

インターネット時代の文字コード： 漢字コードの迷信を打破する！

宮澤 彰

国立情報学研究所

2011-10-05

市民講座

質問1

日本語は、世界で一番難しい言語(のひとつ)だ。

- 1) YES.
- 2) NO.

質問2

今の文字コードには漢字が足りない。

- 1) YES.
- 2) NO.

質問3

漢字はおおよそ何文字ある？

- 1) 2万字。
- 2) 5万字。
- 3) 8万字。

質問4

漢字には、意味がある。

- 1) YES.
- 2) NO.

文字化けしたウェブページ



コンピュータ内部の文字表現

ファイル上のビット列

11100110 10111001

Shift-JIS

EUC-JP

ISO-8859-1

35977

貉

32666

羚

230

æ

185

1

エンコーディング

内部コード(数値)

エンコーディング

- ◆ 多くの種類 (iconv 1.14 というプログラムによれば143種。別名も数えると425種。)
- ◆ 例えば
ASCII, BIG5, GBK, EUC-KR, LATIN1,
MacCyrillic, UTF-8
- ◆ 日本語 (JIS X0208) だけでも
EUC-JP, SHIFT_JIS, CP932, ISO-2022-JP,
ISO-2022-JP-2, ISO-2022-JP-1
- ◆ 正確には, コード表とエンコーディングの差

フォント

内部コード(数値)

35977

MS Pゴシック

YouYuan

貉

貉

32666

MS Pゴシック

YouYuan

羚

羚

70

MS Pゴシック

Wingdings

F



9758

MS Pゴシック



いろいろなフォント

Courier

abcdefghijklmnop

Lucida Sans

abcdefghijklmnop

Times New Roman

abcdefghijklmnop

平成明朝 W7

情報学

Gungshu

情報学

中ゴシック B B B

情報学

HGS教科書体

さいたま

HG明朝B

さいたま

文字とグリフ

文字

グリフ

LATIN SMALL LETTER G

103

g

LATIN CAPITAL LETTER O

79

g

DIGIT ZERO

48

0

いくつかの用語

- ◆ スクリプト:

ラテンアルファベット, キリル文字, アラビア文字,
漢字, ハングル, かな, ...

- ◆ 話し言葉 と 書き言葉

- ◆ 書法

書き言葉を表現するための規則の集合

日本語は, 漢字, かなスクリプトを混用する書法を
標準としている。

漢字とは？

- ◆ 元来が中国語用のスクリプトで、一文字が、中国語の単音節語に対応するという点が最大の特徴。(表意文字ideographではなく、表語文字logographと呼ぶことを主張する人々もいる。)
- ◆ 日本語では、音(当時の中国語の発音を日本語の音韻構造に移したもの)と、訓(中国語から日本語への標準的訳語)を混用しながら表記する書法が発達した。この中で、中国語の単音節語に対応しない漢字(国字)を創作することも行われてきた。

「漢字コード」の歴史

- ◆ 1978: JIS C 6226 情報交換用漢字符号系
 - ❖ 6349漢字(第1水準2965, 第2水準3384)
- ◆ 1983: JIS C6226 改訂 = JIS X 0208 (1987)
- ◆ 1990: JIS X 0212
 - ❖ +5801漢字
- ◆ 1993: ISO/IEC 10646-1 UCS
 - ❖ 21140漢字(統合漢字20902 + 互換漢字238)
- ◆ 2000: JIS X 0213:
 - ❖ + 3685漢字(第3水準1249, 第4水準3390)

「漢字コード」の歴史

- ◆ 2000: ISO/IEC 10646 UCS
 - ❖ + 6583漢字
- ◆ 2003: ISO/IEC 10646 UCS
 - ❖ 計 71226漢字
- ◆ 2011: ISO/IEC 10646 UCS
 - ❖ 計 75616漢字

問題点：何が一つの「漢字」か？

野埜

沢澤澤

搔搔

姫姫

高高

戸戸戸

これらは、ISO10646 UCSでは、すべて別の文字。

問題点 続

葛葛

与与

芸（藝の常用字体）芸（ウン）芸（藝の簡体字）

これらは、ISO10646 UCSでは、同じ文字。

文字学

- ◆ 漢和辞典：諸橋轍次，『大漢和辞典』など
- ◆ 字源学的研究：白川静，『字統』，『字通』など

これらの文字学の成果は，前記の問題点を解決できるか？

野 埜

野

𠩺 後下三・一	埜 克鼎	埜 禽志鼎
埜 說文古文	野 說文・里部	埜 五十二病方 二三五
野 繹山碑	野 睡虎地簡六 ・四五	野 相馬經三 一下
野 武威簡・服 傳二〇	野 漢印	野 白石神君碑
野 魏王基殘碑	野 校官碑	

埜

𠩺 後下三・一	埜 克鼎	埜 古鉢
---------	------	------

yě 《集韻》以者切，上馬以。歌部。

形声 声符は予よ。「説文」十三下に「郊外なり」とあり、

野

のヤシヨ
いなか
いやしい

【野】 11 6712

【埜】 11 4410

【埜】 15 4410

① 野 (11-40133) の古字。

【埜】 5154

【野】 40133

ヤ

ヨ

シヨ

ヨ

ヨ

野 古文 埜

埜 唐韻 野

汉语

字通

大漢和

康熙

姫姫

姫 曉 粹三八六 𠂔 京津五一七 𠂔 師酉簋
 𠂔 作姫簋 𠂔 司寇良父壺 𠂔 魯伯大父簋
 𠂔 陳侯作嘉姫簋 𠂔 吳王光鑑 𠂔 說文·女部
 𠂔 熹·春秋·僖廿五年 𠂔 魏正始石經

《說文》：“姫，黃帝居姫水以為姓。从女，臣聲。”

姫 zhěn 《集韻》止忍切，上軫章。謹慎。《玉篇·女部》：“姫，慎也。”

漢語

姫 集韻 止忍切 音軫 慎也
姫 古文
𠂔

康熙

会意 旧字は女と臣いとに従う。臣*は乳房の象形。

姫 10 〔姫〕 9 ひめ

𠂔 𠂔 𠂔

〔姫〕 6230

𠂔 キ 〔集韻〕居之切 𠂔
 𠂔 イ 〔集韻〕盈之切 𠂔 41-chi
 𠂔 キヨ 〔列子注〕音居

字統

〔姫〕 6229 シン 〔集韻〕止忍切 軫

つつしむ。〔集韻〕姫、慎也。

〔参考〕 もと姫(3-6230)は別字。今、姫の略字として用ひる。〔康熙字典、辨似、二字相似〕姫、音軫、慎也。姫、音基、姓也。

大漢和

姫 用例



光明皇后杜家立成雑書要略
(国会図書館近代デジタルライブラリー) より



大書源
(二玄社)より

高高

高

高 説文・高部
高 老子甲五七

高

高 白石神君碑
陰

【高】 45314
カウ
會〔高、俗作〕高。

【高】 45313
日カウ
日カウ
〔集韻〕居勞切
〔集韻〕居号切
號 豪
KaO'

高 10
コウ(カウ)
たかい・とうとい・すぐれる

高 廣韻古勞切
謂天體也又

高 京の省文と口とに従う。
京は凱旋門。
*けい

高 高 高 高

漢語

大漢和

字統

康熙字典

葛葛



葛

北魏
孫秋生造像記

葛

北魏
元廓平妻王氏墓誌

葛

唐 歐陽詢
九成宮醴泉銘

葛

唐 顏真卿
顏氏家廟碑

葛

唐
玄奘新記明老部

葛

平安 道證寺鐘銘

葛

北宋 黃庭堅
伏波神祠詩卷

葛

明 黃道周

葛

清 陳鴻壽

葛

清 沈友聲



【葛】

13

4472

くカツ
かず
かづら

字通

北越奇談
(早稲田大学古典籍
総合データベース)より

大書源
(二玄社)より

文字学は

- ◆ 残念ながら、漢字の文字という単位の識別に対して一貫した答えは与えてくれない。
- ◆ 白川文字学は、その中で文字単位に対して一貫した立場と用語を貫いているが、「野」と「埜」とが同じ文字の異なる字形であるということになると、実用的立場からは受け入れ難い、という意見が多数である。
- ◆ 実際に書かれている文字は、字書には示されていない字体も多く使われている。

異体字

- ◆ 昔から「古文」、「古字」、「俗字」、「異体」などさまざまな呼び方で広く認められてきた「異なるグリフ」
- ◆ 楷書形への写し方の違いによって生まれる「異なるグリフ」
- ◆ 漢字簡素化など、書法を変更することによって生まれる「異なるグリフ」

字体の「包摂」

- ◆ 異なる字体(グリフ)を一つの文字(コード値)としてみとめること。
- ◆ UCSでは, 付属書Sで, いわゆる unification rule(統合規則)を決めている。
- ◆ JIS X 0208は, 1997版から包摂基準を持っている。

Unification ルール (Unifyの例示)

迂・迂・迂	示・示・示	艮・艮・艮	食・食・食	黄・黄	盥・盥	曷・曷
包・包	青・青	每・每	册・册	争・争	盩・盩	曷・曷
步・步	者・者	臭・臭	并・并	骨・骨	呂・呂	录・录
鼎・鼎	吳・吳・吳	眞・眞・眞	爲・爲	单・单	曾・曾	直・直
專・專	内・内	晉・晉	龜・龜	++		成・成

Unification ルール (Separateの例示)

扌・擴	策・筭	𠂇・𠂇	𠂇・𠂇	𠂇・𠂇	区・區	夾・夾
單・單	雀・雀	戈・戈	𠂇・𠂇	𠂇・𠂇	𠂇・𠂇	間・間
朶・朶	雋・雋	恒・恆	奘・奘	𠂇人・𠂇人	𠂇朶・𠂇朶	𠂇𠂇・𠂇𠂇

JIS X0208-1997の包摂基準(部分)

- a) 方向・曲直などの点画の性質による違い
— 入りの左右

包摂する部分字体

王 壬 壬

耒 耒

室 宰 宰

圭 圭

戸 戸 戸

適用除外(規定)

18-06 31-49
王, 壬

包摂の問題点

- ◆ 包摂規則より先にいくつかの漢字コード表があり、包摂規則を崩している場合がある。
- ◆ 形だけからの包摂規則は、字源から見れば問題となる場合がある。
- ◆ 言語におけるあらゆる規則と同じく、万人の納得する包摂規則は不可能であり、多数が認めるものとして決めざるを得ない。
(例えば、一つの単語という単位でも同様の問題はあある。)

結論として 1

- ◆ 文字は形ではない。

「康熙字典の[一]部には丁, 丂, 七...を録する。...字はその構造的な原理から離れ, その構造的な意味も捨てられて, ただ筆画の形式によって分離配列されている。そこにあるものはすでに文字ではなく, 文字の形である。意味を失っている記号である。」

— 白川静, 「字統の編集について」より

結論として 2

- ◆ しかしながら、漢字についてはどの単位で文字と考えるか、の社会的合意が弱い。

現在のISO10646の包摂で不都合な場合は、アプリケーションソフト側で、（複数の文字を一つの文字として扱うための）正規化をおこなうか、一つの文字の中で異なるグリフを区別するためのフォントの切り替えをおこなう方法がある。

結論として 3

- ◆ 字形のわずかな違いを捉えて、異なる漢字である、としたがる社会的風潮は、データ作成のコストを上げ、知的資源の無駄遣いとなるばかりである。

少なくとも、昭和初期くらいまでの日本人は、字体、字形の異なりに寛容であった。字体の包摂は、明文化しなくとも常識として通用していた。

解題

2011-10-05

市民講座