

Characters and CSS for Fonts

NTT Network Innovation Lab.

KAWABATA, Taichi

Prologue

In this presentation, I will explain 3 character/fonts-related topics that may affect the standardization of CSS3.

- ❑ Character Display
 - ❑ Private Character and Font
 - ❑ IVS and Font Matching
- ❑ Character Search and Comparison
 - ❑ Normalization

Private Characters and Font Formats

Why Private Characters?

▣ Unicode Standard, Section 1.1

Note, however, that the Unicode Standard **does not encode idiosyncratic, personal, novel, or private-use characters, nor does it encode logos or graphics.** Graphologies unrelated to text, such as dance notations, are likewise outside the scope of the Unicode Standard. Font variants are explicitly not encoded. **The Unicode Standard reserves 6,400 code points in the BMP for private use,** which may be used to assign codes to characters not included in the repertoire of the Unicode Standard. **Another 131,068 private-use code points are available** outside the BMP, should 6,400 prove insufficient for particular applications.

What are Private Characters

- ❑ (Company) Logos and Graphics
- ❑ Unencoded characters
 - ❑ Emoticons (Emoji) and Special Symbols
- ❑ Some Kanji (CJK Ideograph) characters

Symbols and Emoticons

- ▣ In Books, special symbols are sometimes used to convey the complex or abstract idea in simpler mannar.
- ▣ Various Emoticons have become popular. Unicode begins to encode some of them, but not (yet) all of them.

NESTLE-ALAND 26TH EDITION		
295	ΚΑΤΑ ΙΩΑΝΝΗΝ	13,11-24
	καθαρός ὅλος· και ὑμεῖς καθαροί ἐστε, ἀλλ' οὐχι πάντες.	15,3
	11 ἦδει γάρ τὸν παραδιδόντα αὐτόν· ὁ δὲ διὰ τοῦτο εἶπεν ὅτι οὐχι πάντες καθαροί ἐστε. ᾤ	6,64
	12 Ὅτε οὖν ἐνίψεν τοὺς πόδας αὐτῶν ὁ [και] ἔλαβεν τὰ ἴματια αὐτοῦ (και ἀνέπεσεν): πάλιν:1, εἶπεν αὐτοῖς· γινώσκετε τί πεποίηκα ὑμῖν; 13 ὑμεῖς φωνεῖτέ με· ὁ διδάσκαλος, και· ὁ κύριος, και καλῶς λέγετε· εἰμι γάρ. 14 εἰ οὖν ἐγὼ ἐνίψα ὑμῶν τοὺς πόδας ὁ κύριος και ὁ διδάσκαλος, και ὑμεῖς ὀφείλετε ἀλλήλων νίπτειν τοὺς πόδας· 15 ὑπόδειγμα ὁ γάρ ἔδωκα ὑμῖν ἵνα καθὼς ἐγὼ ἐποίησα ὑμῖν και ὑμεῖς ποιήτε. 16 ἀμήν ἀμήν λέγω ὑμῖν, οὐκ ἔστιν δοῦλος ὁ μείζων τοῦ κυρίου αὐτοῦ οὐδὲ ἀπόστολος· ὁ μείζων τοῦ πέμψαντος αὐτόν. 17 εἰ ταῦτα οἴδατε, μακάριοι ἐστε ἐὰν ποιήτε αὐτά.	7 Mt 23,8.10 1T 5,10 1J 2,6; 3,16 Mt 10,24; p Jc 1,25 L 10,28.37 6,70; 15,16.19 E 1,4 Ps 41,10 Is 46,10; 43,10· 14,29; 16,1-4 Mt 24,25 p· 8,241 Mt 10,40!
119 X	18 Οὐ περὶ πάντων ὑμῶν λέγω· ἐγὼ οἶδα ἑτίνους ἐξελεξάμην· ἀλλ' ἵνα ἡ γραφή πληρωθῇ· ὁ τρώγων ἔμμου τὸν ἄρτον ἑπιθήσεται ἐμὲ τὴν πτέρναν αὐτοῦ. 19 ἀπ' ἄρτι λέγω ὑμῖν πρὸ τοῦ γενέσθαι, ἵνα ἰπιστεῦσητε ὅταν γένηται ὅτι ἐγὼ εἰμι. 20 ἀμήν ἀμήν λέγω ὑμῖν, ὁ λαμβάνων ἀν τινὰ πέμψω ἐμὲ λαμβάνει, ὁ δὲ ἐμὲ λαμβάνων λαμβάνει τὸν πέμψαντά με.	
221 P	21 Ταῦτα εἰπόν ὁ [ὁ] Ἰησοῦς ἐταράχθη τῷ πνεύματι και ἐμαρτύρησεν και εἶπεν· ἀμήν ἀμήν λέγω ὑμῖν ὅτι εἰς ἐξ ὑμῶν παραδώσει με. 22 ἐβλεπον τ εἰς ἀλλήλους οἱ μαθηταί τ ἀπορούμενοι περὶ τίνος λέγει. 23 ἦν τ ἀνακειμένος εἰς ἐκ τῶν μαθητῶν αὐτοῦ ἐν τῷ κόλπῳ τοῦ Ἰησοῦ, ὃν ἠγάπα ὁ Ἰησοῦς. 24 νεύει οὖν τούτῳ Σίμων Πέτρος	21-30; Mt 26, 21-25 Mc 14,18-21 L 22,21-23 11,33!
221 X		19,26; 20,2; 21, 7,20
<small> 11 □ D ON A Θ f¹⁻¹³ ᾠ e vg txt P⁶⁶ B C L W Ψ pc it • 12 O P⁶⁶ K A C² L Ψ 33, 1241 al it vg^s sy^{s-p} txt B C^s D W Θ f¹⁻¹³ ᾠ lat sy^h (αναπισσαν C² D Θ f¹⁻¹³ ᾠ vg sy^h και αναπ. P⁶⁶ N² A(*h. t.) L Ψ 33, 1241 al it vg^s txt N^s B C^s W pc e p sy^{s-p} [·, et¹⁻¹³] • 14 T ποσω μαλλον D Θ it (sy^{s-p}) • 15 O P⁶⁶ 700 pc d Γ δεδ- P⁶⁶ K A K Ψ f¹⁻¹³ 28, 33, 700, 892, 1241 pm txt B C D L W Γ Δ Θ 1010^s, 1424 pm • 16 □ Θ bo^{ms} O P⁶⁶ • 18 Γ ους P⁶⁶ A D W Θ Ψ f¹⁻¹³ ᾠ; Eus txt K B C L 33, 892, 1241 pc; Or P μετ εμου P⁶⁶ K A D W Θ Ψ f¹⁻¹³ ᾠ lat sy bo; Eus Eriph txt B C L 892 pc (a) sa; Or Γ¹ επρηκν K A Θ W pc O P⁶⁶ B • 19 (2 3 / A D W Θ Ψ f¹⁻¹³ ᾠ it vg^s † -ουητε οτ. γεν. B (fC) txt P⁶⁶ K L pc • 20 T και P⁶⁶ • 21 O t P⁶⁶ K B L txt P⁶⁶ A C D W Θ Ψ f¹⁻¹³ ᾠ; Or • 22 Του P⁶⁶ K^s A D L W Θ f¹⁻¹³ ᾠ sy^h δε pc a sy^{s-p} txt K^s B C Ψ pc e P αυτου P⁶⁶ f¹³ 1241 pc a r¹ sy^s bo • 23 T δε P⁶⁶ K A C² D W Θ f¹⁻¹³ ᾠ latt sy^{s-h*} txt B C^s L Ψ 892, 1424 pc sy^s O P⁶⁶ B </small>		

Kanji (CJK Ideographs)

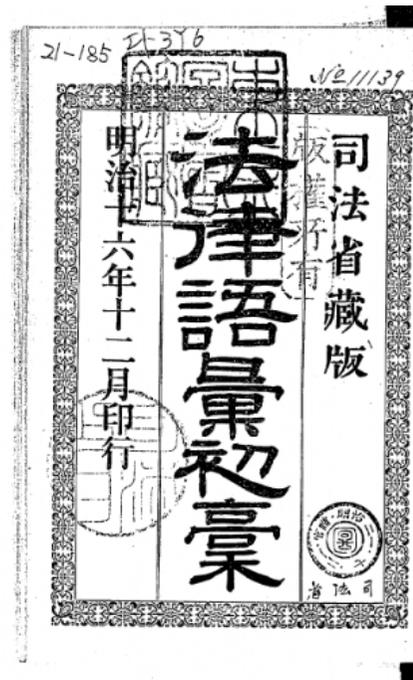
- 75,000+ ideographs are encoded under the Unification Rules.
- However, there are still more ideographs not yet encoded.
 - Misdescribed ideographs
 - Invented ideographs
 - (Very) Local ideographs
 - Historic, short-lived ideographs
 - Rare variants

隆
楚簡
切。隆者漢
祭而祭於陵。既不廟祭。似可不諱。然
劭曰。隆慮山在北。避殤帝名。改曰林
隆者漢殤帝之名。字者五經異義云。

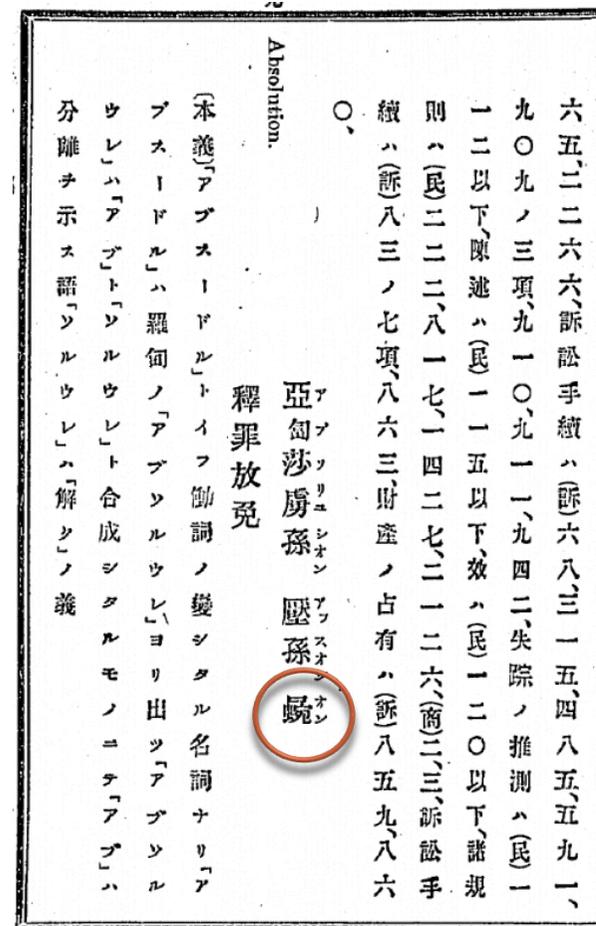
【奚容】⁴⁶ 複姓。奚容蒧を見よ。
【奚容蒧】⁴⁷ 春秋、衛の人。字は子皙。孔子の弟子。〔史記、仲尼弟子傳〕奚容蒧、字子皙。
〔注〕正義曰、衛人。
【奚落】⁴⁸ 面前でけなす。罵る。〔通俗編、言笑、譏落〕高則誠琵琶曲有「奚落語」、奚、蓋譏之。

Invented Ideographs

- Many ideographs are invented, but not widely used, in the past time.



In late 19th century, Meiji Government created hundreds of characters to describe new legal vocabraries, but they were never used.



Private Characters

- ▣ There are several ways to represent private characters

Methods	Prop	Search	Num	Editor	Notes
Assign Private Character to = (U+3013)	○	×	1	×	Not proper use
Assign Private Character to ◀ (U+FFFD)	×	×	1	×	
Assign variants to encoded ideographs	○	○	1+	×	Only applicable to variants.
Use Private Use Area (PUA)	×	×	137,000+	○	
Use inline image ()	×	×	—	×	Hinting, etc. can not be applied.

Font Formats

Font format	Pros	Cons
OpenType	Most popular. Various tools available.	Large Size. Piracy problem
WOFF	Small Size. Tailored for Web use.	You can't use it on OS.
SVG	gradation, color, animation Embeddable in HTML text. Can inherit CSS.	Not all browsers support it.

Emoticon Characters use colors.



GlyphWiki: Create/Share Private Characters

- <http://glyphwiki.org/>
- Everyone can edit/register his/her own characters.
- Fonts can be generated by creating independent wiki pages.
- Nearly 100,000 characters are registered.
- Hanazono-Mincho (花園明朝) is only one free font which covers all UCS/AJ1 ideographs.

GlyphWiki

u26f97 (@3)

出典: フリーグリフデータベース『グリフウィキ(GlyphWiki)』

蓋 蓋 蓋 (SVG画像 ポリゴン パス) (EPS画像)

文字コード関連情報

- CHISE character description (説明)
- chise_linkmap (説明)
- IPSJ-TS 0008:2007 (説明)
- Unicode Unihan Database (説明)

あなたのブラウザでの表示: 蓋

前の符号位置: 蓋 (u26f96)

次の符号位置: 蓋 (u26f98)

関連グリフ

関連字: 蓋 (u26f97) 蓋 (u26f97)

異体字: 苙 (u8369) 苙 (u8369)
蓋 (u85ce) 蓋 (koseki-366050)
蓋 (u26cd2) 蓋 (u26cd2)
蓋 (u270e4) 蓋 (u270e4)

Private Chars and Text Orientation

- ❑ Character's behavior on different text orientation is determined by character's East Asian Width property.
- ❑ PUA characters do not have any property.
- ❑ One possible solution is to look at VORG/vert/vrt2 tables to determine character orientation.

When 'text-orientation' is 'vertical-right', set all characters upright (using vertical font settings if available) unless otherwise specified above.

In OpenType, vertical font settings are provided by the vhea, vmtx, and VORG tables, as well as the vert and vrt2 GSUB features. If any of these are present, the font is considered to have vertical font settings available. [CSS3 writing-mode]

IVS and Font Selection

Ideographic Variation Sequence (IVS)

- IVS enables to display Ideographic Variants by attaching “Variation Selector” after the “Base Character”.

e.g. “U+559D U+E0101” → 喝
 “U+559D U+E0100” → 喝

- Character code only specifies **abstract character**, but variation selector can specify **concrete glyph** for a character.

Ideographic Variation Database (IVD)

▣ UTS (Unicode Technical Standard) #37

registrant



registrar



IVD



Register collection



IVD_Collections.txt
will be updated.

Register glyphs (1)



VS are assigned.
IVD_Sequences.txt
will be updated.

Register glyphs (2)



IVD	Ideographic Variation Database
IVC	Ideographic Variation Collection
IVS	Ideographic Variation Selector
VS	Variation Selector

Ideographic Variation Collections

- Currently, Two collections are registered.

Collection Name	Purpose of Collection	Implemented Fonts available
Adobe-Japan1	To use with Japanese Desktop Publishing Systems.	Kozuka Mincho (小塚明朝)/Gothic(小塚ゴシック), etc
Hanyo-Denshi (汎用電子)	To use with Administration Systems of Japanese Government	IPAmjm Mincho

Glyph Correspondence

□ 「邊」 (AJ1: 15 glyphs, Hanyo-Denshi : 15 glyphs)



Not all variants are common between two collections

From <http://d.hatena.ne.jp/NAOI/20100406/1270550459>

Why IVS?

- ▣ Two main usage
 - ▣ To show “archaic style”
 - ▣ To correctly display human or proper names.

(Demonstrations)

CSS3 Font Matching Algorithm

- CSS font-family designation:

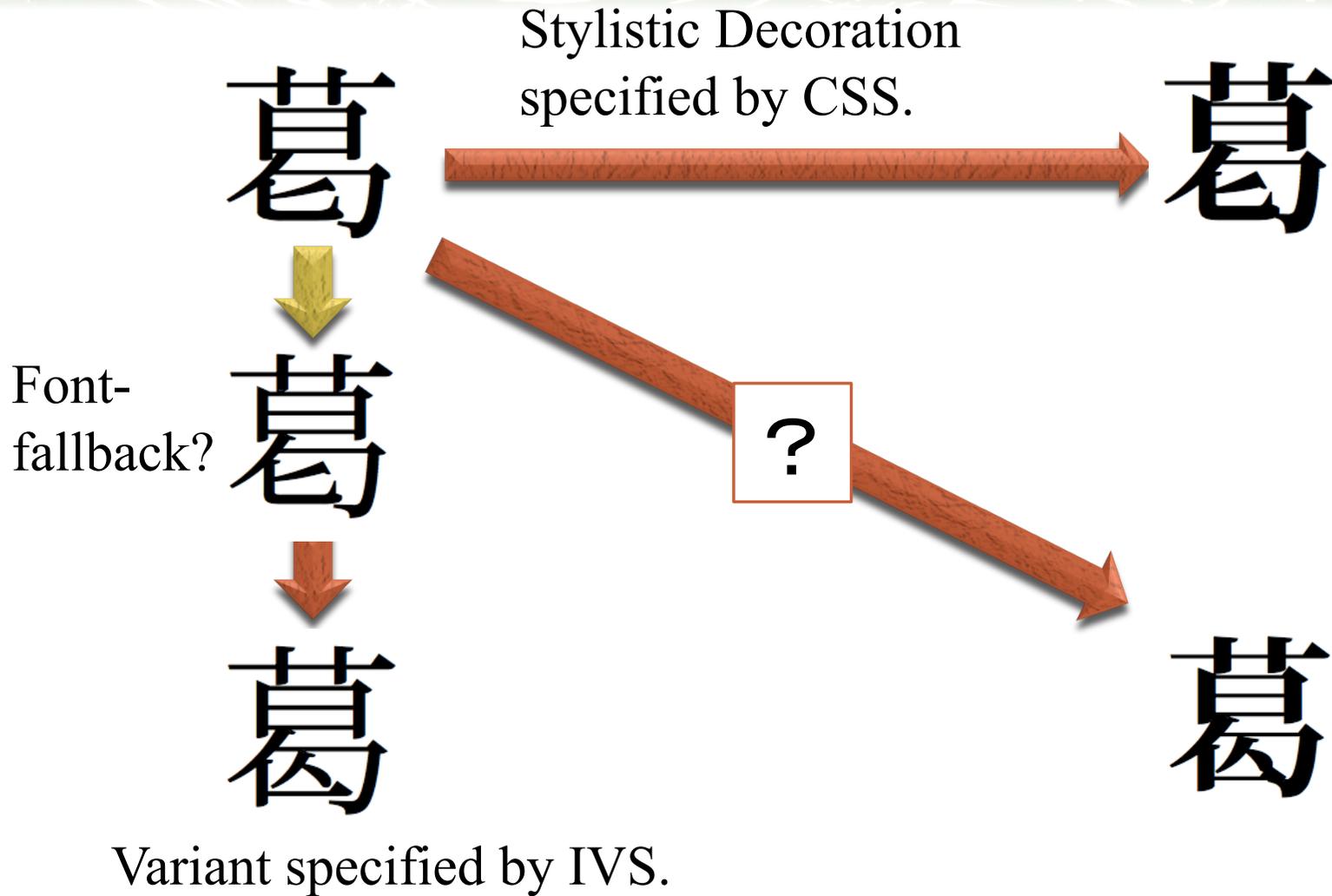
```
font-family: font-A, font-B, font-C;
```

- Character (Grapheme Cluster) Sequence

```
C1 C2 C3 C4 C5.....
```

- If fonts in *font-A* family contains all glyphs in C₁, then best font is selected by the following font spec order.
 - font-stretch (condensed/stretched)
 - font-style (italic/oblique)
 - font-weight (bold/light)

Two different text decorations



Font Matching vs. IVS

- Consider CSS font-family designation:

font-family: *font-A*, *font-B*, *font-C*;

where *font-A*: supports only base characters.

font-B: supports IVC_X

font-C: supports IVC_Y

- Consider Character (IVS) Sequence:

$C_1 IVS_X (\in IVC_X) IVS_Y (\in IVC_Y) C_2 \dots$

- In which font family should IVS_X/IVS_Y be rendered?

- Option A: $C_1 BASE_X BASE_Y C_2$ (all by *font-A*)

- Option B: $C_1 IVS_X IVS_Y C_2$ (by each IVS font.)

Which is better?

- ❑ Option A
 - ❑ Pro: Whole text has a consistent font-family.
 - ❑ Con: Multiple IVC fonts can not be supported.
- ❑ Option B
 - ❑ Pro: Each IVS will be rendered with their fonts.
 - ❑ Con: Text may be displayed in inconsistent font.

- ❑ Under option A, it is **difficult** for user to display each IVS(font-family must be specified in each IVS).
 - ❑ Font is specified by CSS font-family spec.
- ❑ Under option B, it is **easy** for user to display only base character (just remove VS characters) .
 - ❑ Variant is specified with VS character

Normalization

What is Normalization?

- ❑ Unicode has MANY ways to describe the same/similar character. Therefore, bitwise comparison do not guarantee correct result.

- ❑ e.g. U+1EB6 (NFC)

U+1EA0 U+0306

U+0041 U+0323 U+0306 (NFD)

U+0102 U+0323

U+0041 U+0306 U+0323



are ALL the same character.

- ❑ If text is NORMALIZED, then text is guaranteed for comparison.
 - ❑ NFD ... diacritical marks are all decomposed.
 - ❑ NFC ... diacritical marks are composed in unique manner
 - ❑ NFKC/NFKD ... squared/circled chars are also decomposed.

Normalization Problem

- ❑ Implementation is cumbersome
 - ❑ NFC requires complex algorithm
- ❑ Singleton Decomposition
 - ❑ Different character will be folded to the same character. (It may go beyond specific subset)
 - ❑ Å (U+212B / JIS X 0208) → Å (U+00C5 / ISO8859-1)
- ❑ Compatibility Ideographs
 - ❑ All compatibility ideographs will be transformed to corresponding unified ideographs.
 - ❑ Apple proposed to exclude Compatibility Characters from normalization, but it was not adopted. (10 years ago)

Compatibility Ideographs

- Ideographs are unified by unification rule specified in “ISO/IEC 10646 Annex S.”
 - Unifiable ideographs specified **before 1992** are separately encoded into **Unified Ideographs**.
(Source Code Separation Rule)
e.g. 飲 (U+98F2) vs. 飲 (U+98EE)
 - Unifiable ideographs specified **after 1992** are assigned to **Compatibility Ideographs**.
- For example, Japanese compatibility ideographs include Name characters specified by the ordinance of the Ministry of Justice of Japan.
e.g. 「社」 vs. 「社」 / 「者」 vs. 「者」

Normalization : When and Where?

- ❑ Early Unicode Normalization (2005)
 - ❑ **Character Model for the World Wide Web 1.0: Normalization**
 - ❑ Since then, it has become widely known that EUN has various problems.
 - ❑ <http://www.w3.org/2011/05/04-i18n-minutes.html>
- ❑ HTML Text should not be early normalized, but **ID** and **Class** values in HTML5 may/should be normalized.
 - ❑ XML specification Appendix J
- ❑ On web browser, “**NFKC**” normalization is more useful.
 - ❑ e.g. “シ ュース” can be searched by “シ ュース”
 - ❑ e.g. “リッ” can be searched by “リットル”
- ❑ Old Kanji characters can not be searched by new characters. New normalization must be conceived.
 - ❑ e.g. 「小澤」 can not be searched by 「小沢」

Thanks for Listening!