## Shaping Ideographic Description Sequences (v1) — TSWG 2020-01

Fredrick R. Brennan < copypaste@kittens.ph > July 21, 2020

I have one main desire for text shaping: the creation of a standard for an OpenType table that implements a Chinese character description language. Despite the excellent and wonderful work done by Ken Lunde, Qin Lu, and many, many others over the years, it is still not possible to type an ideograph not in Unicode until it gets in Unicode, which is a long process for most ideographs. You might think, surely no popular ideograph is left out. Actually, many popular ideographs remain unencoded. Indeed, an ideograph so popular it has a lengthy description of itself on Wikipedia (not Wiktionary), □ (U+030EDE; Pinyin: *biáng*), has only been officially part of Unicode since 13.0. This year! It was only proposed in 2015, in L2/15-223.¹ I have been lamenting this in public since at least October 2018.²

Right now, thanks to <u>a long series of proposals</u> to the ISO/IEC JTC 1/SC 2 and UTC beginning with one from China on November 7, 1995, X3L2/95-111,<sup>3</sup> "Ideographic Structure Symbol (additional request)", Ideographic Description Characters (IDC's) exist.

We can see from an October 22, 1998 proposal, <u>SC2 N3186</u>, from Japan's liaison to SC2, "The IDS describes the ideograph in the abstract form. It is not interpreted as a composed character and does not have any rendering implication." However, this has changed. <u>Chapter 18</u>, Section 2, page 735<sup>4</sup> of The Unicode Standard v13.0 reads in part: "An implementation may render a valid Ideographic Description Sequence either by rendering the individual characters separately **or by parsing the Ideographic Description Sequence** 

1

<sup>&</sup>lt;sup>1</sup> This document is for CJK Unified Ideographs Extension H, but U+030EDE ended up in Extension G.

<sup>&</sup>lt;sup>2</sup> Brennan, Fredrick (2018). "Answer to 'What are the disadvantages of typography?'". Quora.

<sup>&</sup>lt;sup>3</sup> This proposal seems unfortunately lost to time.

<sup>&</sup>lt;sup>4</sup> In linked PDF, page 28.

and drawing the ideograph so described." [emphasis mine] There was a lot of precedent, for this, though. Indeed, John Jenkins was already doing it in 1999.<sup>5</sup>

Indeed, no less than Dr. Lunde himself has released fonts which do this, first being one just for <i>biáng</i> on March 24, 2014. <sup>6</sup> That font replaces the IDC
sequence □辶⊟穴□月□□□□幺長⊟言馬⊟幺長刂心 with <i>biáng</i> . However, this is very much an "old-style" GSUB simple substitution, stuck in ccmp, and not what l
propose. I propose extending OpenType shaping to make a best effort on arbitrary IDC's. Of course, fonts would need to opt in, and as the Unicode
Standard says, encodings are always preferred to IDC's. Indeed, human designers adding ccmp's are preferred, but unrealistic given an unlimited number
of IDC's exist. I will save more details for later if requested, but here's a sketch.
The shaper needs to do three main things, using □□口止□久口 (路) as an example:

Set up necessary glyph classes for all IDC's. Essentially, one class for the eleven IDC "functions", one mega class of all "normal" IDC arguments the font supports containing the variation that takes up the full ideograph. (口, 止, …) Then, mega classes containing the same number of glyphs for all IDC argument types. So, for 口:

	o argument types. 66, for $\mathbf{H}$ .	
0	□_Plain.	
0	Ш: □_Left; □_Right.	
0	⊟: □_Above; □_Below.	
0	Ш: □_LeftOf3; □_HMiddleOf3; □_RightOf3.	
0	⊟: □_TopOf3; □_VMiddleOf3; □_BottomOf3.	
0	□: □_Surrounding; □_Surrounded.	
0	□: □_SurroundingFromAbove; □_SurroundedFromAbove.	
0	□: □_SurroundingFromBelow; □_SurroundedFromBelow.	
0	□: □_SurroundingFromLeft; □_SurroundedFromLeft.	
0	□: □_SurroundingFromUpperLeft; □_SurroundedFromUpperLeft.	
0	☐: □_SurroundingFromUpperRight; □	
	_SurroundedFromUpperRight.	
0	□: □_SurroundingFromLowerLeft; □_SurroundedFromLowerLeft.	

<sup>&</sup>lt;sup>5</sup> John H. Jenkins (1999). "New Ideographs in Unicode 3.0 and Beyond". San Jose, CA: 15th International Unicode Conference. pp. 12–21.

<sup>&</sup>lt;sup>6</sup> Lunde, Ken (2014). "IDS + OpenType: Pseudo-encoding Unencoded Glyphs". Adobe, CJK Type Blog.

- □: □\_OverlaidA; □\_OverlaidB.
  Do necessary substitutions:
   sub □ @\_Plain' by @\_Above;
   sub @\_Above @\_Plain' by @\_Below;
   sub □ □ @\_Above' @\_Below' by @\_AboveInTop @\_AboveInBottom;
   ... and so on.
   It's immediately clear that there's some arbitrary rehere, as we must define classes for all types, and it
  - It's immediately clear that there's some arbitrary recursion limit here, as we must define classes for all types, and in any event, make glyphs for them. This is where we must extend the world of normal OpenType for best results.
- On □, defined with width 0, and with its two anchors, "L" and "R", attach □ \_\_InLeft to "L" and □\_InRight to "R". On □\_InLeft, with its two anchors, "T" and "B", attach □\_TopInLeft to "T" and 止\_BottomInLeft to "B". On □ \_\_InRight, attach 夂\_TopInRight and □\_BottomInRight. Voilà, 路.

It would be nice if we could determine the start and end of a discrete IDC describing one character as our font should be able to handle rows of them.  $\square$  is essentially a "function" with two "position arguments". So, in C notation, do  $\{\Box (\Box (\Box, \bot), \Box (\nabla, \Box))\}$  while (0).

It would also be nice if we didn't need so many classes. My friend Simon Cozens wrote a program called <u>fontFeatures</u> which implements a file format analogous to, but better than, Adobe Feature Files (.fea), called "FEE". This gets us much of the way there, at least in terms of removing the tedium of writing these deeply nested class definitions like

@\_LeftInRightInBottomInBottomOf3InOverlaidA. \*! Let's not forget our good friend biáng. In an OpenType font that could shape arbitrary IDS's, assuming it had no special glyph for biáng, (presumably, many such fonts would have special glyphs for common IDS's that look better,) and assuming that we change nothing about the standard to help, the glyph sequence would need to be so long I put it in its own section. (§ Biáng as input to a shaper)

<sup>&</sup>lt;sup>7</sup> In case you forgot, a do {} while (cond) loop is guaranteed to run at least once regardless of while condition.

<sup>&</sup>lt;sup>8</sup> Font Engineering (with) Extensibility

It would be nice if we could tell Fontconfig, CoreText, etc., what runs our font can actually handle, perhaps via a new OpenType name ID, (in the form of, I don't know, a PECL compatible regexp? Some kind of formal grammar? Comments welcome.), so we get skipped even if we have all the glyphs for a run but we can't handle it. We don't want Fontconfig to pass HarfBuzz a run our font can't handle. For example, imagine a simple version of the font with everything implemented up to a maximum recursion depth of 2. Fontconfig will just see we have glyphs for  $\square$ ,  $\boxminus$ , et cetera, and all the radicals, and assume we should handle the run. So, we will spit up junk back to the user when really we should've just been skipped.

It would be *very* nice if OpenType had an ability to do simple linear scales of referenced glyphs. If I'm not mistaken, PostScript Type 3 could. Now, of course, I'm *not* saying most classes should be implemented via scales. They shouldn't! However, we'd need far fewer glyphs if we could do even the most simple of rectangular scales of even up to ±20%. But, my ask is really, if we're dreaming big dreams: *allow us to define anchor classes of a glyph at a certain point along their linear interpolation gradient*. Imagine a variable font with an axis for each of the twenty-six basic ways to set any given ideograph with IDS's. That is to say, when I add an anchor to the base glyph, besides position, I can also set any or all of its variable font axes.

## Biáng as input to an arbitrary IDS shaper

Note: This is just an example. Per the Unicode Consortium, IDS's should not be used for existing characters in documents once they have encodings.



**Expected output:** 

Input IDS:9 □辶□穴Ⅲ月Ⅲ□幺長□言馬□幺長刂心

## List of glyphs after OpenType shaping:

- ☐ i SurroundingFromLowerLeft
  - $\exists$   $\dot{\pi}_{\text{TopOf3InSurroundedFromLowerLeft}}$ 

    - - 長 BelowInLeftOf3InVMiddleOf3InHMiddleOf3InSurroundedFromLowerLeft
      - 日言 AboveInVMiddleOf3InVMiddleOf3InHMiddleOf3InSurroundedFromLowerLeft
        - \_\_\_\_\_ 馬 BelowInVMiddleOf3InVMiddleOf3InHMiddleOf3InSurroundedFromLowerLeft
      - $\exists$  & AboveInRightOf3InVMiddleOf3InHMiddleOf3InSurroundedFromLowerLeft
        - ${\bf \xi}\_{\tt BelowInRightOf3InVM} iddleOf3InHM iddleOf3InSurroundedFromLowerLeft$
      - $\hbox{\tt II} \hbox{\tt RightOf3InHMiddleOf3InSurroundedFromLowerLeft}$

<sup>ightharpoonup</sup> BottomOf3InSurroundedFromLowerLeft

<sup>&</sup>lt;sup>9</sup> A normalization step is required. Details TBD.