

**From:** Peter Edberg, Mark Davis, Andy Heninger  
**Subject:** Segmentation & Linebreak  
**Date:** 2011-05-09

We have the following action:

125	A099	Andy Heninger, Peter Edberg, Mark Davis		Create a WD of a PU UAX #14 which addresses the general issue of breaks after dashes.
-----	------	---	--	---

Here is a breakdown of the issues and our recommendations.

**1. Hebrew linebreak:** With <hebrew hyphen non-hebrew>, there is no break on either side of the hyphen.

We recommend that this be done in the following way:

Split AL → (AL | HL) with the following redefinitions

HL to be the current AL ∩ [:script=hebrew:]

AL to be the current AL ∖ [:script=hebrew:]

Add new rule:

*Don't break after Hebrew + Hyphen*

LB21a HL (HY|BA) ×

**Issue:** BA includes the following.

U+00AD ( ) SOFT HYPHEN

U+058A ( ) ARMENIAN HYPHEN

U+1400 ( ) CANADIAN SYLLABICS HYPHEN

U+2010 ( - ) HYPHEN

U+2027 ( · ) HYPHENATION POINT

U+2E17 ( ) DOUBLE OBLIQUE HYPHEN

The reason we need BA is to get Hyphen (2010). It doesn't hurt to include Armenian or CA hyphen, or the double oblique. The question is whether we need to split BA in two, because of either Soft Hyphen or Hyphenation point. The same is true for #3 below.

**2. Hebrew word break:** with <hebrew quote hebrew> there is no break on either side of the quote (single or double)

The recommendation is to add “ (0022) to MidLetter, so that it behaves like Gershayim and Apostrophe, because it is often used for that. This shouldn't cause any problems, because SA and Ideographs are not linked by it. So these would behave the same:

U+05F4 ( " ) HEBREW PUNCTUATION GERSHAYIM

U+0022 ( " ) QUOTATION MARK

**3. Finnish linebreak (and others):** with <letter space hyphen letter>, we need to allow a break before the hyphen but not after the hyphen.

Currently, we break as follows. The → shows the proposed change. We believe that this wouldn't disturb other usage.

1. a - b *break before and after*
2. a -b *break before and after* → *break before*
3. a- b ***break after***
4. a-b *break after*

Proposed Rule

LB21b

SP (HY|BA) × !SP

**Issue:** we know this will not be straightforward to implement as-is in ICU. We could approximate it with ÷ (HY|BA) × !SP (that is, don't break after if there was a break before).

**4. Spanish (and others):** with <letter space emdash letter>, don't break between emdash & letter

The issue is that some languages use emdashes to set off a parenthetical, and you don't want to break the surrounding ones from the contained text. In that usage, there is a space on the side where it can be broken. This doesn't conflict with symmetrical usages (spaces either before or after).

Currently, we break as follows. The → shows the proposed change.

1. a — b *break before and after*
2. a —b *break before and after* → ***break before, but not after***
3. a— b *break before and after* → ***break after, but not before***
4. a—b *break before and after*

Proposed Rule:

LB21c

SP B2 × !SP  
!SP × B2 !SP

## 5. Resubmit L2/09-263

**6. Current Kinsoku too restrictive.** Submit changes from #3571 to make the current rules “normal”.

Proposal: Move normal (not halfwidth) small kana and prolonged sound mark from linebreak class NS (non-starter) to ID (ideographic), like other kana. This applies to the following characters:

Small hiragana: [3041 3043 3045 3047 3049 3063 3083 3085 3087 308E 3095 3096]

Small katakana: [30A1 30A3 30A5 30A7 30A9 30C3 30E3 30E5 30E7 30EE 30F5 30F6 31F0-31FF]

KATAKANA-HIRAGANA PROLONGED SOUND MARK: 30FC

Rationale: CSS Text Level 3 (which supports Japanese line layout) defines three distinct values for its line-break behavior:

- strict (corresponds to current UAX #14 behavior), typically used for long lines.
- normal (CSS default), the behavior typically used for books and documents.
- loose, typically used for short lines such as in newspapers.

These have different sets of “kinsoku” characters which cannot be at the beginning or end of a line; strict has the largest set, while loose has the smallest. The motivation for the smaller number of kinsoku characters is to avoid triggering justification that puts characters off the grid position.

The default UAX #14 behavior should be changed to align more closely with the CSS “normal” behavior. The changed behavior is described in the first reference below, as follows:

“Following breaks be forbidden in ‘strict’ line breaking and allowed in ‘normal’:

- breaks before Japanese small kana
- breaks before the KATAKANA-HIRAGANA PROLONGED SOUND MARK (U+30FC)

“Additionally, if the language is known to be Chinese or Japanese, breaks before hyphens (U+2010, U+2013, U+301C, U+30A0) may be allowed in ‘normal’.”

Since UAX #14 cannot make assumptions about text language, it should make just the change corresponding to the bulleted items (the other change is optional anyway). This can be accomplished by moving those characters from class NS to ID.

References:

- <http://www.w3.org/TR/css3-text/#line-break>
- [http://www.w3.org/TR/jlreq/#en-subheading2\\_1\\_7](http://www.w3.org/TR/jlreq/#en-subheading2_1_7)

**7. Issue found while doing this document:** LB indicates that the following are punctuation, but they are Alphabetic. This seems incorrect.

#### **Line\_Break=Close\_Punctuation**

U+1325B ( ) EGYPTIAN HIEROGLYPH Ooo6D  
...{1}...U+1325D ( ) EGYPTIAN HIEROGLYPH Ooo6F  
U+13282 ( ) EGYPTIAN HIEROGLYPH Oo33A  
U+13287 ( ) EGYPTIAN HIEROGLYPH Oo36B  
U+13289 ( ) EGYPTIAN HIEROGLYPH Oo36D  
U+1337A ( ) EGYPTIAN HIEROGLYPH Vo11B  
U+1337B ( ) EGYPTIAN HIEROGLYPH Vo11C

#### **Line\_Break=Open\_Punctuation**

U+13258 ( ) EGYPTIAN HIEROGLYPH Ooo6A  
...{1}...U+1325A ( ) EGYPTIAN HIEROGLYPH Ooo6C  
U+13286 ( ) EGYPTIAN HIEROGLYPH Oo36A  
U+13288 ( ) EGYPTIAN HIEROGLYPH Oo36C  
U+13379 ( ) EGYPTIAN HIEROGLYPH Vo11A