

การรองรับภาษาไทยในซอฟต์แวร์โอเพนซอร์ส

เทพพิทักษ์ การุญบุญญานันท์
theppitak@gmail.com

Thai Linux Working Group

กันยายน 2564

- รหัสอักขระ
 - ยุค 8 บิต: TIS-620, ISO-8859-11
 - ยุค multilingual: Unicode, ISO/IEC 10646
- ข้อกำหนดท้องถิ่น (POSIX locale, ISO/IEC 14652, Unicode CLDR)
 - string collation (LC_COLLATE, ISO/IEC 14651, UTS #10)
 - date/time format
 - numeric format
 - currency & format
 - etc.

การแสดงข้อความ

- การจัดเรียงสระและวรรณยุกต์

- หลบหาง ป ฝ ฟ (พ)

น้ำนี้ฟี่ฟ้าปี่ป่าฟ้าฝุ่น

- หลบหาง/แปลงรูป ฎ ฏ

กฏุมพื ตริกฏุก กุฏฏฏฏฏ

กฏุมพื ตริกฏุก กุฏฏฏฏฏ

- แปลงรูป ฎ ฐ

กัตถุญญ ทิฏฐจุชุกรรม

- คาถาบาลี-สันสกฤต
 - ญู ฐู ตัดเชิง

ปถุณายติ ทิฏฐจา

- รongรับการช้อนนคหิตเหนือสระอิ (ที่ไม่ใช่สระอี)

จกขุสมปี

- Implementation

- ยุค Windows XP และก่อนหน้า: PUA glyphs
 - glyph หลายชุดสำหรับวางในตำแหน่งต่างๆ หรือแปลงรูป
 - rendering engine ต้องรู้ว่าจะใช้ชุดไหนเมื่อไร
 - ฟอนต์รองรับได้เท่าที่ rendering engine ทำ
 - Windows กับ Mac แยกใช้รหัสคนละชุด → ฟอนต์ใช้ข้ามระบบไม่ได้
 - Pango: detect และรองรับฟอนต์จากทั้งสองแพลตฟอร์ม

- Implementation

- ยุคหลัง Windows XP: OpenType
 - GSUB เลือก/เรียงเรียง glyph

เลือกวรรณยุกต์ต่ำ-สูง

ตัดเชิง ญ

แยกส่วนสระอำ

คูคี

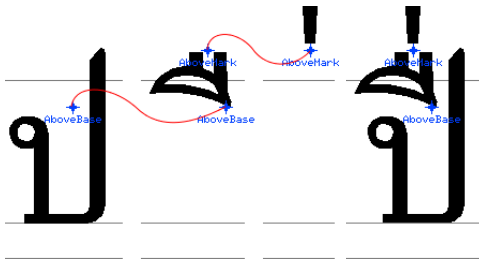
วิญญู

ป่า

- ฟอนต์สามารถเพิ่มรูปแบบได้อย่างอิสระ

- Implementation

- ยุคหลัง Windows XP: OpenType
 - GPOS กำหนด anchor สำหรับวางซ้อนอักขระ



- ไม่จำเป็นต้องเพิ่ม glyph เพื่อการจัดตำแหน่ง
- ปรับตำแหน่งโดยละเอียดได้

ประเด็นเปิดของฟอนต์

- ความสูงของตัวอักษร
 - Scaled down (Windows, Mac)
 - TH Sarabun* 16pt + Arial 12pt
 - MS Word, LibreOffice ต้องแยกฟอนต์ CTL, Latin
 - Latin-compatible (Linux, L^AT_EX, Google Fonts)
 - Kinnari 12pt + Arial 12pt
 - ผสมกับฟอนต์สากลได้โดยไม่ต้องแยกขนาด

ประเด็นเปิดของฟอนต์

- ความสูงของตัวอักษร (ต่อ)
 - นิยามของ point size?



(ที่มา: เพจเฟซบุ๊ก “แกะรอยตัวพิมพ์ไทย”)

ประเด็นเปิดของฟอนต์

- ความสูงของตัวอักษร (ต่อ)
 - OS/2 Metrics: TypoAscent, TypoDescent, WinAscent, WinDescent



- spec ของไมโครซอฟท์ที่ฟอนต์ไทยในวินโดวส์ไม่ใช่???

ประเด็นเปิดของฟอนต์

- การรองรับคาถาบาลี-สันสกฤต
 - การตัดเชิง ญ ฐ
 - mark up language เป็นบาลีหรือสันสกฤต
 - ใช้ stylistic alternative หรือ stylistic set
 - ปัจจุบัน:
HTML/CSS, LibreOffice, L^AT_EX → ทำได้
MS Office → ไม่รองรับ
 - การซ้อนนิกหิตเหนือสระอิ
 - GPOS
 - ปัจจุบัน:
HTML/CSS, LibreOffice, L^AT_EX, MS Office → ทำได้

ประเด็นเปิดของฟอนต์

- การรองรับภาษาชาติพันธุ์
 - การยืมอักขระไปใช้ต่างหน้าที่ → ซ่อนอักขระได้ไม่จำกัด
ปะเต็ล โล็ญ บัวฮ ท็อง เป็ว มุ่ย
เต็ง เจ๋อ เปริ่ห้ โจ้ เป็ย โทร ม็อง เต็ง อ้า ย้า
ปี่ปี่ปี่ปี่
จ็อรุ การุ
 - ฟอนต์:
 - เตรียม anchor สำหรับการซ่อน
 - รองรับอักขระพิเศษ เช่น ตัวขีดเส้นใต้ (Macron)

- Frameworks
 - XIM/XKB
 - GTK/Qt IM Module
 - SCIM
 - iBus
 - uim
 - fcitx

- keyboard layout
 - เกษมณี, มอก.820-2538
 - ปัดตะโชติ
 - มนูญชัย!
- sequence check/correction
 - วทท 2.0
 - ภาษาชาติพันธุ์

- Line break
 - word boundary
 - UAX #14 Unicode Line Breaking Algorithm
 - dictionary-based engine
 - hyphenation
 - T_EX hyphenation patterns
- Word break
 - cursor by word
 - text double-click

- กระบวนการ
 - peer review ก่อน submit
 - อภิปรายเพื่อเลือกคำแปล
 - พัฒนา glossary เพื่อคำแปลที่สอดคล้องกัน
 - ผู้ใช้รายงานปัญหาเพื่อปรับคำแปล
- ข้อสังเกต
 - มีผู้เปิดใช้คำแปลไทยมากขึ้น
 - งานแปลที่ไม่ใช่โอเพนซอร์สก็ดูมีคุณภาพมากขึ้น ;-)

Applications

- Spell check
- Soundex
- ...