

CSSにおける文字とフォント

NTT未来ねっと研究所

川幡 太一

最初に

文字について、CSSの標準化に影響しそうな**3トピック**（外字とフォント、IVSとフォント選択、正規化）について、現在の状況と課題を、例を交えて説明する。

- 文字の表示
 - 外字とフォント
 - IVSとフォント選択
- 文字の照合
 - 正規化

外字とフォント

外字はなぜ必要か

□ Unicode標準 第1章1節

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.

外字の種類

- 会社のロゴやシンボルマーク
- 未符号化文字
 - 絵文字・記号類
- 未符号化の漢字

記号類・絵文字

□ 繰り返し現れる項目や、複雑な概念を抽象化・簡潔化して表記するため、独自文字を使用。

□ 日本の携帯電話の絵文字

NESTLE-ALAND 26TH EDITION		
295	ΚΑΤΑ ΙΩΑΝΝΗΝ	13,11-24
	καθαρός ὁλος· και ὑμεῖς καθαροὶ ἔστε, ἀλλ' οὐχὶ πάντες.	15,3
	11 ἦδει γάρ τὸν παραδιδόντα αὐτόν· ὁ δὲ αὐτοῦ εἶπεν ὅτι οὐχὶ πάντες καθαροὶ ἔστε. ὁ	6,64
	12 Ὅτε οὖν ἐνίψεν τοὺς πόδας αὐτῶν ὁ [και] ἔλαβεν τὰ ἴματα αὐτοῦ (και ἀνέπεσεν): πάλιν:1, εἶπεν αὐτοῖς· γινώσκετε τί πεποίηκα ὑμῖν; 13 ὑμεῖς φωνεῖτέ με· ὁ διδάσκαλος, και· ὁ κύριος, και καλῶς λέγετε· εἰμὶ 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 ποιῆτε. 16 ἀμὴν ἀμὴν λέγω ὑμῖν, οὐκ ἔστιν δοῦλος ὁ μείζων τοῦ κυρίου αὐτοῦ οὐδὲ ἀπόστολος· ὁ μείζων τοῦ πέμψαντος αὐτόν. 17 εἰ ταῦτα οἴδατε, μακάριοι ἔστε ἐὰν ποιῆτε 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 ἄτινας ἐξελεξάμην· ἀλλ' ἵνα ἡ γραφή πληρωθῇ· ὁ τρώγων ἔμω τὸν ἄρτον ἢ ἐπὶ ἑνὶ ἡμέρῃ πτέρων αὐτοῦ. 19 ἀπ' ἄρτι λέγω ὑμῖν πρὸ τοῦ γενέσθαι, ἵνα ἰπιστεύσητε ὅταν γένηται ὅτι ἐγὼ εἰμι. 20 ἀμὴν ἀμὴν λέγω ὑμῖν, ὁ λαμβάνων ἄν τινα πέμψω ἐμὲ λαμβάνει, ὁ δὲ ἐμὲ λαμβάνων λαμβάνει τὸν πέμψαντά με.	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!
121 P	21 Ταῦτα εἰπόν ὁ [ὁ] Ἰησοῦς ἐταράχθη τῷ πνεύματι και ἐμαρτύρησεν και εἶπεν· ἀμὴν ἀμὴν λέγω ὑμῖν ὅτι εἰς ἐξ ὑμῶν παραδώσει με. 22 ἐβλεπον τ εἰς ἀλλήλους οἱ μαθηταὶ 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!
123 X		19,26; 20,2; 21, 7,20
<p>11 0 D ON A Θ f¹⁻¹³ 22 e vg txt 266 B C L W Ψ pc it • 12 0 266 N A C² L Ψ 33, 1241 al it vg^{s-p} txt B C² D W Θ f¹⁻¹³ 22 lat sy^h (αναπισσαν C² D Θ f¹⁻¹³ 22 vg sy^h και αναπ. 266 N² A(*h. t.) L Ψ 33, 1241 al it vg^s txt N* B C² W pc e p sy^{s-p} [; et¹⁻¹] • 14 Τ ποσω μαλλον D Θ it (sy^{s-p}) • 15 0 266^s 700 pc d Γ δεδ- 266 N 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 0 Θ bo^{ms} 0 266^s • 18 Γ ους 266 A D W Θ Ψ f¹⁻¹³ 22; Eus txt N B C L 33, 892, 1241 pc; Or Γ μετ εμου 266 N A D W Θ Ψ f¹⁻¹³ 22 lat sy bo; Eus Eriph txt B C L 892 pc (q) sa; Or Γ επηρκεν N A Θ W pc 0 266^s B • 19 Γ 2 3 / A D W Θ Ψ f¹⁻¹³ 22 it vg²¹ Γ -ουητε οτ. γεν. B (fC) txt 266 N L pc • 20 Τ και 266^s • 21 0 Γ 266^s N B L txt 266^s A C D W Θ Ψ f¹⁻¹³ 22; Or • 22 Του 266 N* A D L W Θ f¹⁻¹³ 22 lat sy^h δε pc a sy^{s-p} txt N² B C Ψ pc e τ αυτου 266 f¹³ 1241 pc a r¹ sy^s bo • 23 T δε 266 N A C² D W Θ f¹⁻¹³ 22 lat sy^{s-h*} txt B C² L Ψ 892, 1424 pc sy^s 0 266^s B</p>		

漢字

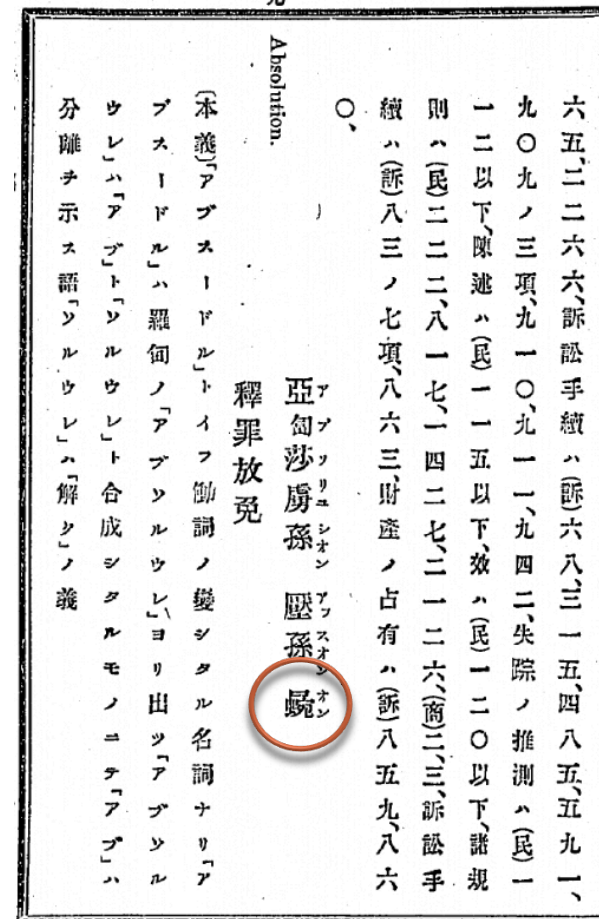
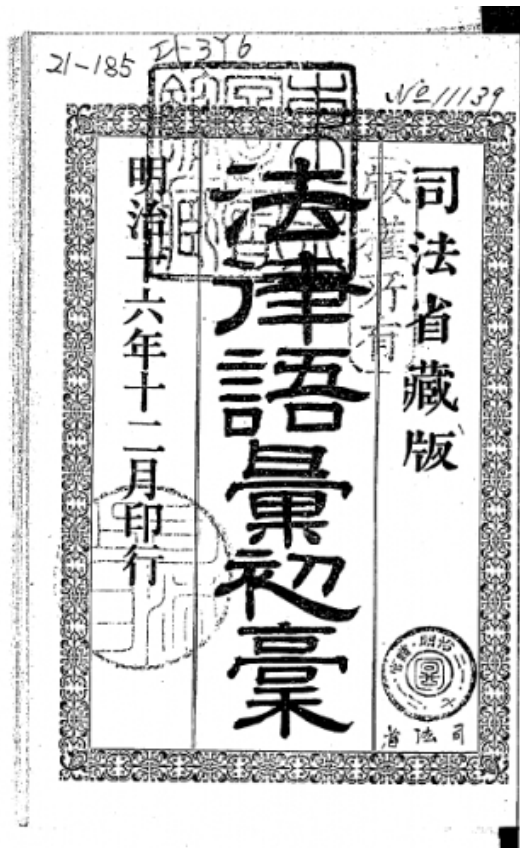
- 7万5千字以上が「統合基準」に基づいて符号化
- まだ多くの漢字が未符号化。
 - 誤字
 - 人工文字
 - 特定の地域でのみ使われる文字
 - 歴史的な文字
 - 珍しい異体字

隆 楚簡 隆 豐大也。隆者漢 陽帝之名。切。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。

【奚容】⁶⁶ 複姓。奚容蔵を見よ。
 【奚容蔵】⁴⁷ 春秋、衛の人。字は子皙。孔子の弟子。〔史記、仲尼弟子傳〕奚容蔵、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。

新造漢字

□ 漢字はしばしば、新規に作字される。



「法律語彙初稿」には、明治政府が法律語彙用に試作した数百の新造漢字が掲載されている。

外字

- CSS/HTML で外字を表現する方法はいくつかある。

方法	属性	検索	数	エディタ	メモ
= (U+3013)に外字を割り当てる	○	×	1	×	正しい使い方ではない。エディタで文字が見えない。
◆ (U+FFFD)に外字を割り当てる	×	×	1	×	縦書きかが判別しない。
符号化された漢字に異体字を割り当てる	○	○	1	×	異体字にしか適用できない。
私用文字領域 (PUA) を使う	×	×	137,000+	○	文字の属性が不明
インライン画像 ()を使う	○	×	—	×	ヒント情報が使えない。

外字フォント

フォント	メリット	デメリット
OpenType	ツールが充実	サイズが大きくなりがち ライセンス上、Webに埋め込むのは困難な場合がある。
WOFF	サイズが小さい。 Webで利用しやすい	OSではそのままでは使えない。
SVG	グラデーション・色・アニメが可能 HTMLに埋め込み可能。 CSSを継承することが可能。	ヒントや、複雑なリガチャ・カーニングなどに未対応 EPUBで必ずしもサポートされない。



GlyphWiki: 外字の共有

- ❑ 誰でもグリフを作成・編集して登録可能
- ❑ フォントは新規Wikiページを作成することで生成。
- ❑ 100,000 以上のグリフ（主に漢字）が登録済み
- ❑ 花園明朝は世界で唯一、全UCS漢字とAdobe-Japan1グリフをサポートしているフリーフォント

The screenshot shows the GlyphWiki interface for the character '蓋' (u26f97). The page includes a navigation menu, a search bar, and a sidebar with links to main page, recent updates, and help. The main content area displays the character '蓋' in three different styles: a large black font, a smaller black font, and a smaller grey font. Below the character, there is a section for '文字コード関連情報' (Text code related information) listing various character sets and their descriptions. The '関連グリフ' (Related glyphs) section lists several related characters, including '蓋' (u26f97), '苙' (u8369), '蓋' (u85ce), '蓋' (u26cd2), and '蓋' (u270e4), each with its corresponding Unicode code point and a link to its page.

外字と縦書き

- 縦書きの時に、文字を倒すかは、文字の EAW 属性などを使って決定する。
- 符号位置から、文字の性質が分からない場合、縦書き時に正立するか？

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.

IVSとフォント選択

漢字字形指示列 (IVS)

- IVSは、漢字の基底文字の符号の直後に、異体字選択子 (VS) を付加することで、異体字の表示を可能にする。

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

- 文字コードは、抽象的な文字しか指示しないが、異体字選択子は直接、具体的なグリフを指示する。

漢字字形データベース (IVD)

□ UTS (Unicode Technical Standard) #37

registrant



registrar



IVD



Register collection

Register glyphs (1)

Register glyphs (2)

IVD_Collections.txt
will be updated.

VS are assigned.
IVD_Sequences.txt
will be updated.

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

異体字コレクション (IVC)

- 現在、2つのコレクションが登録されている。

コレクション名	コレクションの目的	実装済のフォント
Adobe-Japan1	日本語のDTPでの使用	小塚明朝 小塚ゴシック
Hanyo-Denshi (汎用電子)	日本の行政システムでの使用	IPAmjm 明朝

各コレクションのグリフ対応例

□ 「邊」 (AJ1: 15 グリフ, 汎用電子: 15 グリフ)



<http://d.hatena.ne.jp/NAOI/20100406/1270550459>
より引用

なぜIVS?

- 2つの主な用途
 - テキストを「古風」に見せる（スタイル）
 - （従来のOpenTypeのjp78のような使い方）
 - 人名や固有名詞を適切に表示する。

(Demonstration)

CSS3 フォント選択アルゴリズム

- CSS font-family 指定

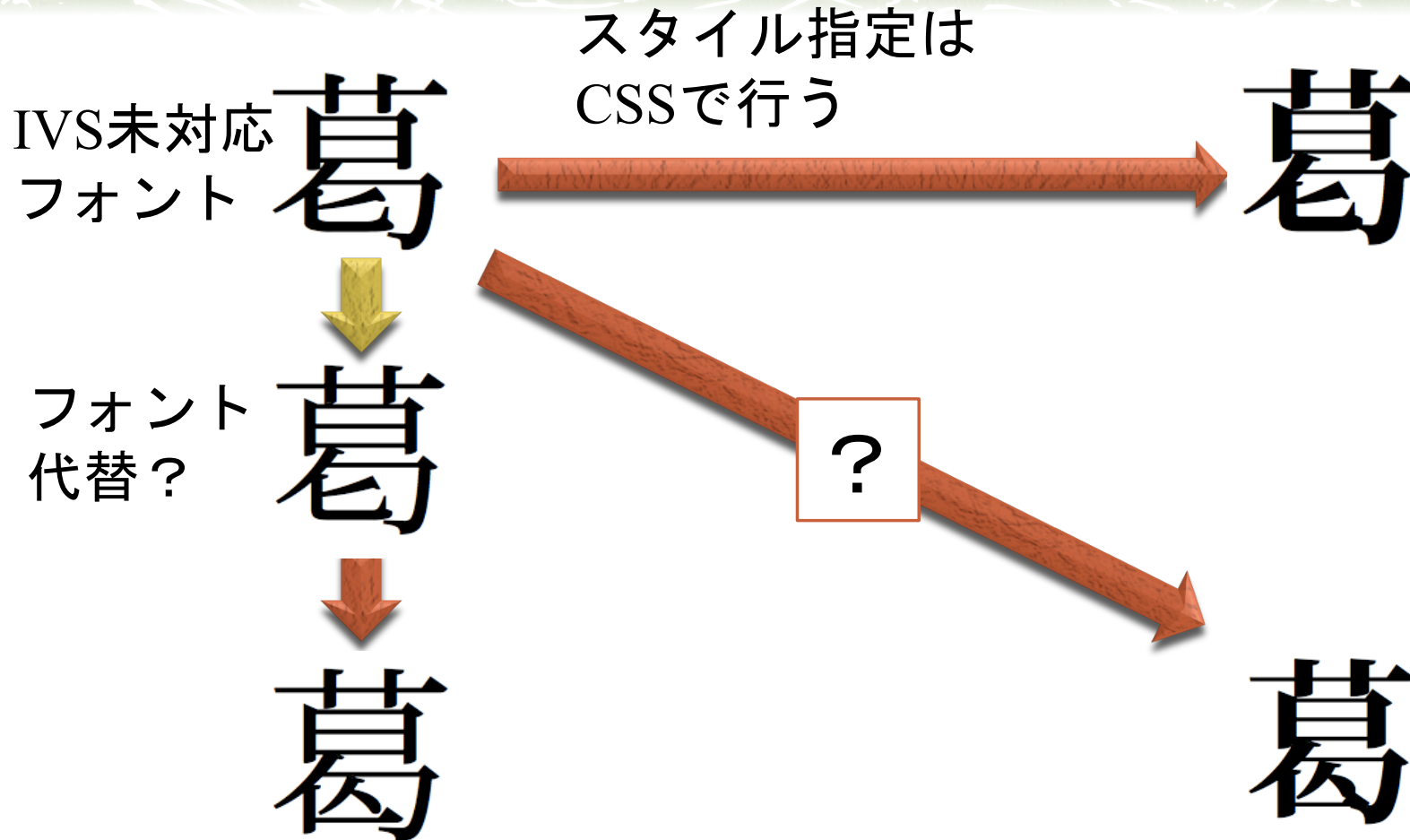
```
font-family: font-A, font-B, font-C;
```

- 文字（書記素クラス）列

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

- *font-A* ファミリで指定されたフォント群が、 C_1 に含まれる全グリフをサポートしている場合、以下の優先度でフォントが選択される。
 - font-stretch (condensed/stretched)
 - font-style (italic/oblique)
 - font-weight (bold/light)

2つの異なるテキスト装飾



異体字はIVSで指定する。

IVSとフォント選択

- CSS font-family 指定:

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

font-A: 基底文字のみ搭載

font-B: 異体字コレクションXをサポート

font-C: 異体字コレクションYをサポート

- 文字・書記素 (IVS) 列:

```
C1 IVSX (∈IVCX) IVSY (∈IVCY) C2 ...
```

- どのフォントファミリで、IVS_X/IVS_Y は表示されるべきか？

- A案: C₁ BASE_X BASE_Y C₂ (all by *font-A*)

- B案: C₁ IVS_x IVS_y C₂ (by each IVS font.)

IVSはどちらで表示すべきか？

□ A案

- 利点：テキスト全体でフォントの一貫性が保たれる。
- 欠点：複数のIVCフォントを混在できない。

□ B案

- 利点：各IVSはそれぞれに対応したフォントで表示可能。
- 欠点：テキスト全体でフォントの一貫性が欠けるかも。

□ A案では、ユーザが複数のIVCを混在させるのは**困難**。 (IVS毎にfont-familyを指定する必要がある)。

- フォントを指定するのはCSS。

□ B案では、ユーザが基底文字のみを表示するのは**簡単**。 (単にVSを削除すれば良い。)

- 異体字を選択するVSは文字符号。

正規化について

正規化とは何か？

- Unicodeは、同じ文字を表現するのに多数の方法がある。そのため、文字列をビット単位で比較しても、正しい比較はできない。

- 例：U+1EB6 (NFC)

U+1EA0 U+0306

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

U+0102 U+0323

U+0041 U+0306 U+0323



はすべて同じ文字。

- 事前に正規化すれば、ビット比較が正しいことが保証。
 - NFD ... 補助記号類はすべて分解される。(かなの濁音等)
 - NFC ... 補助記号類は同一性を保証する形で合成される。
 - NFKD/NFKC ... 丸付き文字 (㊦など) や組文字 (𪛗など) も分解 (+再合成) される。

なぜ正規化は必要か？

- 文字列を正しく比較できる。
 - 識別子で使う文字は、正規化しないとマッチングが取れない。
- HTML/CSSでは特に、id値、class名などは正規化が必要。
- 正規化は万能か？

正規化の課題

- 実装（特にNFC）が複雑
- 単一分解 (Singleton Decomposition)
 - 異なる文字が、同じ文字に畳み込まれる場合がある。(その結果、指定された部分集合外に変換される場合がある。)
 - 例：Å (U+212B / JIS X 0208) → Å (U+00C5 / ISO8859-1)
- 互換漢字
 - すべての互換漢字は、正規化によって対応する統合漢字に変換される。
 - この挙動には問題が多く、Apple等が除外を求めたが受け入れられず。

互換漢字

- UCSの漢字は、ISO/IEC 10646 追補S で記述された規則にしたがって統合される。
 - 1992年以前に規格化された統合可能な漢字は、統合漢字で分離して符号化（原規格分離規則）
 - 例：「飲」と「飲」は本来統合されて1符号になるはずが、JISの都合により分離し符号化。
 - 1992年以降の統合可能な漢字は、互換漢字で符号化。
- 互換漢字の例
 - 法務省が省令で定めた人名用漢字表の漢字（JIS X 0213）
 - 例：「社」と「社」や、「者」と「者」

正規化は、いつ、どこで行うか？

- 事前正規化...Early Unicode Normalization (2005)
 - **Character Model for the World Wide Web 1.0: Normalization**
 - HTMLテキストを必ず事前に正規化する提案だが、現在は課題が多くの人に認識されつつある。
 - <http://www.w3.org/2011/05/04-i18n-minutes.html>
- テキストの事前正規化は不適切だが、HTMLの**ID値**と**Class値**の正規化は適当かもしれない。
 - XML specification Appendix J
- ブラウザ側で“**NFKC**”正規化が便利。
 - 例：「ジュース」が「ジューズ」で検索できる。
「ㇿ」が「リットル」で検索できる。
- 漢字の旧字は未だに新字で検索不可。何か仕組が必要
 - 例：「小澤」は「小沢」で検索できない。

Thanks for Listening!