numeric
6. Alphabetic
7. Ideographic
8. Emoji\_Modifier\_Base
9. Emoji\_Modifer
10. Glue\_After\_Zwj\*
11. E\_Base\_GAZ\*

## Sentence Break

1. Grapheme\_Extend
2. General\_Category=Spacing\_Mark, Format, Titlecase\_Letter, Open\_Punctuation, Close\_Punctuation
3. White\_Space
4. Lowercase
5. Uppercase
6. Alphabetic
7. Line\_Break=Numeric, Quotation
8. Sentence\_Terminal

## Line Break

Line\_Break is a primary property whose values are not expressed in terms of other properties. They are assigned by some unwritten heuristics, but not formally. We do have that, for instance, lb=EM is equivalent with Emoji\_Modifier=Yes (because it would make no sense otherwise), but there is no formal dependency.

The original assignments had started from some rules, but by now we have probably accumulated irregularities or deviations from those rules. For example, see <http://www.unicode.org/L2/L1999/99179.pdf>

That being said, were we to have explicit dependencies (with perhaps exceptions), they would probably be the following:

1. General\_Category
2. Hangul\_Syllable\_Type
3. Emoji\_Modifier\_Base
4. Emoji\_Modifier
5. Regional\_Indicator
6. East\_Asian\_Width