

**Title:** Regularize Extended\_Pictographic  
**From:** Mark Davis, CLDR-TC  
**To:** UTC  
**Date:** 2017-03-22  
**Draft:** [link](#)

We created extended segmentation rules in CLDR to handle emoji when we didn't have time to make the right changes in UAX#14 and UAX#29. That means that anyone who wants to properly handle Emoji in segmentation must use the CLDR additions: customized LDML rules and the additional CLDR data in Extended\_Pictographic.txt. Since we'd like everyone to be future-proofed, that effectively means that everyone should ignore the plain rules and use the CLDR ones (which also go into ICU).

Once E5.0 (UTS #51 Emoji 5.0) is approved, we can fix that situation in subsequent versions of the Unicode Standard and CLDR. This is probably too substantial for people to handle in U10.0, but we can have a clear direction for the future. *However, it would be worth putting out a PRI early to get feedback on this.*

## Proposal

### UTS #51

Move **Extended\_Pictographic** from CLDR into the Emoji data files, for the next version of UTS #51 after Emoji 5.0 (Emoji 6.0 or perhaps a sooner small update Emoji 5.1, whatever timing is needed).

That is, add data for new Emoji property with the following contents:

```
Extended_Pictographic =
  Extended_Pictographic.txt (from CLDR)
  + [:Emoji:]
  - [:Emoji_Component:]
```

- Note that for each new version of E5.0, we'd need to review any recently-added Unicode symbols to see whether they should also be added.

### CLDR

Deprecate **Extended\_Pictographic.txt** and remove the data (leave a readme in the data file).

In LDML, point to the emoji data for the replacement values for Extended\_Pictographic.txt, and remove the customized rules that override UAX#14 and UAX#29.

### [UAX#14](#) and [UAX#29](#)

Add text along the lines of the following:

The rules for segmentation may may be augmented by use of UnicodeSet notation using properties or literals outside of the main property associated with the algorithm, such as `[:General_Category=Letter:]` or `[\uFE0F]`. The properties may either be UCD character properties, or properties associated with a UTS such as *UTS #51, Unicode Emoji*.

Modify the segmentation rules in UAX#14 and UAX#29 based on LDML (updated somewhat). The **old rules** are presented for comparison, while the new rules are presented with a prime' mark.

### Grapheme Cluster Break

*GB10*    *(E\_Base | EBG) Extend\**    ×    *E\_Modifier*

