

Unicode font challenges in linguistics

Berthold Crysman

Laboratoire de linguistique formelle

`crysmann@linguist.jussieu.fr`

1 Scripts

In linguistics, we need support for various scripts, ideally covering the whole of unicode. While scripts from different traditions may follow entirely different design principles (compare Latin/Greek/Cyrillic vs. Arabic vs. CJK) and can therefore be recruited from fonts with different design without too much trouble, the design of scripts within a script family should match.

Thus, for Latin base fonts, we want matching, or better identical glyph designs for Greek and Cyrillic, and, of course the IPA.

The following 2 fonts appear to satisfy these conditions:

- Linux Libertine (French renaissance antiqua)
- FreeSerif (=Times; Baroque antiqua)

Charis SIL has very good unicode coverage for Latin and Cyrillic, including sophisticated accent placement, but, unfortunately no Greek (in text mode).

- Charis SIL (=Charter; Linear Antiqua)

Samples:

- Greek
 - Monotonic
Γλυκά θροεῖ η κουκουναριά σης ρεμματιάς το πλάι, ὅμως και συ, γιδοβοσκέ, γλυκειά φλογέρα παίζεις· δῶρο σου πρέπει δεύ-τερο, ὕστερ' ἀπό τον Πάνα.
 - Polytonic
Ἄνδρα μοι ἔννεπε, μοῦσα, πολύτροπον, ὃς μάλα πολλὰ πλάγχθη, ἐπεὶ Τροίης ἱερὸν πτολίεθρον ἔπερσεν: πολλῶν δ' ἀνθρώπων ἴδεν ἄστεα καὶ νόον ἔγνω, πολλὰ δ' ὃ γ' ἐν πόντῳ πάθεν ἄλγεα ὃν κατὰ θυμόν, ἀρνύμενος ἦν τε ψυχὴν καὶ νόστον ἐταίρων.
- Cyrillic
 - Bulgarian: Всички хора се раждат свободни и равни по достойнство и права.
 - Russian: Все люди рождаются свободными и равными в своем достоинстве и правах.

– Ukrainian: Всі люди народжуються вільними і рівними у своїй гідності та правах.

- IPA

ˌɪntəˈnæʃnəl fəˈnɛtɪk ˈælfəbet

In addition to the IPA, many languages with latin-based scripts contain additional characters that are not merely accented variants of existing base characters. E.g. West-African languages have capitals for e.g. implosives, as in Hausa:

- (1) Barawo ya saci nama.
thief 3.sg.compl steal meat/animal
‘The thief stole the meat.’

These glyphs are typically found in section Latin Extended B.

Unicode only caters for some precombined letters with diacritics. Combinations without their own code position can be produced with combining diacritics. Precombined and combining diacritics should be optically indistinguishable. This is typically achieved with OpenType rules.

- (2) Bàrāwō yā s̄aci nāmā.
thief 3.sg.pot steal meat/animal
‘The thief might steal the meat.’
- (3) Bàrāwō yā s̄aci nāmā.
thief 3.sg.pot steal meat/animal
‘The thief might steal the meat.’

Many fonts lack the appropriate mark-to-mark anchors for recursive stacking of diacritics: the SIL fonts are a notable exception. I have created versions of FreeSerif and Linux Libertine that do support this. Accent placement for Linux Libertine is done mostly automatically; for FreeSerif, I heavily take advantage of existing accent marks.

And finally, some Vietnamese:

Tuyên ngôn toàn thế giới về nhân quyền của Liên Hợp Quốc. Tất cả mọi người
i sinh ra đều được tự do và bình đẳng về nhân và quyền. Mọi con người i đều được
tạo hoá ban cho lý phải đối xử với nhau trong tình bằng trí và lương tâm và
cân hữu.

2 Math

Empirical linguistics and formal linguistics alike need good math support. Modulo kerning (which is necessarily different in math and text modes), multilingual text fonts should have accompanying math fonts. Math support includes not only symbols for statistics, set theory or logic but also stretchable operators for e.g. AVMs.

Owing to the limited availability of unicode math fonts, it seems necessary to rely on non-unicode LaTeX math packages, such as Mathdesign Charter for Charis, and newtxmath for Libertine. For FreeSerif only, there is an equally Times-based unicode math font, i.e. TG Termes Math.

- Comparison of text vs. math italics (compare character weight).

aa bb ff

Charter appears slightly lighter than Charis SIL, although the design matches.

- Some inline maths:

- (4) a. For any leaf type $t_1[\text{mud } \mu_1, \text{morsyn } \sigma]$, $t_2[\text{mud } \mu_2, \text{morsyn } \sigma \wedge \tau]$ is a morphological competitor, iff $\mu_1 \subseteq \mu_2$.
- b. For any leaf type t_1 with competitor t_2 , expand t_1 's morsyn σ with the negation of t_2 's morsyn $\sigma \wedge \tau$: $\sigma \wedge \neg(\sigma \wedge \tau) \equiv \sigma \wedge \neg\tau$.

- AVMs with math operators

$$(5) \quad \text{word} \rightarrow \left[\begin{array}{l} \text{morphs} \quad \boxed{e_1} \circ \dots \circ \boxed{e_n} \\ \text{morsyn} \quad \boxed{0} (\boxed{m_1} \uplus \dots \uplus \boxed{m_n}) \\ \text{rules} \quad \left\langle \left[\begin{array}{l} \text{morphs} \quad \boxed{e_1} \\ \text{mud} \quad \boxed{m_1} \\ \text{morsyn} \quad \boxed{0} \end{array} \right], \dots, \left[\begin{array}{l} \text{morphs} \quad \boxed{e_n} \\ \text{mud} \quad \boxed{m_n} \\ \text{morsyn} \quad \boxed{0} \end{array} \right] \right\rangle \end{array} \right]$$

- Old-style figures (or not) in attribute names:

$$(6) \quad \left[\begin{array}{l} \text{c-cont} \quad \left[\begin{array}{l} \text{hook} \quad \left[\begin{array}{l} \text{index} \quad \boxed{e} \\ \text{ltop} \quad \boxed{l} \end{array} \right] \\ \text{rels} \quad \left\langle \left[\begin{array}{l} \text{pred} \quad \text{hortative-rel} \\ \text{lbl} \quad \boxed{l} \\ \text{arg1} \quad \boxed{h} \end{array} \right] \right\rangle \\ \text{hcons} \quad \langle \boxed{h} =_q \boxed{l} \rangle \end{array} \right] \\ \text{dtrs} \quad \left\langle \left[\text{ss} | \text{l} | \text{cont} \quad \left[\begin{array}{l} \text{hook} \quad \left[\begin{array}{l} \text{index} \quad \boxed{e} [\text{tam} \quad \text{subj}] \\ \text{ltop} \quad \boxed{l} \end{array} \right] \\ \text{rels} \quad \left\langle \dots \left[\begin{array}{l} \text{arg0} \quad \boxed{e} \\ \text{lbl} \quad \boxed{l} \end{array} \right] \dots \right\rangle \end{array} \right] \right] \right\rangle \end{array} \right]$$

Note: Although the unicode-math package permits substitution of math letters with glyphs from the unicode text font, this is typographically questionable, due to incorrect metrics.

3 Summary

If you want widest Unicode coverage, go for FreeSerif. Kerning is sometimes not optimal (e.g. Test).¹ In addition to LGC you even get Ethiopic, Arabic, Hebrew Devanagari, and lots more with an essentially Timesish design. Plus: TG Termes Math gives you a 100% matching unicode math font. Minus: The FreeSans font has some kerning issues (Test). There are no oldstyle figures.

¹You may like to experiment substituting certain character ranges with TG Termes, but, watch out (!), there are no anchors for diacritics (yet). Of course, these may be added by script at first, but I do not know how this could be fed into the TeX Gyre workflow: in the long run, one would want to fine-tune anchor positions by hand.

Linux Libertine has very broad unicode coverage, although only LGC + Hebrew. The newtx-math package provides adequate math support, though not unicode math. Plus: there are 3 weights (normal semi-bold, extra-bold), a display font and a nice companion sans font. Unicode coverage for extra bold sans is reduced, though. The font provides old-style figures, as well as calligraphic ligatures.

Charis has very good general support for Latin and Cyrillic with sophisticated diacritic placement, but no Greek. Another minus: the Mathdesign Charter fonts are slightly too light, though one might adjust the stem width in FontForge by a few points, and hope that this does not mess up the metrics too much.

As for sans serif and mono fonts with good unicode coverage, there are the DejaVu Sans and DejaVu Sans Mono.² There is even good non-unicode math support via the arev math package. I am not so sure that DejaVu goes that well with Times, though, because of their reduced ascenders.

²Some anchors for diacritics are present in DejaVu Sans. I shall provide a version soon with extensive anchors.