

文字コードの変遷

一般社団法人 文字情報技術促進協議会 理事・事務局長
一般社団法人情報処理学会 情報規格調査会 規格役員 (ISO/IEC JTC1 NB)
ISO/IEC JTC1 SC2 (Coded character sets) 委員

田丸 健三郎

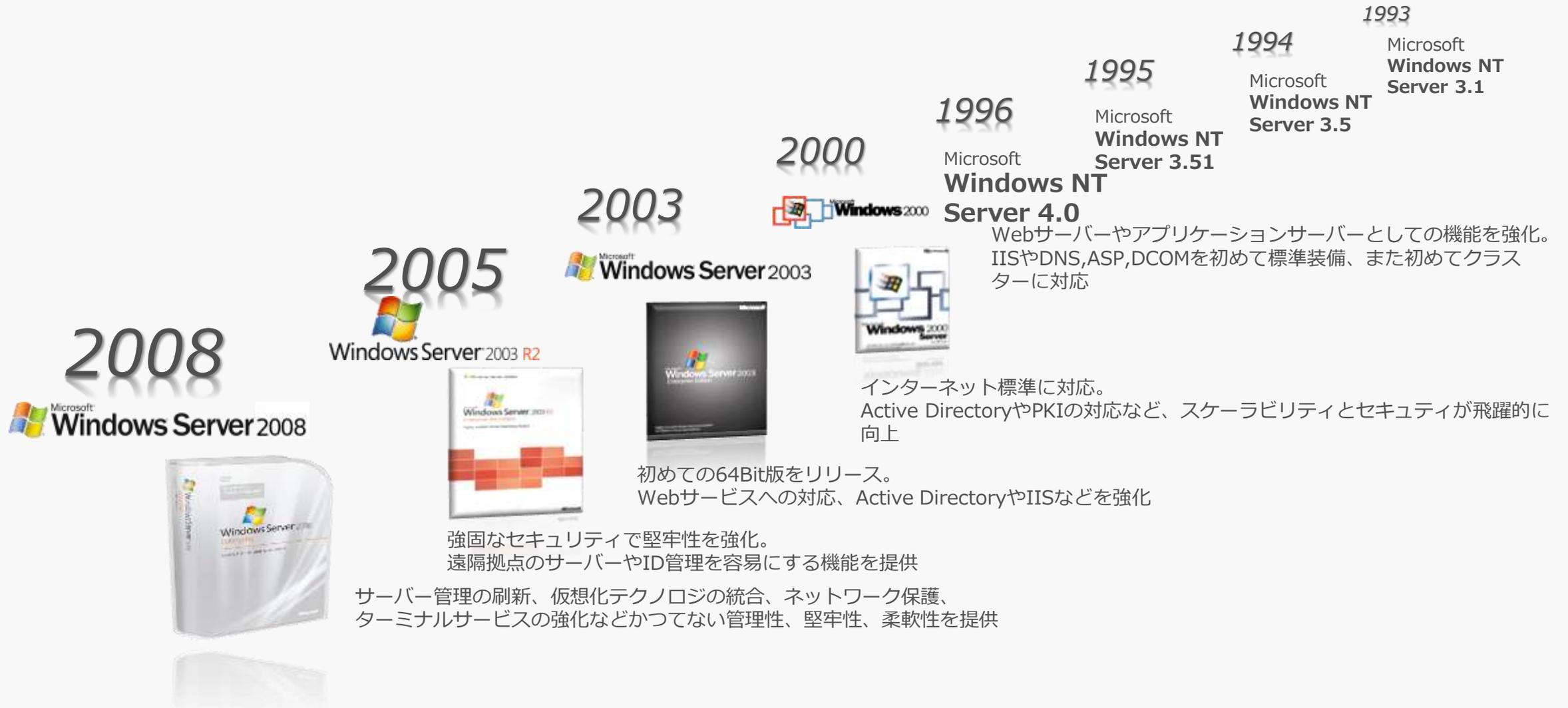
E-mail: k.tamaru@moji.or.jp

シフトJIS、マイクロソフト標準キャラクタセット、JIS X
0213からUnicode サロゲートペアまで



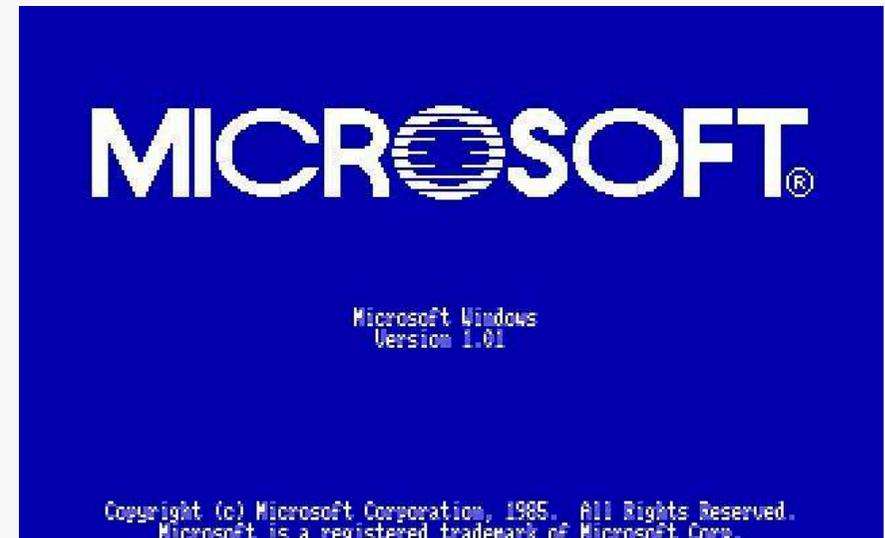
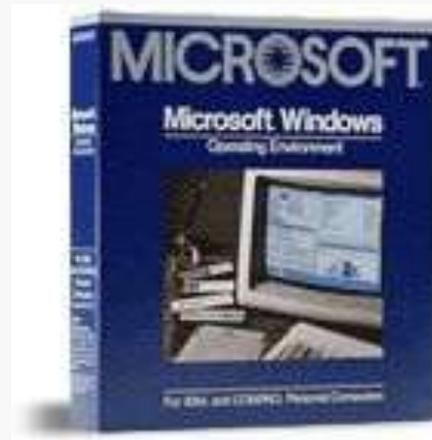
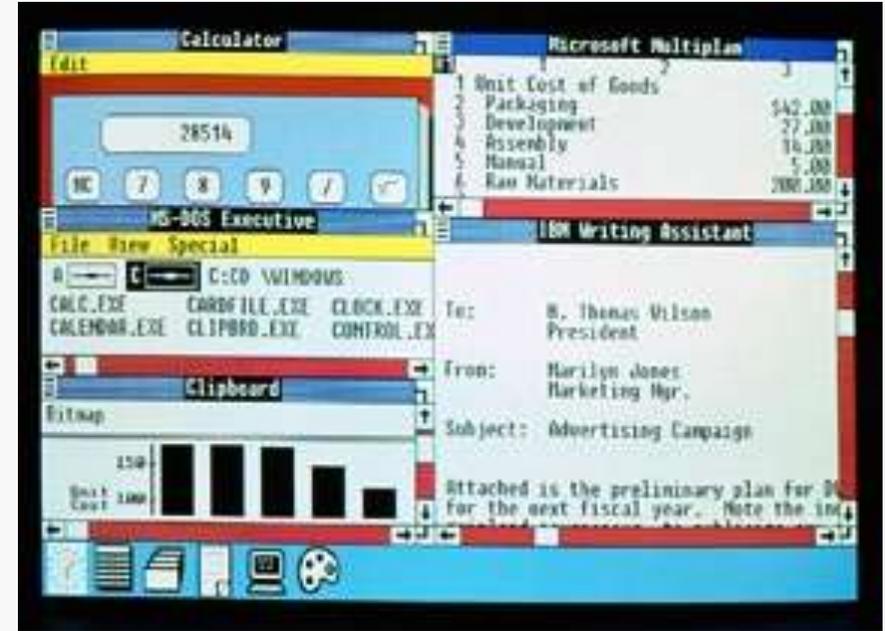
文字コードの変遷

Windows の変遷



1985年 - Windows 1.0

- 16bit OS
- 初のGUI
- アプリケーションの同時実行

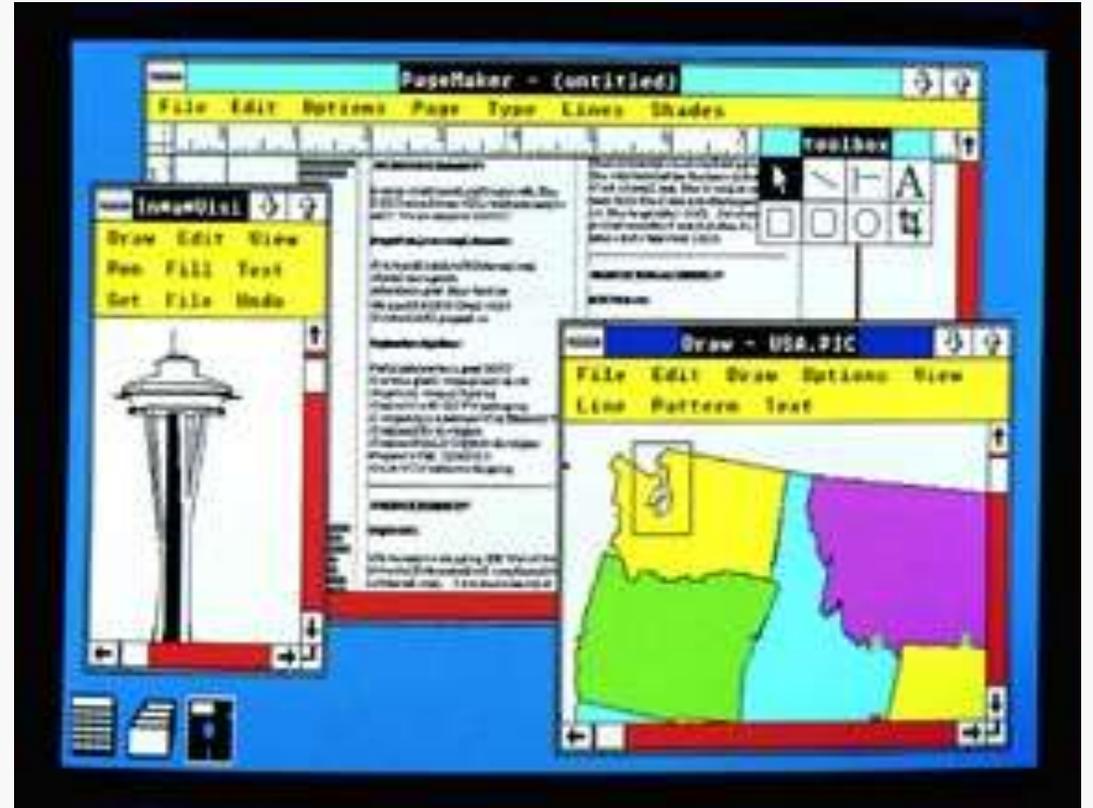


1987年 - Windows 2.0

- キーボードショートカット
- i286プロセッサの採用し物理メモリ空間が**16MB**に拡張。1GBの仮想記憶も使用可能に。
(多くのPCは1MB未満)
- DDE (Dynamic Data Exchange)
- ウィンドウの重なり

CPUは、フラグレジスタのbit15-12で判定

- 1111bから変更可：8086世代
- 0000bから変更不可：80286
- それ以外：80386以降



WIFE (Windows Intelligent Font Environment)

DBCSフォントを扱うできるようにするために開発された『インストール可能なフォント・ドライバを実現するサブシステム』

1990 2月5日 DDS詳細仕様書 V0.1

1990 6月7日 WIFE詳細仕様書 V0.2

1990 6月27日 WIFE詳細仕様書 V0.21

外字等の仕様を追加

1990 7月2日 WIFE詳細仕様書 V0.31

文字の回転（縦書）の仕様を追加

W I F E : Windows Intelligent Font Environment

詳細仕様書

Version 1.10

1991年7月24日

マイクロソフト株式会社

1 概要

1-1 ドキュメントの定義

本ドキュメントは、MS-WINDOWS 3. Xにおいて、インストール可能なフォント・ドライバを実現するサブシステムの仕様書です。

これは、互いに関連する一群のモジュールによって構成される一つのフォント・サポートの環境であり、この環境の総称を、WIFE (Windows Intelligent Font Environment) と呼びます。

WIFE環境では、SBCS文字とDBCS文字の区別なくフォントを利用でき、いままで統一されたインストール可能なフォント・システムを持たなかったDBCS圏でインストール可能なシステムが利用できるようにするものです。

WIFEは、欧米のSBCS圏のピュアなMS-WINDOWS環境に、いくつかのモジュールを付加することによって実現されます。このドキュメントでは、これらの付加されるモジュールと、それによって実現されるWIFE独自の環境のみに関して記述しています。そのため、このドキュメントを読むためには、MS-WINDOWSのGDIモジュールとデバイス・ドライバとのインターフェース、及び、デバイス・ドライバの文字列出力機能の詳細に関する知識を必要とします。

1-2 開発の目標

WIFE環境の目標とするのは、次の二つです。

- 1) DBCSフォントを扱うことができる
- 2) オペレーティング・システムのレベルでサポートする
- 3) ファイル及びそれ以外の任意の物理メディアによるフォントの提供を可能とする
- 4) 任意の論理フォーマットによるフォントの提供を可能とする

これらの目標を達成するために、インストール可能なフォント・ドライバという概念を導入します。フォント・ドライバとは、フォントデータを読み込み、システムが理解できる標準形式のビットマップに変換するプログラムです。このプログラムをユーザーが自由にインストールできるようにすることにより、任意のメディアかつフォーマットによるフォントを利用できるようにします。フォント・ドライバはフォントデータの提供者が作成し、フォントデータとともに配布することを基本とします。

このフォント・ドライバを使用可能とするためのサポートモジュール群を含めたシステム全体をWIFE環境と呼びます。

このWIFE環境を達成するにあたって、次の副目標を設定します。

- a) WINDOWSアプリケーションは、WIFE環境をサポートするために一切の変更を必要としない。(同じアプリケーションをWIFE環境の有無にかかわらず使用できること)
- b) WINDOWSカーネル・モジュール群(KERNEL、USER、GDI)は、WIFE環境の追加によって、一切の変更を必要としない。
- c) WINDOWSデバイス・ドライバ(従来から存在するディスプレイ及びプリンターのドライバ)は無変更でWIFE環境で使用できること。その際には、パフォーマンスの低下は容認される。
- d) 前項と関連して、WINDOWSデバイス・ドライバ(ディスプレイ及びプリンターのドライバ)は、WIFE環境に適応するように修正を加えることによって、パフォーマンスを改善できる。
- e) WINDOWSシステムの持つ既存の全てのフォントを使用できること。

f) 将来的にWINDOWSシステムに組み込まれる可能性のあるフォント・システムと可能な限り共存できるよう配慮する。

g) 新規に追加されるモジュールであるフォント・ドライバは、OS/2のフォントシステムと可能な限り互換性を持たせる。(ソースレベルの互換性)

h) コントロールパネルは、WIFE独特の機能の追加に対応して変更を加える。

i) 特定の言語(日本語等)に依存しない。

j) 一つのフォント・データで、一つのキャラクターセットの文字を提供する。一つのフォントを、SBCSとDBCSで別々のデータとして提供するようなことはしない。

以上。

1-3 互換性

本節では、欧米のSBCS圏のピュアなMS-WINDOWS環境に対する、WIFE環境を含んだMS-WINDOWSシステムの互換性について述べます。

この二つの環境の関係は、以下の図で示すように、ピュアなMS-WINDOWSの全てを含み、それを外側から包むような形で、WIFE環境が存在します。そのため、WIFE環境では、ピュアなMS-WINDOWS環境の機能は全て使用できます。

ですから、WIFE環境は、ピュアなWINDOWS環境の上位環境であると位置づけられます。



なお、上の図は、フォントの取り扱いの概念について説明したものです。

実際には、USバージョンなどのピュアなMS-WINDOWSに、モジュールを追加することで、WIFE環境を実現することは不可能です。なぜなら、ピュアなMS-WINDOWS環境を構成するモジュールには、DBCSキャラクターを正しく扱うための能力がないためです。DBCSフォントをインストール可能にすることと、DBCSキャラクターを正しく扱うことは、全く別の問題です。本仕様書は、前者の「DBCSフォントをインストール可能にすること」についてのみ、対象とします。

1-3-1 APIの互換性

アプリケーション・プログラムから呼び出すことができるAPIには、一部の例外を除き、ピュアなMS-WINDOWS環境に対して、追加や変更はありません。

つまり、既存のAPIによって、WIFE環境の全ての機能が利用可能です。

アプリケーションから見ると、WIFEフォントは、デバイス・ドライバに固有のデバイス・フォントと同等に見えます。そのため、アプリケーションは、ピュアなMS-WINDOWSが持っている全てのデバイス・フォントを利用するためのAPIを利用し、WIFEフォントにアクセスすることができます。もちろん、本来のデバイス・フォントも平行して利用できます。

相互運用性の喪失

各ベンダーによりJISが独自拡張されたことにより、異なるメーカー間での文字情報の交換・共有が困難になる。



「マイクロソフト標準キャラクタセット」

各ベンダーによる独自拡張された文字を整理統合し、異機種・プログラム間の文字情報の相互運用を可能に。

Microsoft Windows™ version 3.1
日本語版

マイクロソフト標準キャラクタセット
仕様書

Revision 1.0
Dec 25, 1992

Microsoft Co., Ltd.
Systems R&D

図 A-1 DBCS コードレイアウト

ShiftJIS	JIS (区)		
8140	1	JIS 非漢字 X0208 JIS90	
:	:		
8440	7		
:	8		
:	:	未使用 (376 文字)	
8740	13	NEC 特殊文字 (83 文字)	
:	:	未使用 (188 文字)	
8840	15	JIS 漢字 X0208 JIS90	
:	16		
:	:		
:	:		
:	:		
EA40	83		
:	84		
EB40	85		未使用 (376 文字)
:	:		
ED40	89		
:	90	NEC IBM 選定文字 (374 文字)	
EE40	91		
:	92		
EF40	93	未使用 (188 文字)	
:	94		
FO40	95	外字 (1880 文字)	
:	:		
F940	113		
:	114		
FA40	115	IBM 拡張文字 (388 文字)	
:	:		
FC40	119		
:	:		

US-ASCIIの穴を縫って作られたシフトJIS

文字を指すポインタが文字の先頭か中間か容易に判別できない。

E n g l i s h	と	日	本	語
45 6E 67 9C 69 73 68	82C6	93FA	967B	8CEA



0xC6は文字の先頭か中央か？

JIS X 0201 / JIS X 0208 / JIS X 0213	第1バイト	第2バイト
1バイト文字 ASCII (SBCS / Single Byte Character Set)	0x21 ~ 0x7E	
1バイト文字 半角カタカナ (SBCS / Single Byte Character Set)	0xA1 ~ 0xDF	
2バイト文字 JIS X 0208 / JIS X 0213(1面) (DBCS / Double Byte Character Set)	0x81 ~ 0x9F 0xE0 ~ 0xEF	0x40 ~ 0x7E 0x80 ~ 0xFC
2バイト文字 ユーザー定義文字 (DBCS / Double Byte Character Set)	0xF0 ~ 0xFC	0x40 ~ 0x7E 0x80 ~ 0xFC
2バイト文字 JIS X 0213(2面) (DBCS / Double Byte Character Set) JIS X 0213にて規定されていますが、 Windowsをはじめとする多くのプラットフォーム、 ツールではJIS X 0213の文字としては認識せず、 外字として認識します。	0xF0 ~ 0xFC	0x40 ~ 0x7E 0x80 ~ 0xFC

実装コストが高く、しばしばバグの原因に

- 文字列のポインタをそのまま使用するバグる
- 1文字前のポインタを取得するだけでも計算資源を浪費する



文字列の先頭まで辿らないと文字の先頭が「二」(0xC6)か判別できない

```
LPSTR SzAlign(LPSTR szStart, LPSTR szSource)
{
    LPSTR sz;
    CPINFO cpinfo;

    if (!szStart || !szSource)
        return NULL;

    if (szStart == szSource)
        return szStart;

    sz = szSource;
    // 以下のコードでは、ポインタが確実にDBCSの第1バイトでない
    // 文字が現れるまでポインタを前方移動します。これは、DBCSの
    // 第1バイトと第2バイトの範囲に重複があるためです。
    for (;;)
    {
        if ( sz==szStart || !IsDBCSLeadByte(*(sz-1)) )
            return (LPSTR) ( szSource -
                (((ULONG_PTR) sz ^ (ULONG_PTR) szSource) & 1) );
        sz--;
    }
}
```

0x5C 問題

バックslashはデリミタとして広く使用されている

\ <Folder> \ <FileName>

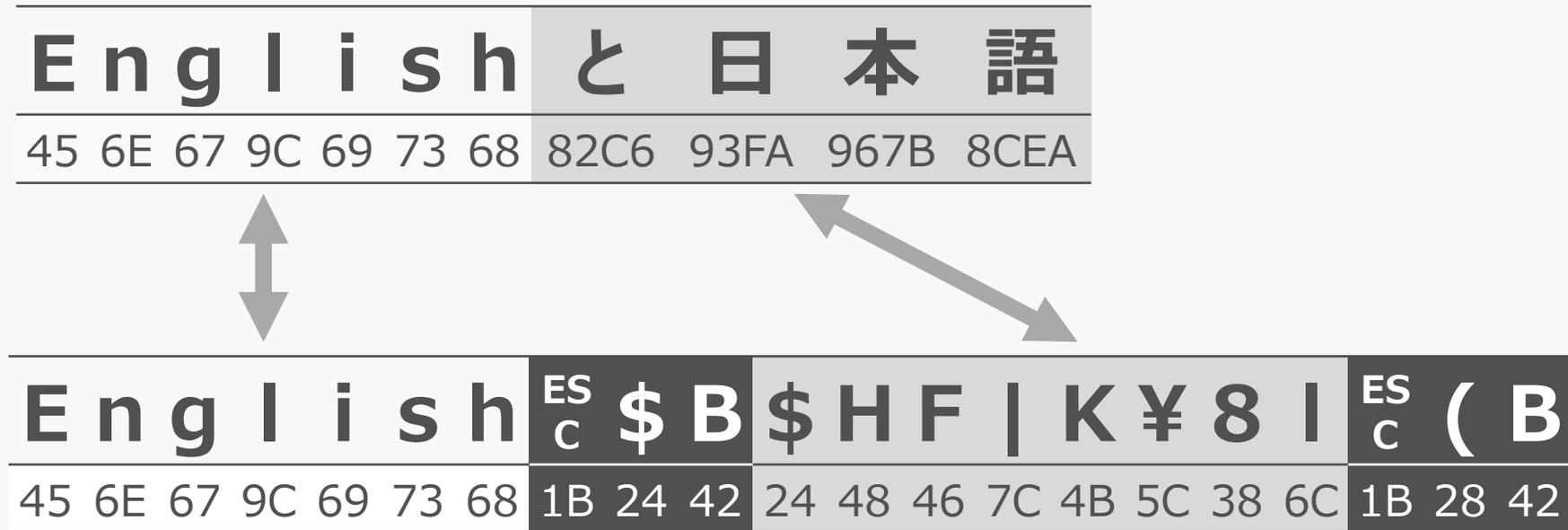
同じコードが日本語環境だと

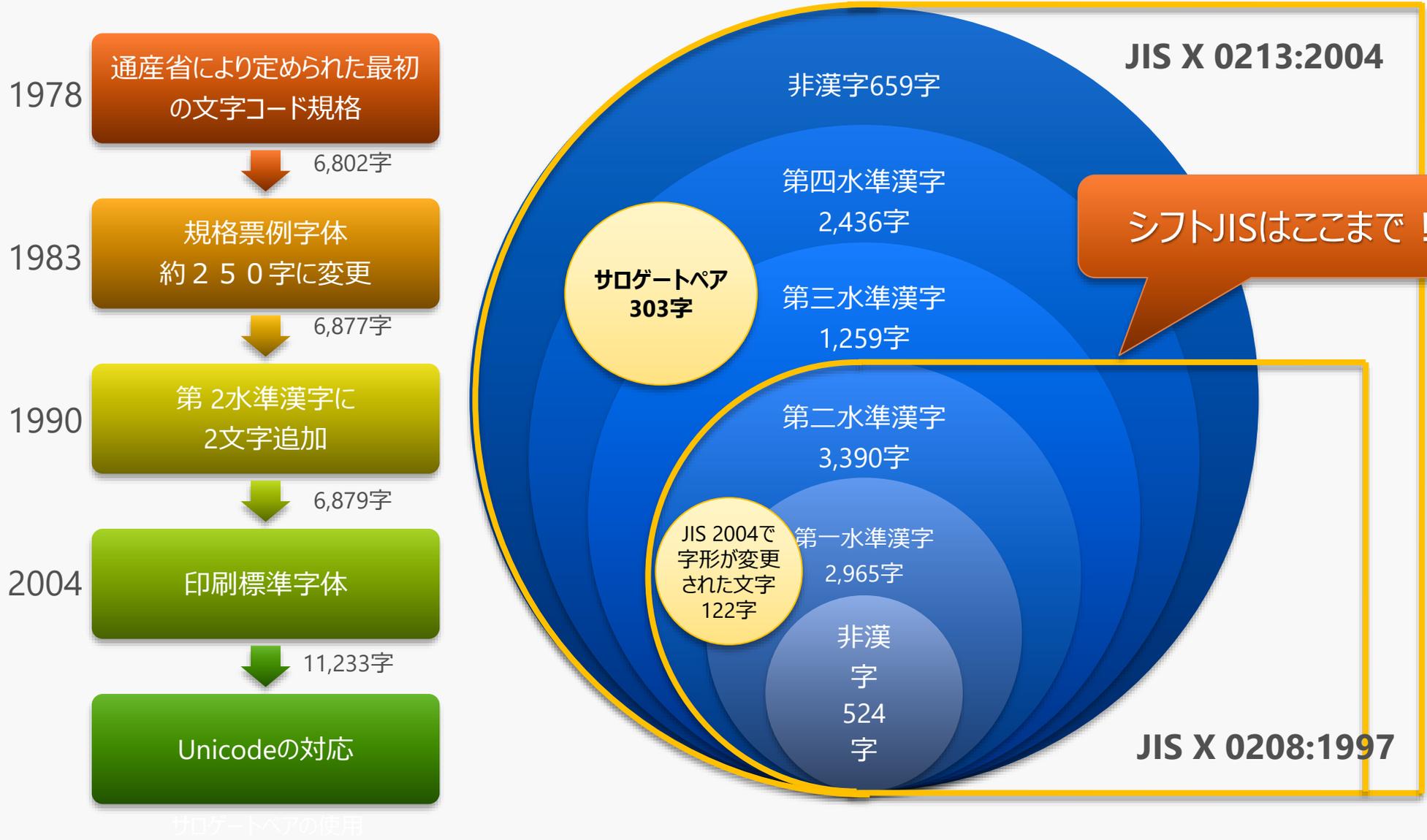
¥

JIS X 0201 / JIS X 0208 / JIS X 0213	第 1 バイト	第 2 バイト
1 バイト文字 ASCII (SBCS / Single Byte Character Set)	0x21 ~ 0x7E	
1 バイト文字 半角カタカナ (SBCS / Single Byte Character Set)	0xA1 ~ 0xDF	
2 バイト文字 JIS X 0208 / JIS X 0213 (1 面) (DBCS / Double Byte Character Set)	0x81 ~ 0x9F 0xE0 ~ 0xEF	0x40 ~ 0x7E 0x80 ~ 0xFC
2 バイト文字 ユーザー定義文字 (DBCS / Double Byte Character Set)	0xF0 ~ 0xFC	0x40 ~ 0x7E 0x80 ~ 0xFC
2 バイト文字 JIS X 0213 (2 面) (DBCS / Double Byte Character Set) JIS X 0213 にて規定されていますが、 Windows をはじめとする多くのプラットフォーム、 ツールでは JIS X 0213 の文字としては認識せず、 外字として認識します。	0xF0 ~ 0xFC	0x40 ~ 0x7E 0x80 ~ 0xFC

メールで一般的だった 7bit JIS / ISO-2022-JP

メールの送信ではエンコード、受信ではデコード処理が必要。
当時は方言、送信側バグなどへの対応により高コストの象徴







Unicodeとサロゲートペア

Windows 95 (1995)

- 32bit OS
- Intel Pentium Pro 150~200MHz
- **平均4~16MB メモリ**
- MP (Multi Processor)



二つの議論

第一世代 (1983 -) Shift JIS

(JIS78 or JIS83) + メーカー拡張

MS-DOS 2 - , MS OS/2 1.x, MS-Windows 2.x, 3.0

1993年夏
Windows 95
でUCS ?

第二世代 (1993 -) マイクロソフト 標準 キャラクタ セット

JIS90, 10646

Windows 3.1, NT 3.1/3.5, 95, NT 3.51/4.0

第三世代 (1998 -) 補助漢字をUCS拡張

補助漢字

Windows 98, NT 4.0 SP4, 2000, Me, XP

2000年冬
BMP か
Surrogateか ?

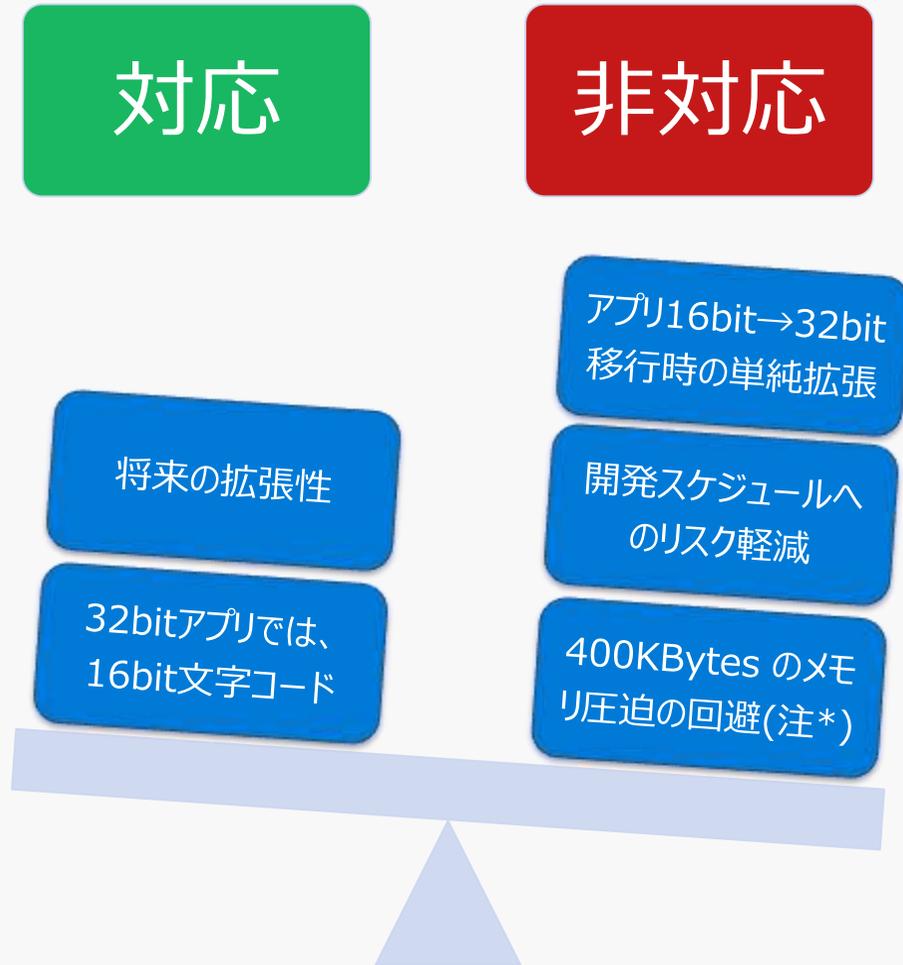
第四世代 (2007 -) 国語審議会答申に基づく最新規格への対応

JIS2004

Windows Vista, Windows Server
Windows XP (with JIS2004 Pack)

Windows 95 での UCS 対応

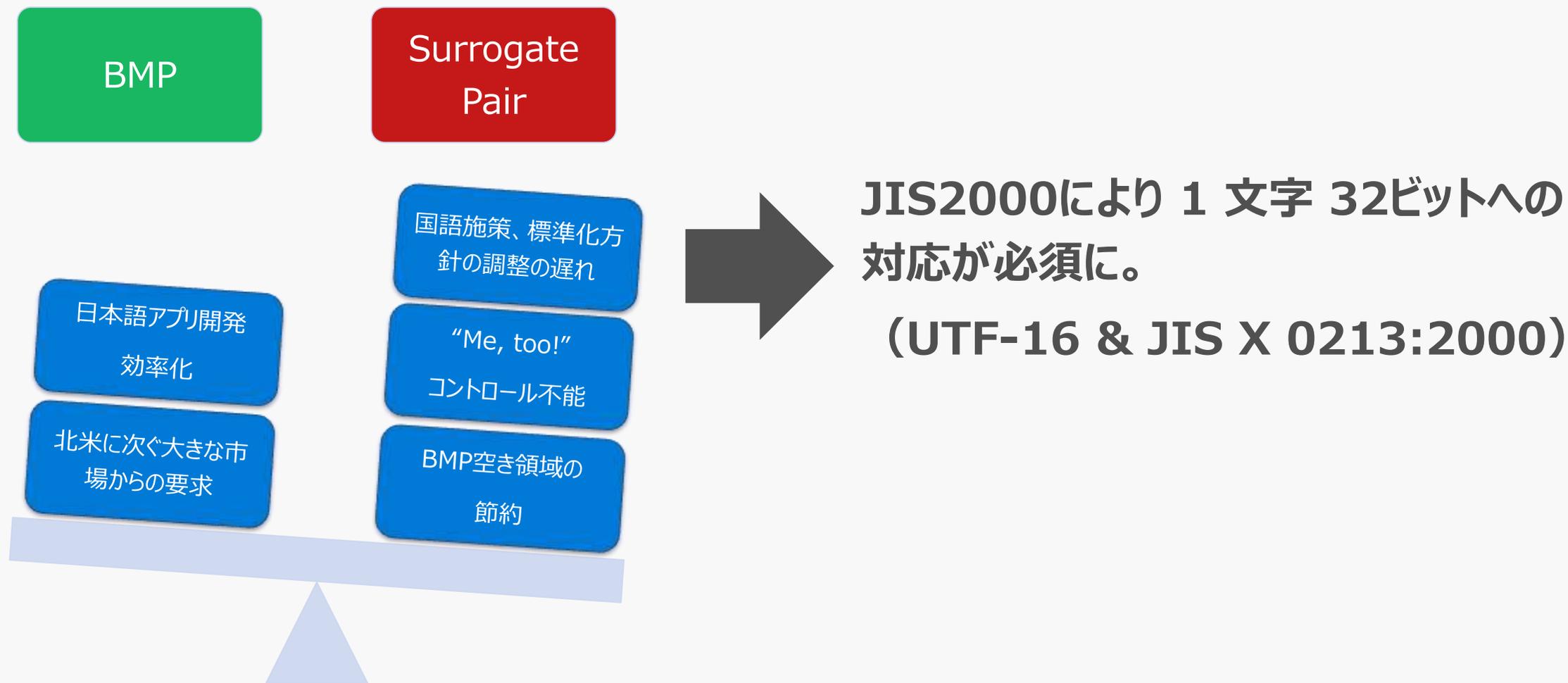
1993年夏 @ 米国 ワシントン州



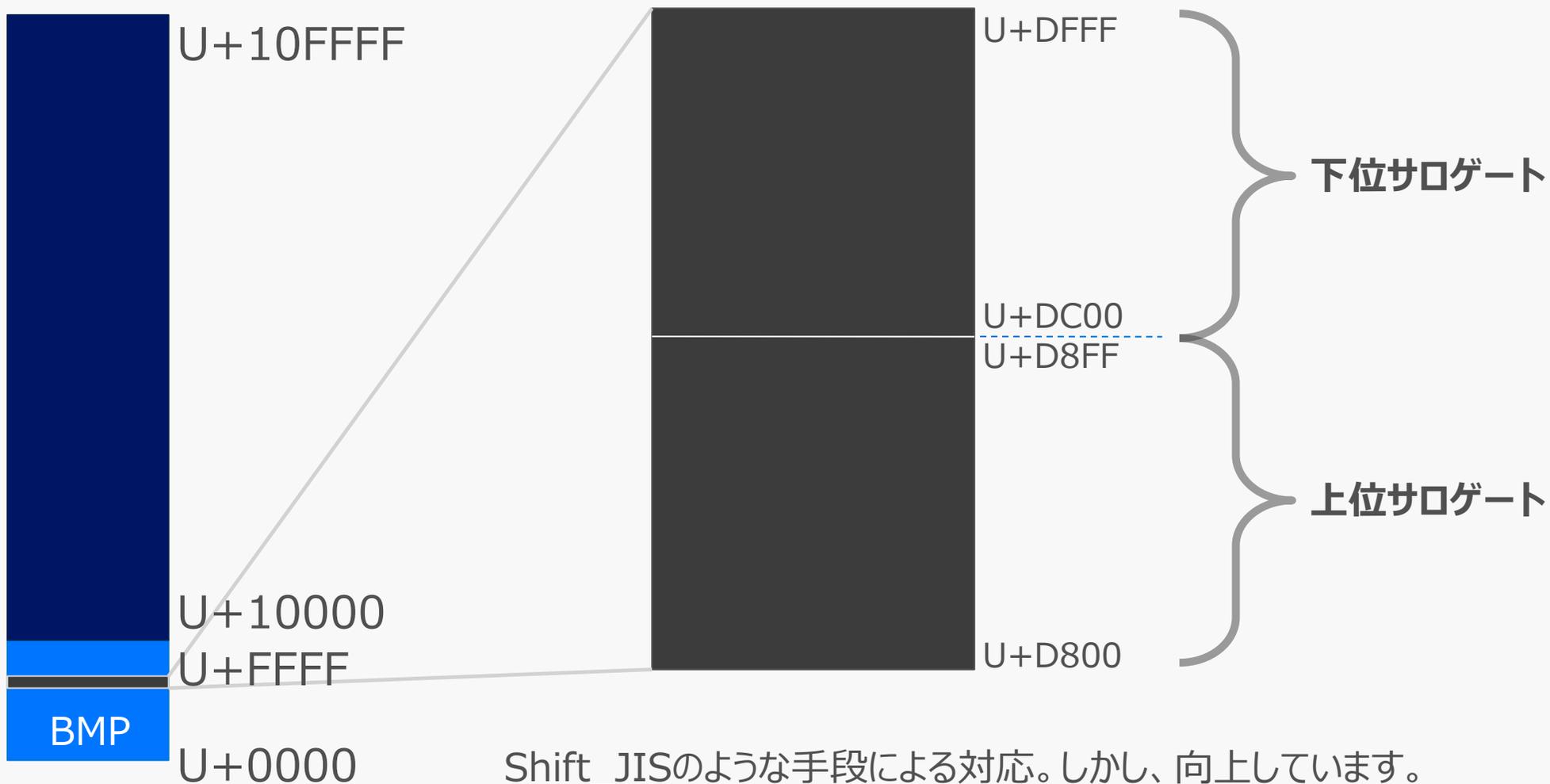
注*: 2 ~ 4Mbytes の必要メモリを想定

JIS2000 の 303文字の行方

2000年 冬 @ 米国 カリフォルニア州、中国 北京

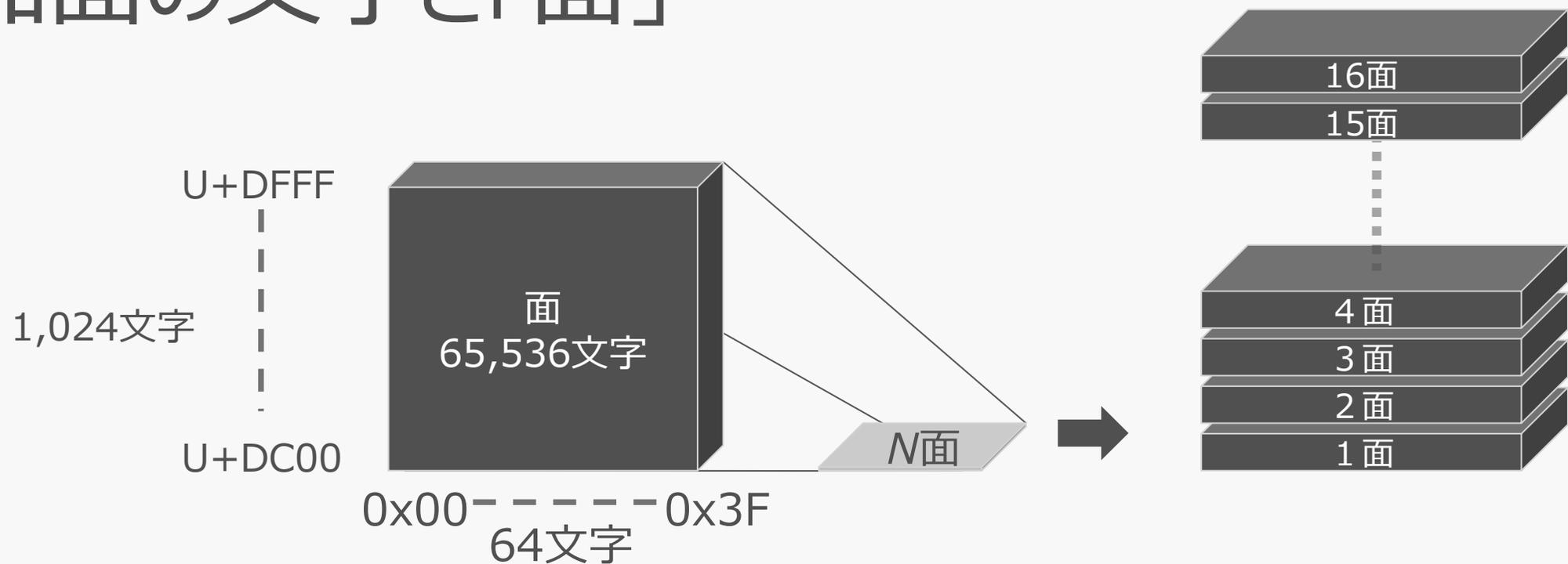


追加面の文字 (サロゲートペア)



Shift_JISのような手段による対応。しかし、向上しています。
値から上位サロゲート、下位サロゲートの判断が可能。

追加面の文字と「面」



「	丈	土	壑	」	は	追	加	面	の	文	字
300C	D840 DC0B	D844 DE3D	D844 DF1B	300D	306F	8FFD	52A0	9762	306E	6587	5B57

上位サロゲート (U+D800 ~ U+DBFF)	X	下位サロゲート (U+DC00 ~ U+DFFF)	=	合計拡張文字数
1,024		1,024		1,048,576文字

JIS2004における大きな変更

- 字形の変更
- 追加面の文字（サロゲートペア）の使用
- Windows Vista, Windows Server 2008 はJIS 2004の字体をUIに使用。
- アプリケーションはJIS90, JIS04両方を使用可能(OpenType)

JIS X 0213:2004 Windows Vista Windows Server 2008以降	JIS X 0208:1990 Windows XP以前
味 噌	味 噌
葛 飾 区	葛 飾 区
祇 園	祇 園
進 捗	進 捗
噂	噂
芦 屋	芦 屋
飴	飴

Windows Vista, Windows Server 2008以降に搭載されているOpenTypeフォントは、両方の字形をサポート

書籍では

	使用されている字形数	Unicode		該当無
		Plain	IVD(AJ1-6)	
萬葉集	3,531	99.4%	99.8%	0.2%
梶井基次郎全集	4,563	97.3%	99.8%	0.2%
中島敦全集	3,543	97.6%	99.8%	0.2%

- JISでサポートされている文字（約2,000）だけでは書籍で使用されている字形全てを表現することができない。
- 約97%以上はシフトJIS（CP932）に含まれる
- Unicode、AJ1-6で約99.8%を網羅

平成22年度書籍等デジタル化推進事業

（デジタル・ネットワーク社会における出版物の利活用推進のための外字・異体字利用環境整備事業）

Unicodeの文字数

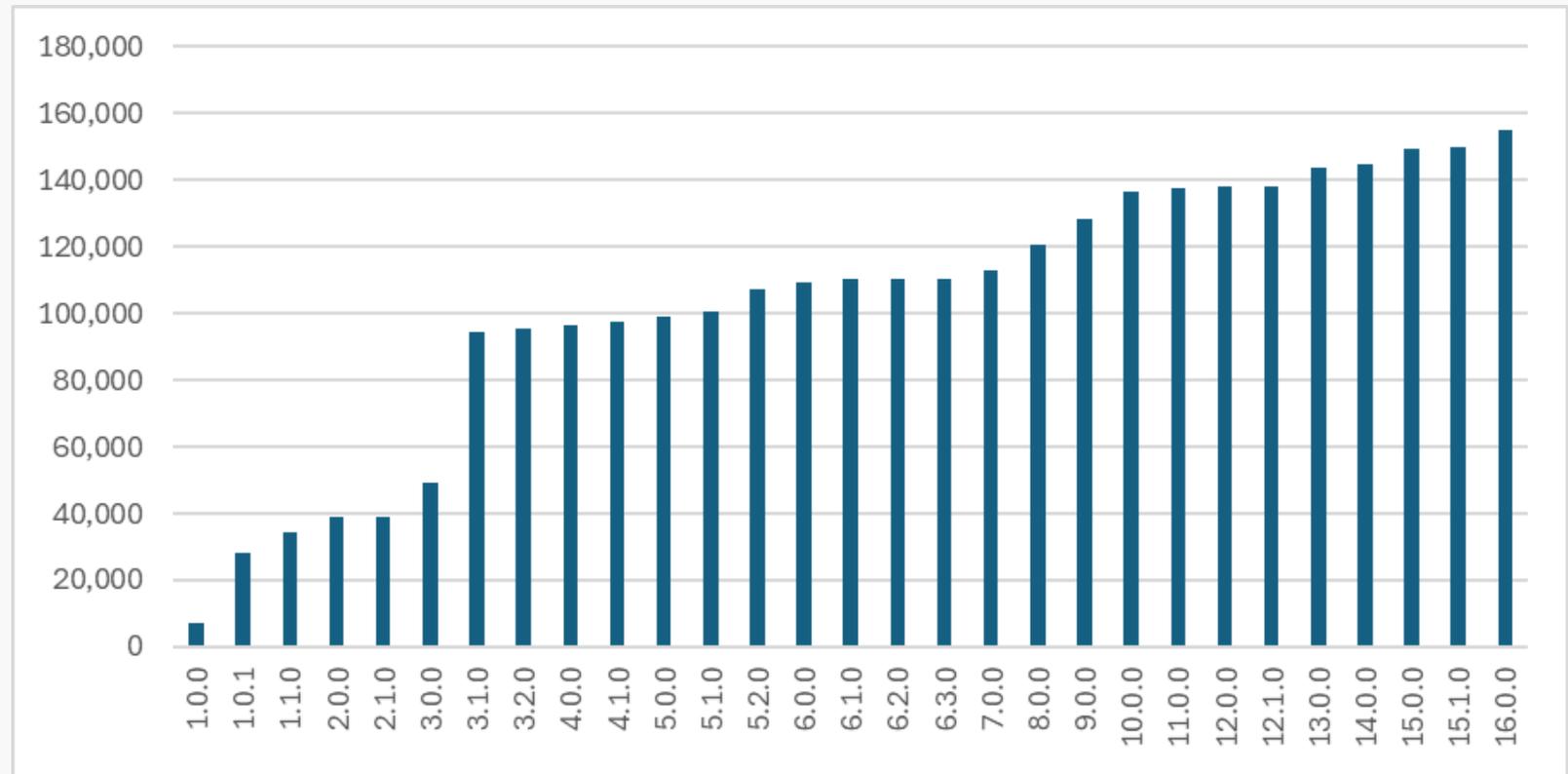
Unicodeバージョン	制定年	文字数	詳細
1.1	1992	34,233	<ul style="list-style-type: none">• JIS X 0208とJIS X 0212を含む Unicode のバージョン
2.0	1996	38,950	<ul style="list-style-type: none">• サロゲートペアを技術仕様として採用 (この時点では文字は未定義であり、3.1にて実装)• ハングル文字の移動 (Unicode 1.1 と互換性消失)• (技術仕様としては、JIS X 0213:2004 に対応)
2.1	1998	38,952	<ul style="list-style-type: none">• ユーロ通貨記号追加、多数の記号定義変更
3.0	1999	49,259	<ul style="list-style-type: none">• CJK 統合漢字拡張 A、漢字 6,582 文字追加
3.1	2001	94,205	<ul style="list-style-type: none">• サロゲートペア 303 文字を追加• JIS X 0213:2000 一部対応、言語タグ追加• CJK 統合漢字拡張 B ブロック追加• CJK 統合漢字拡張 B、漢字 42,711文字追加
3.2	2002	95,221	<ul style="list-style-type: none">• JIS X 0213:2000 および JIS X 0213:2004 に正式対応、• 異体字セレクタ 1 ~ 16 追加• (JIS X 0213:2004 の追加 10 文字は、すでに存在)• CJK 互換漢字ブロックに追加された JIS X 0213:2000 漢字の 59 文字および追加丸付き数字 (~ ⑤) などの非漢字を追加
4.0.0	2003	96,447	<ul style="list-style-type: none">• 異体字セレクタ 17 ~ 256 追加
5.0.0	2006	99,089	<ul style="list-style-type: none">• BMP(基本多言語面) 領域にバリ文字など追加• サロゲート領域にフェニキア文字など追加
6.0	2010	109,449	<ul style="list-style-type: none">• ISO/IEC 10646:2010• 絵文字の追加



それでも足りない文字

Unicodeの登録文字数の変遷

制定年月日	バージョン番号	収録文字数
1991年10月	1.0.0	7,161
1992年6月	1.0.1	28,359
1993年6月	1.1.0	34,233
1996年7月	2.0.0	38,950
1998年5月	2.1.0	38,952
1999年9月	3.0.0	49,259
2001年3月	3.1.0	94,205
2002年3月	3.2.0	95,221
2003年4月	4.0.0	96,447
2005年3月31日	4.1.0	97,720
2006年7月14日	5.0.0	99,089
2008年4月4日	5.1.0	100,713
2009年10月1日	5.2.0	107,361
2010年10月11日	6.0.0	109,449
2012年1月31日	6.1.0	110,181
2012年9月26日	6.2.0	110,182
2013年9月30日	6.3.0	110,187
2014年6月16日	7.0.0	113,021
2015年6月17日	8.0.0	120,737
2016年6月21日	9.0.0	128,172
2017年6月20日	10.0.0	136,690
2018年6月5日	11.0.0	137,374
2019年3月5日	12.0.0	137,928
2019年5月7日	12.1.0	137,929
2020年3月10日	13.0.0	143,859
2021年9月22日	14.0.0	144,697
2022年9月13日	15.0.0	149,186
2023年9月12日	15.1.0	149,813
2024年9月10日	16.0.0	154,998



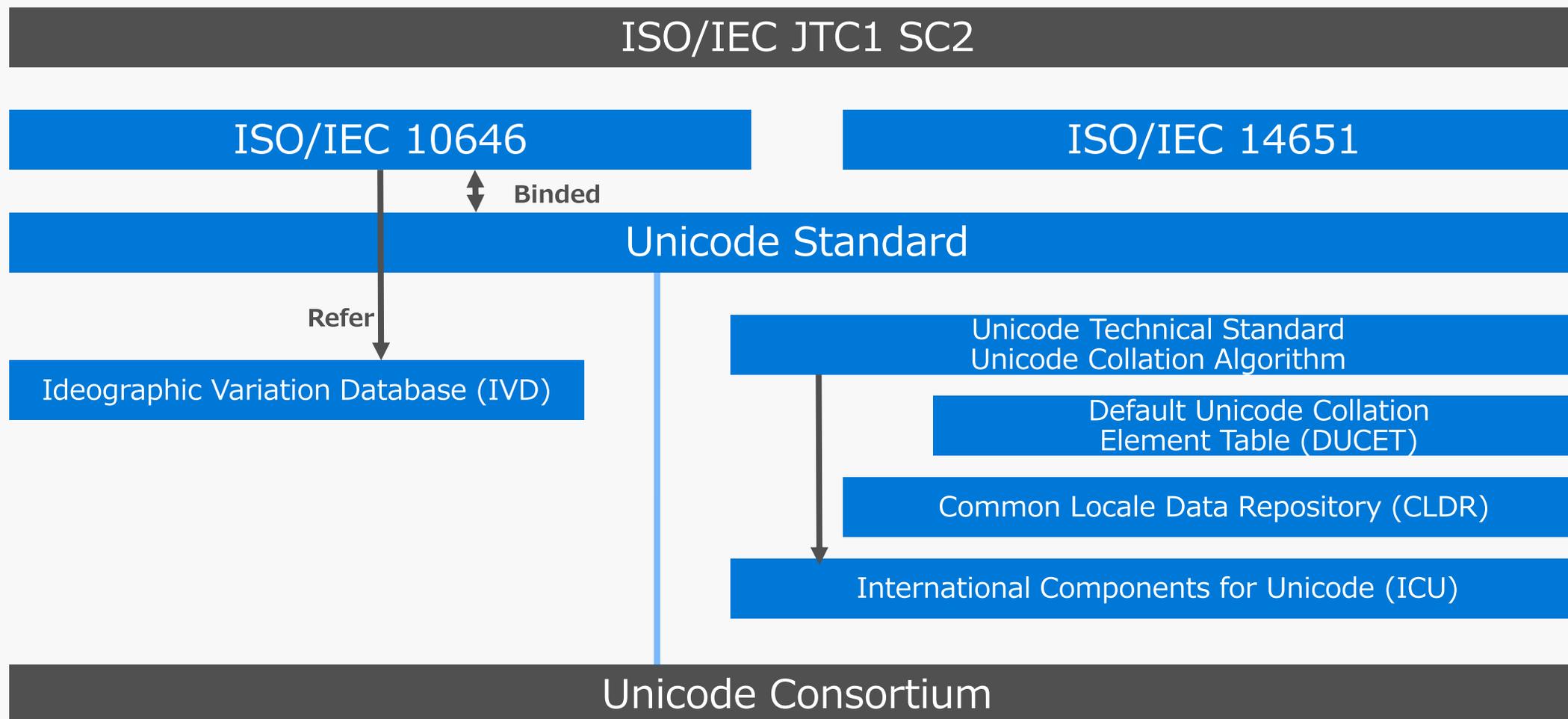
基本多言語面 (Basic Multilingual Plane) 追加面 (Supplementary Plane)



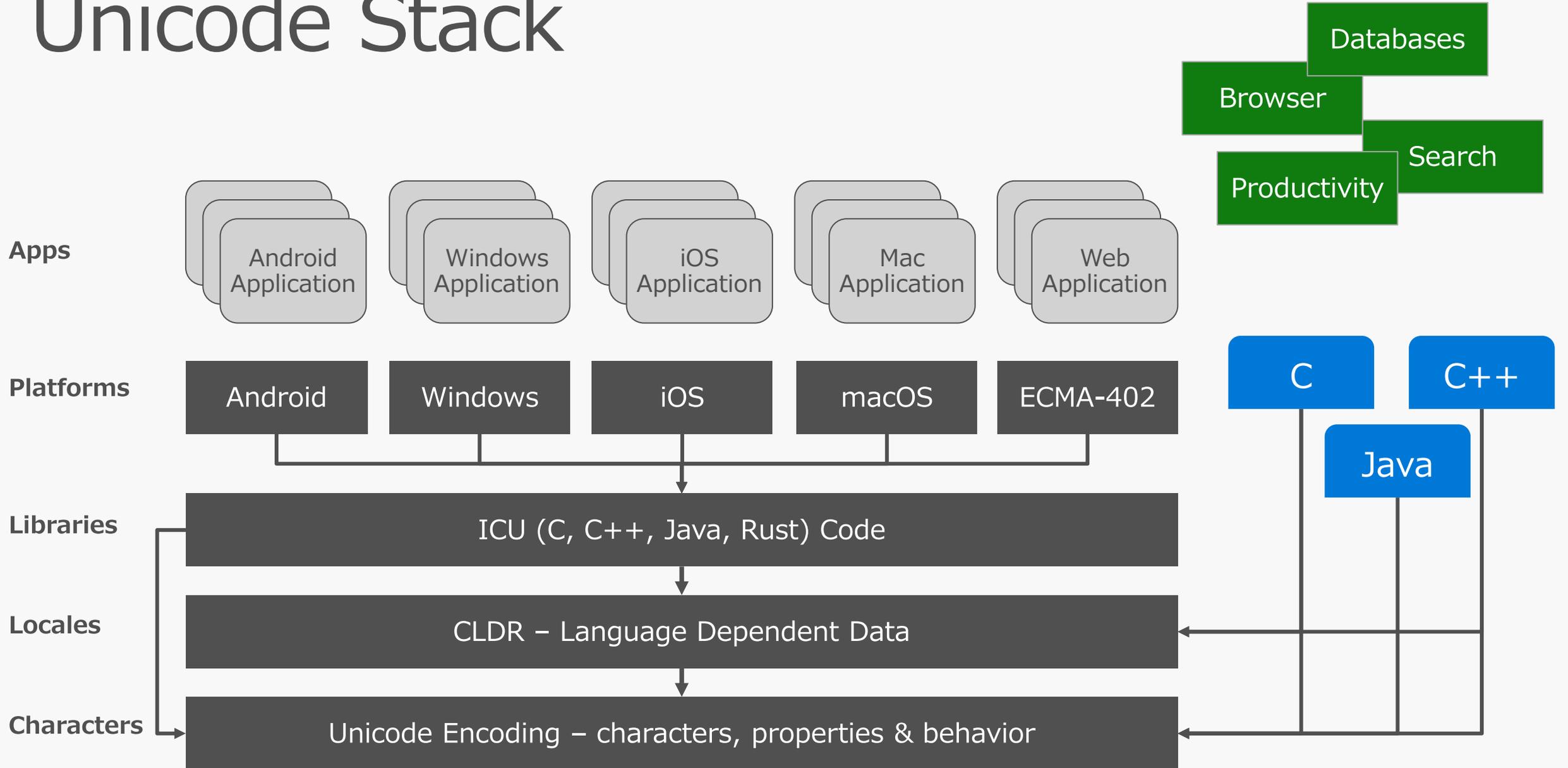
Unicodeのアドレス空間は**21ビット**

文字と規格の関係

- 国際規格（デジュール）と業界標準（フォーラム）が連携
- 文字に関する規格は、ICUに集約



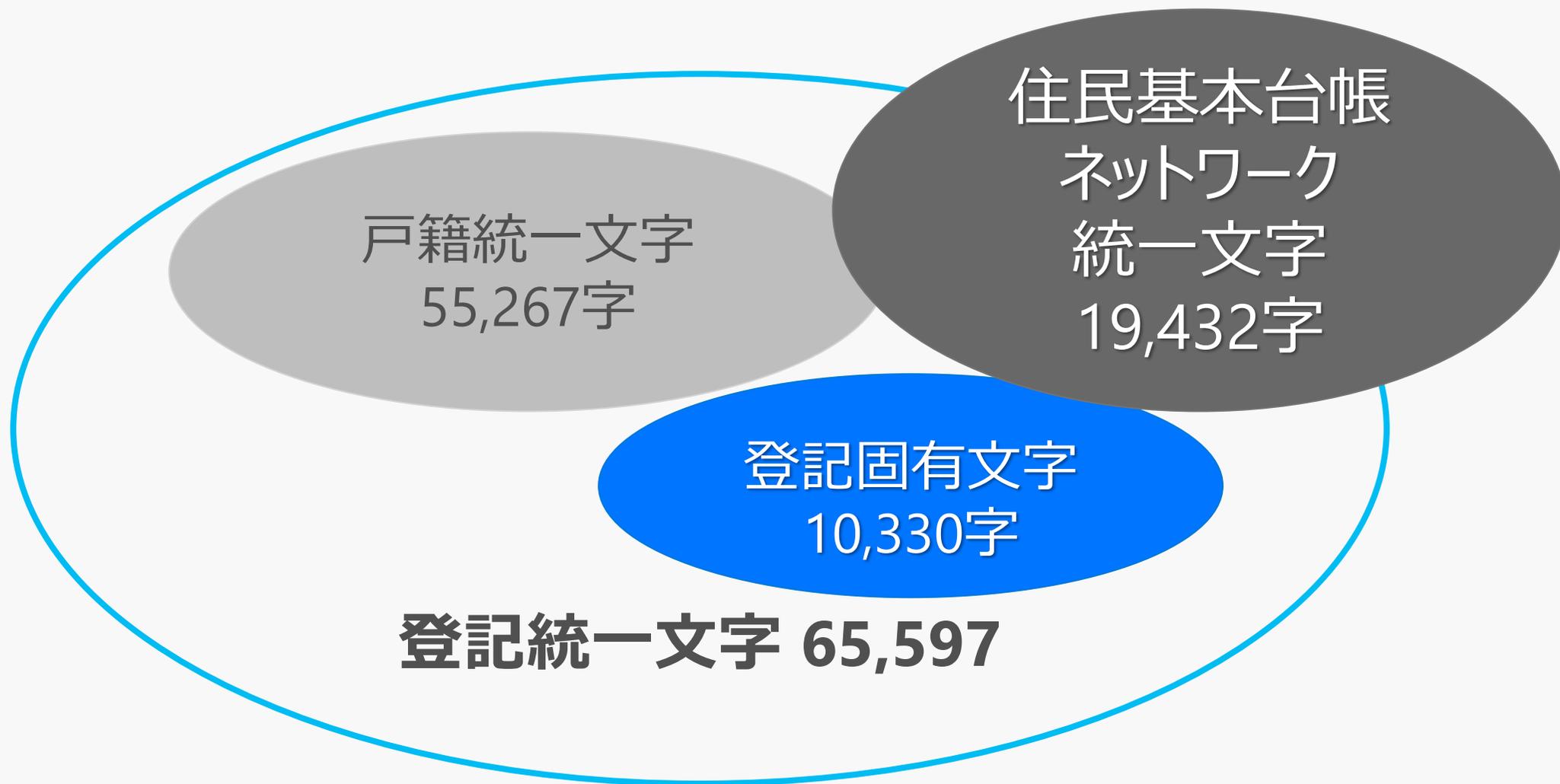
Unicode Stack





外字問題とIVD

行政システムでは



※それぞれ漢字以外の文字を除く

外字を使用する背景・・・

戸籍統一文字
55,267字

住民基本台帳
ネットワーク
統一文字
19,432字

登記固有文字
10,330字

登記統一文字 65,597字

無い文字は外字
で対応・・・



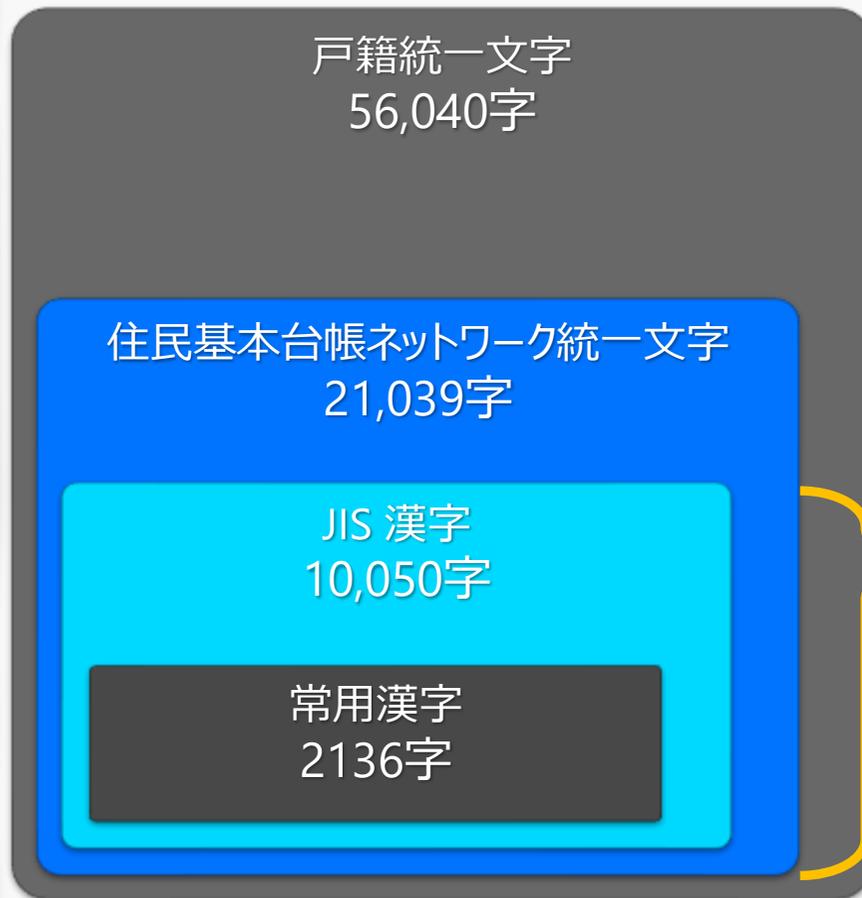
相互運用性

一般に・・・

- JIS、Unicodeを初め標準に含まれない字
- 改定常用漢字表
- JIS X 0208, JIS X 0213
- この他、フォントに含まれず独自に追加した字

外字を必要とする背景とは

9089	邊	邊	邊	邊
	E0100 Adobe-Japan1 CID+6930	E0101 Adobe-Japan1 CID+13407	E0102 Adobe-Japan1 CID+14241	E0103 Adobe-Japan1 CID+14242
	邊	邊	邊	邊
	E0104 Adobe-Japan1 CID+14243	E0105 Adobe-Japan1 CID+14244	E0106 Adobe-Japan1 CID+14245	E0107 Adobe-Japan1 CID+14246
908A	邊	邊	邊	邊
	E0108 Adobe-Japan1 CID+14247	E0109 Adobe-Japan1 CID+14248	E010A Adobe-Japan1 CID+14249	E010B Adobe-Japan1 CID+14250
	邊	邊	邊	
	E010C Adobe-Japan1 CID+14251	E010D Adobe-Japan1 CID+14252	E010E Adobe-Japan1 CID+20233	
908A	邊	邊	邊	邊
	E0100 Adobe-Japan1 CID+6929	E0101 Adobe-Japan1 CID+14235	E0102 Adobe-Japan1 CID+14236	E0103 Adobe-Japan1 CID+14237
	邊	邊	邊	邊
	E0104 Adobe-Japan1 CID+14238	E0105 Adobe-Japan1 CID+14239	E0106 Adobe-Japan1 CID+14240	E0107 Adobe-Japan1 CID+20234



外字を使用

UTF-8、UTF-16により符号化が可能

7ビット J I S による符号化が可能

BMPの外字使用と外字を使うことの深刻な課題



外字を使用した運用



私的コードの割当

1. コードの割り当て、管理
2. グリフの作成
3. フォントファイルへの追加
4. 更新したフォントの配布
5. 他システム連携のためのマッピングテーブルの更新



字体の作成



フォントファイルへの追加



フォントの配布

■ 文字のメンテナンス要する
膨大なコスト

■ 外部システムとの相互運用が
できない**閉じた文字情報空間**

東北大震災からの教訓

被災自治体



支援自治体



同じ字形であっても外字使用により
コードが異なり、データ連携ができない



9089	邊 E0100 Adobe-Japan1 CID-6930	邊 E0101 Adobe-Japan1 CID-13407	邊 E0102 Adobe-Japan1 CID-14241	邊 E0103 Adobe-Japan1 CID-14242
	邊 E0104 Adobe-Japan1 CID-14243	邊 E0105 Adobe-Japan1 CID-14244	邊 E0106 Adobe-Japan1 CID-14245	邊 E0107 Adobe-Japan1 CID-14246
	邊 E0108 Adobe-Japan1 CID-14247	邊 E0109 Adobe-Japan1 CID-14248	邊 E010A Adobe-Japan1 CID-14249	邊 E010B Adobe-Japan1 CID-14250
	邊 E010C Adobe-Japan1 CID-14251	邊 E010D Adobe-Japan1 CID-14252	邊 E010E Adobe-Japan1 CID-20233	
908A	邊 E0100 Adobe-Japan1 CID-6929	邊 E0101 Adobe-Japan1 CID-14235	邊 E0102 Adobe-Japan1 CID-14236	邊 E0103 Adobe-Japan1 CID-14237
	邊 E0104 Adobe-Japan1 CID-14238	邊 E0105 Adobe-Japan1 CID-14239	邊 E0106 Adobe-Japan1 CID-14240	邊 E0107 Adobe-Japan1 CID-20234

Unicode IVS/IVD

Windows 8 の IVSサポート基盤

- **文字情報基盤構築事業で整備された文字、約58,000文字を扱う基盤のサポート**
 - Unicode UTS #37 (IVS/IVD)
 - JIS X 0221:2013 追補版で対応予定
 - 印刷・出版業界で使用されているAJ1-6対応フォントを使用可能に
- **JIS X 0208:1990、JIS X 0213:2000、JIS X 0213:2004の文字全てを同時にサポート**
 - JIS90フォントパック使用することなく、JIS90の字体をIMEを使用して入力、表示可能に

グリフデータベース基盤



邊 E0100 Adobe-Japan1 CID=6930	邊 E0101 Adobe-Japan1 CID=13407	邊 E0102 Adobe-Japan1 CID=14241	邊 E0103 Adobe-Japan1 CID=14242
邊 E0104 Adobe-Japan1 CID=14243	邊 E0105 Adobe-Japan1 CID=14244	邊 E0106 Adobe-Japan1 CID=14245	邊 E0107 Adobe-Japan1 CID=14246
邊 E0108 Adobe-Japan1 CID=14247	邊 E0109 Adobe-Japan1 CID=14248	邊 E010A Adobe-Japan1 CID=14249	邊 E010B Adobe-Japan1 CID=14250
邊 E010C Adobe-Japan1 CID=14251	邊 E010D Adobe-Japan1 CID=14252	邊 E010E Adobe-Japan1 CID=20333	

様々なグリフをIdeographic Variation Databaseとして管理、利用可能にすることで、固有のIDによりグリフの識別を実現

フォントサイズの推移

組み込みなどOSの用途が広がる中で文字、
フォントを増やすことは容易ではない

	MSゴシック	MS明朝
Windows 3.1	2,500,692	3,446,204
Windows 2000 SP4	8,272,028	9,135,960
Windows Vista SP2	9,165,480	10,056,872
Windows 7	9,172,500	10,057,108
Windows 8	9,209,540	10,080,360
Windows 10	9,214,692	10,081,800

Windows 8 でIVSに正式対応

- IVS（異体字セレクタ）に対応
 - 文字の入力
 - 文字の出力（表示、印刷など）

Windows 8



文字情報基盤の成果

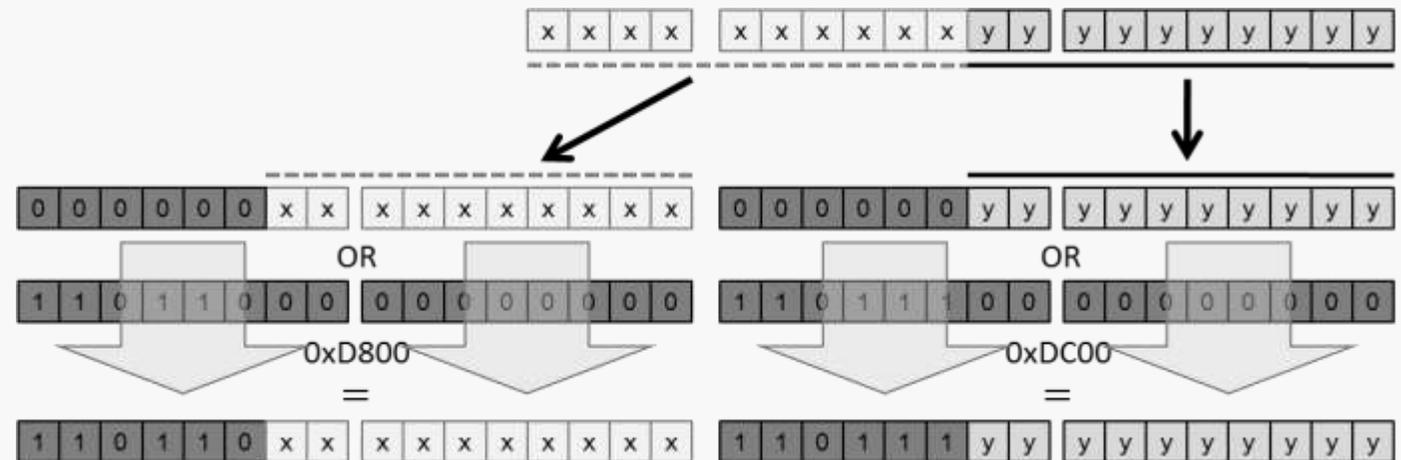
AJ1



Windows 8で異体字が使用可能に！

1文字 = 16ビットではない (UTF-16)

Unicodeのアドレス空間は**21ビット**
UTF-16は、**1文字を16ビット**で符号化
化する方式



レンダリング文字数とデータサイズは一致しない

「	丈	土	叱	」	は	追	加	面	の	文	字
300C	D840 DC0B	D844 DE3B	D842 DF9F	300D	306F	8FFD	52A0	9762	306E	6587	5B57
U+300C	U+2000B	U+2123D	U+20B9F	U+300D	U+306F	U+8FFD	U+52A0	U+9762	U+306E	U+6587	U+5B57

12文字 v.s. 15文字

C#におけるSurrogate Pairの識別

- Char.IsSurrogate(char)
- Char.IsSurrogate(string, int)

C#における見た目上の文字数の取得

- StringInfo.LengthInTextElements

異体字セレクタの仕組み

~Ideographic Variation Sequence~

- 基本となる文字と、その字形を指定する事が可能に
- 見た目の異なる文字であっても、基本となる文字により、意味上の一意性を確保



邊
9089

基本となる文字

U+9089のバリエーション

- IVD
 - Ideographic Variation Database
- IVS
 - Ideographic Variation Selector
 - Ideographic Variation Sequence

異体字セレクタ



U+9089

Unicode

9085	邊
9086	邊
9087	邊
9088	邊
9089	邊
908A	邊
908B	邊
908C	邊
908D	邊
908E	邊
908F	邊

グリフDB

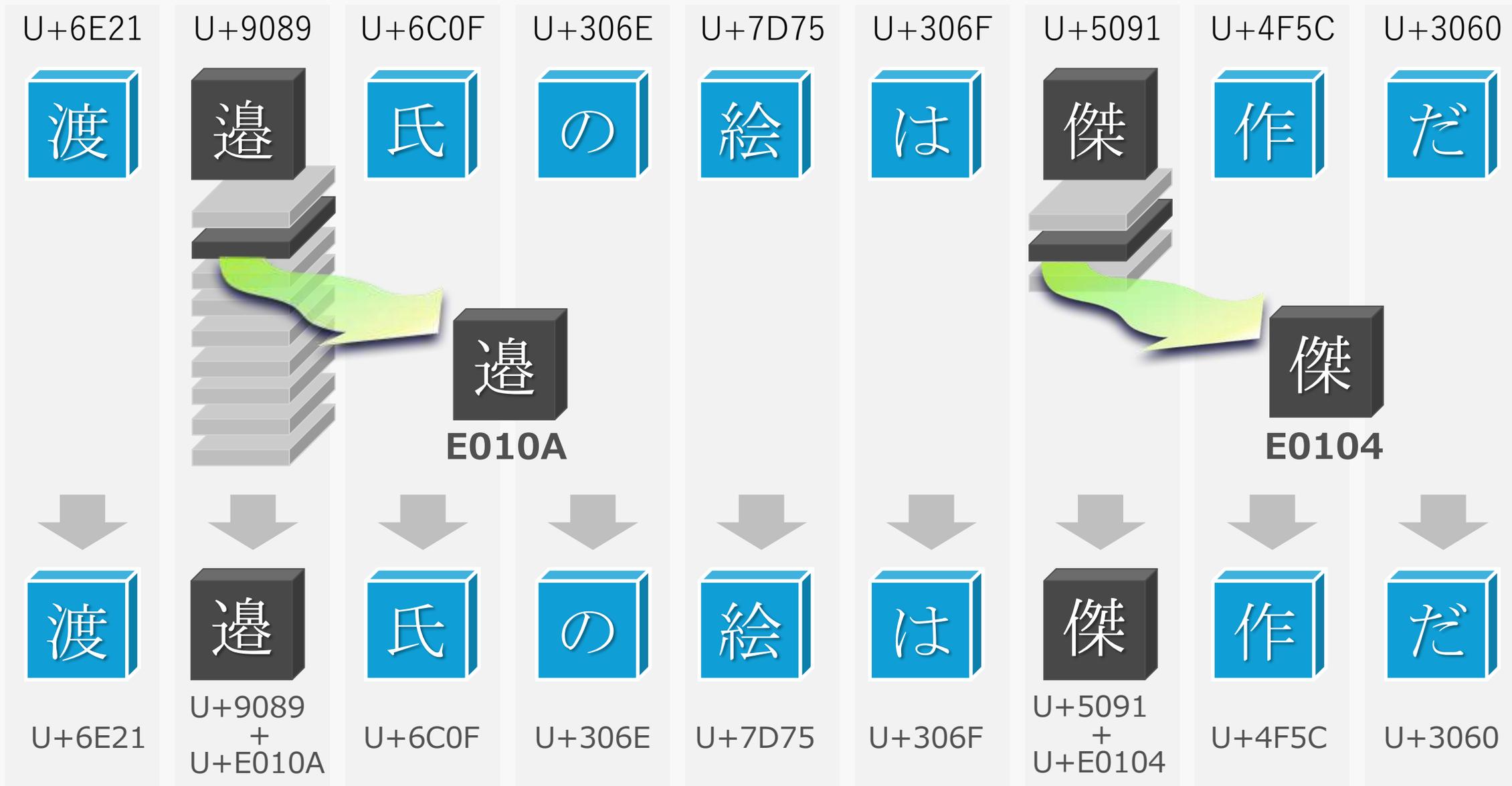
E010F	邊
E0110	邊
E0111	邊
E0112	邊
E0113	邊
E0113	邊
E0114	邊
E0115	邊
E0116	邊
E0117	邊
E0118	邊



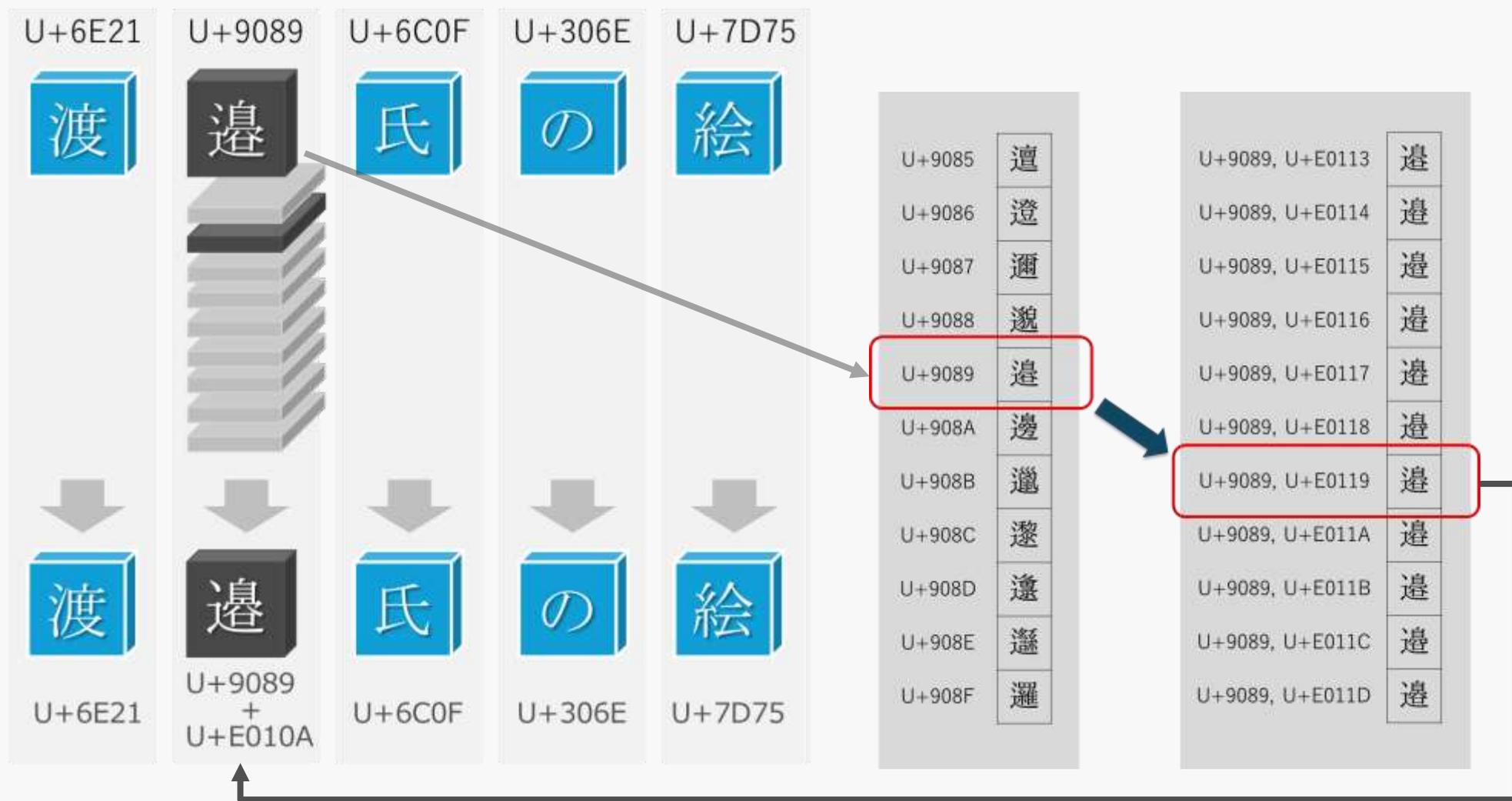
U+9089 U+E0116



グリフの指定が可能に



グリフの指定とIVD



16ビットから64ビット (UTF-16の場合)

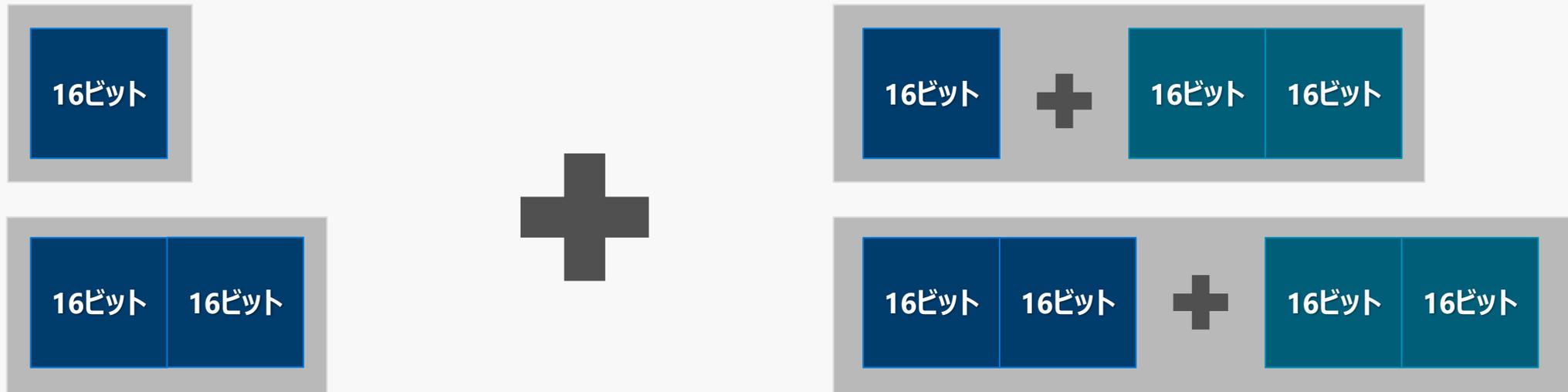
- BMP
 - 16ビット
- サロゲートペア
 - 32ビット (16ビット + 16ビット)
- JIS X 0213、IVSともにサロゲートエリアに割り当てられている



= 1文字の場合も

実装上の違い

	符号長/ 1 文字	不足する文字への対応
シフトJIS	可変 (8、1 6ビット)	外字
Unicode 基本多言語面	固定 (1 6ビット) *1	外字
Unicode サロゲートペア	可変 (1 6、3 2ビット) *1	外字
Unicode IVS/IVD	可変 (1 6、3 2、4 8、6 4ビット) *1	国際標準

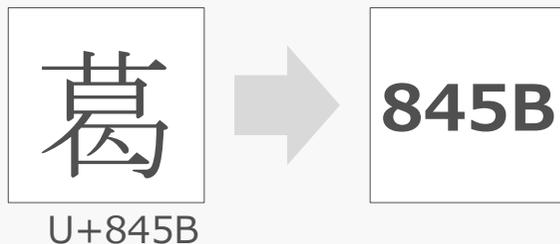


 = 1 文字

*1 - UTF-16符号化方式の場合

文字によって異なるバイト数 (UTF-16)

基本多言語面の文字



追加面の文字



基本多言語面の文字と
異体字セレクタ

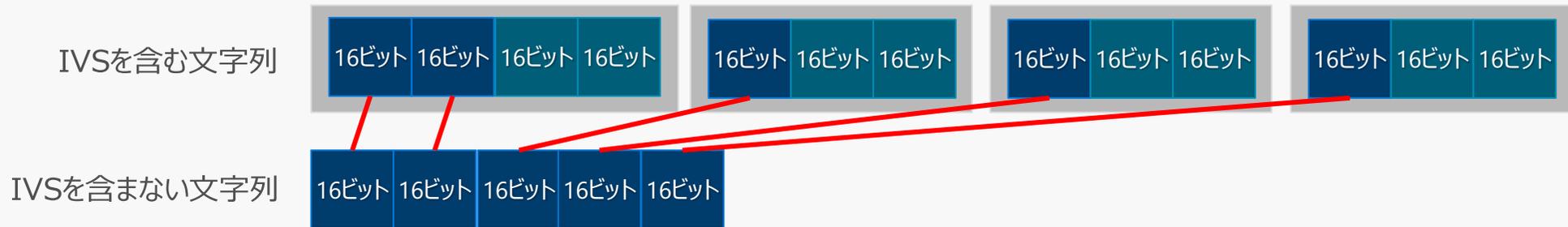


追加面の文字と
異体字セレクタ



各データベースにおける対応状況

- データベースにより異なる対応が異なる文字コード対応
- 照合順序がIVSに対応していないデータベースでは、親文字から構成される文字列による検索ができない。



	Upper/Lower	Accent	かな・カナ	Half/Full	IVS	文字情報基盤対応
PostgreSQL 14	citext extension	Unaccent extension				×
MySQL 8.0	CI/CS	AI		M		△
Oracle 21c	ci/cs	ai/as	ks	unicode		△
SQL Server	CI/CS	AI/AS	KS	WS	VSS	◎ *1

※1 - SQL Serverでは、「_VSS」が含まれる照合順序が異体字セレクタに対応



文字情報基盤

文字情報基盤整備事業

～内閣官房情報通信技術（IT）担当室、経済産業省、IPAにより推進～

平成14～17年度にかけて、戸籍統一文字、住民基本台帳ネットワーク統一文字を対象として、人名、地名等に使用される文字を調査を実施。平成20年度～登記に使用される文字を調査。

情報システムへの実装に向けて、ISO/IEC 10646に必要な文字の追加提案を行う。

これまでの経緯

年		発行主体等	内容
1978	(S53)	日本規格協会	JIS C 6226「情報交換用漢字符号系」
1994	(H6)	法務省	戸籍事務の電算化開始
1997	(H9)	閣議決定	「行政情報化推進基本計画（改定）」JIS3,4水準制定により解消。残る外字の交換ルール策定。
2000	(H12)	日本規格協会	JIS X0213「情報交換用符号化拡張漢字集合」（JIS第3,4水準）
2002	(H14)	総務省	住民基本台帳ネットワーク統一文字
2002	(H14)	経済産業省他	汎用電子情報交換環境整備プログラム開始（戸籍、住民基本台帳ネットワーク、登記を検証）
2003	(H15)	IT戦略本部決定	「e-Japan 戦略Ⅱ」公開用文字情報データベースの構築。文字コード規格を整備。
2004	(H16)	法務省	戸籍統一文字（戸籍手続オンラインシステムの構築のための標準仕様書）
2005	(H17)	IT戦略本部決定	「IT政策パッケージ-2005」
2008	(H20)	IT戦略本部決定	「IT政策ロードマップ」
2010	(H22)	経済産業省、IPA	文字情報基盤プロジェクト開始

これまでの経緯

年		発行主体等	内容
2010	(H22)	経済産業省、IPA	文字情報基盤プロジェクト開始
2011	(H23)	IT戦略本部決定	「電子行政推進に関する事本方針」文字情報基盤の活用
2011	(H23)	経済産業省、IPA	文字情報基盤プロジェクト成果公開 IPAmj 明朝フォント、文字情報対応表
2013	(H25)	閣議決定	「 世界最先端IT国家創造宣言 」文字情報基盤を原則
2014	(H26)	総務省	「 電子自治体の取組を加速するための10の指針 」
2014	(H26)	CIO連絡会議	「 電子行政分野におけるオープンな利用環境整備に向けたアクションプラン 」 導入ガイドの活用、縮退マップの整備

世界最先端IT国家創造宣言 内閣閣議決定

従来政府が担っていたサービスの提供機能を民間にも開放し、官民の協働によって、より利便性の高い公共サービスを創造する。国民がステークホルダーとして積極的に参加できるよう、クラウドを活用したオープンな利用環境を、データ・フォーマット、用語、コード、文字等の標準化・共通化、アプリケーション・インターフェイス（API）の公開等を行いつつ整備する。特に文字の標準化・共通化に関しては、今後整備する情報システムにおいては、国際標準に適合した**文字情報基盤**を活用することを原則とする。

電子行政分野におけるオープンな利用環境整備に向けたアクションプラン

各府省情報化統括責任者（CIO）連絡会議決定

今後整備する政府情報システムにおいては、国際標準に適合した**文字情報基盤**を活用することを原則とする。

各府省は、経済産業省を中心に策定した「文字情報基盤導入ガイド」を参照しつつ、導入を推進する。また、経済産業省は、各府省の円滑な導入を支援するため、文字情報基盤の文字（約6万文字）を、市販コンピュータで特別な設定無しで活用できるJIS範囲の文字への変換を行う際に参照する縮退マップの提供を、平成26年度中に実施する。

IPAmj明朝

初版：2011年10月

新フォント (IPAmj明朝) を開発

異体字のさまざまなパリエーションの例：「辺」、「齊」

邊邊邊邊邊邊邊邊邊邊
邊邊邊邊邊邊邊邊邊邊
邊邊邊邊邊邊邊邊邊邊

齊齊齋齋齊齊齊齋齊齊
齋齋齋齋齋齋齋齋齋齋
齋齋齋齋齋齋齋齋齋齋

文字情報基盤概要

- 我が国で行政業務上必要とされる人名漢字等約6万文字を、国際標準に則り提供
 - IPAmj明朝フォント
 - 文字情報一覧表
 - <http://mojikiban.ipa.go.jp/>

文字情報基盤 (58,814字)

戸籍統一文字 (漢字55270字)

戸籍のオンライン手続に使用することを目的として整理した文字 (辞書をベースに整理)

住民基本台帳ネットワーク システム統一文字 (漢字19563字)

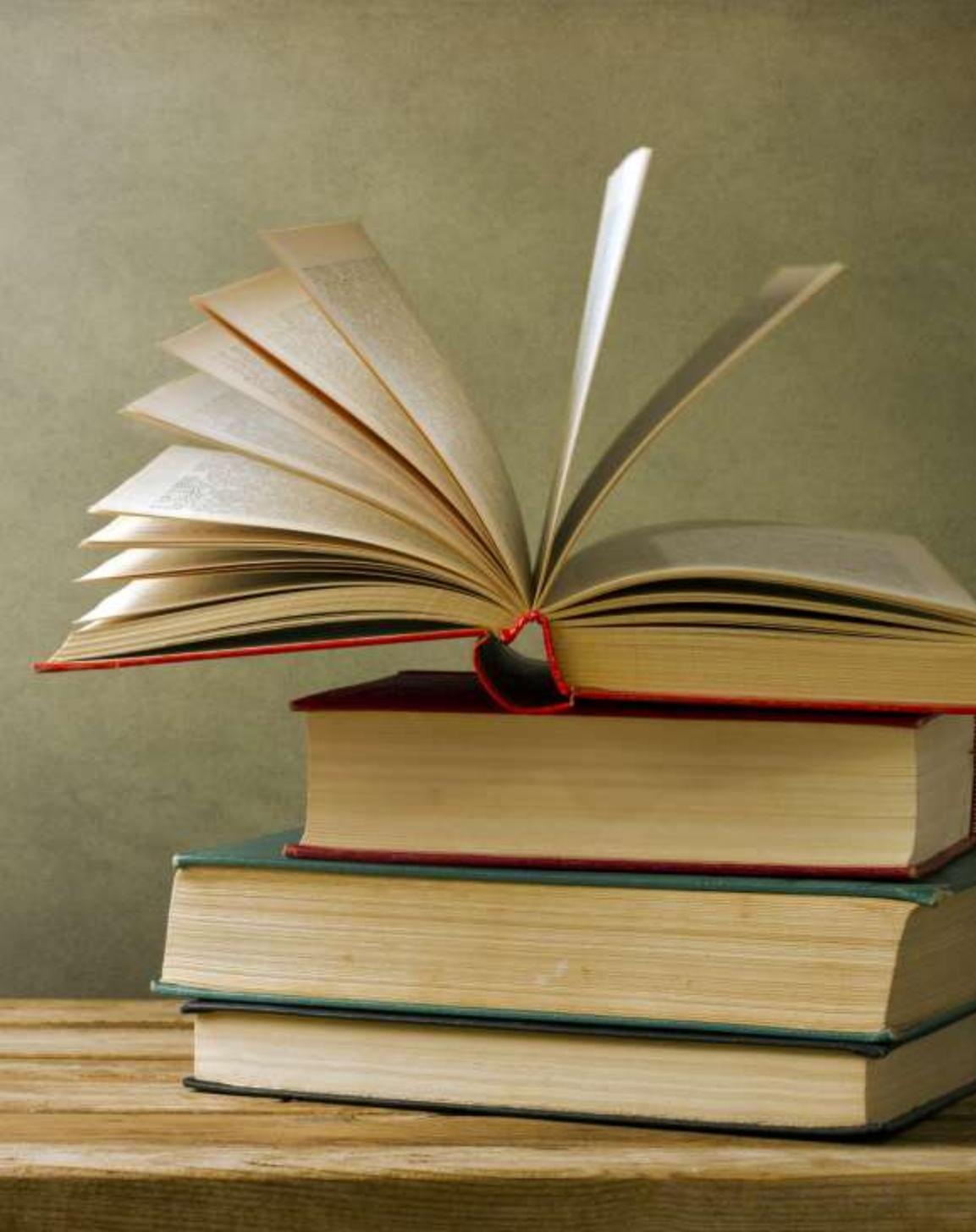
多くの住民が氏名に使う文字を整理

JIS漢字コード (10050字)

実用上の情報交換の必要性から、出現頻度などを元に文字を選定 (JIS X 0213:2004)

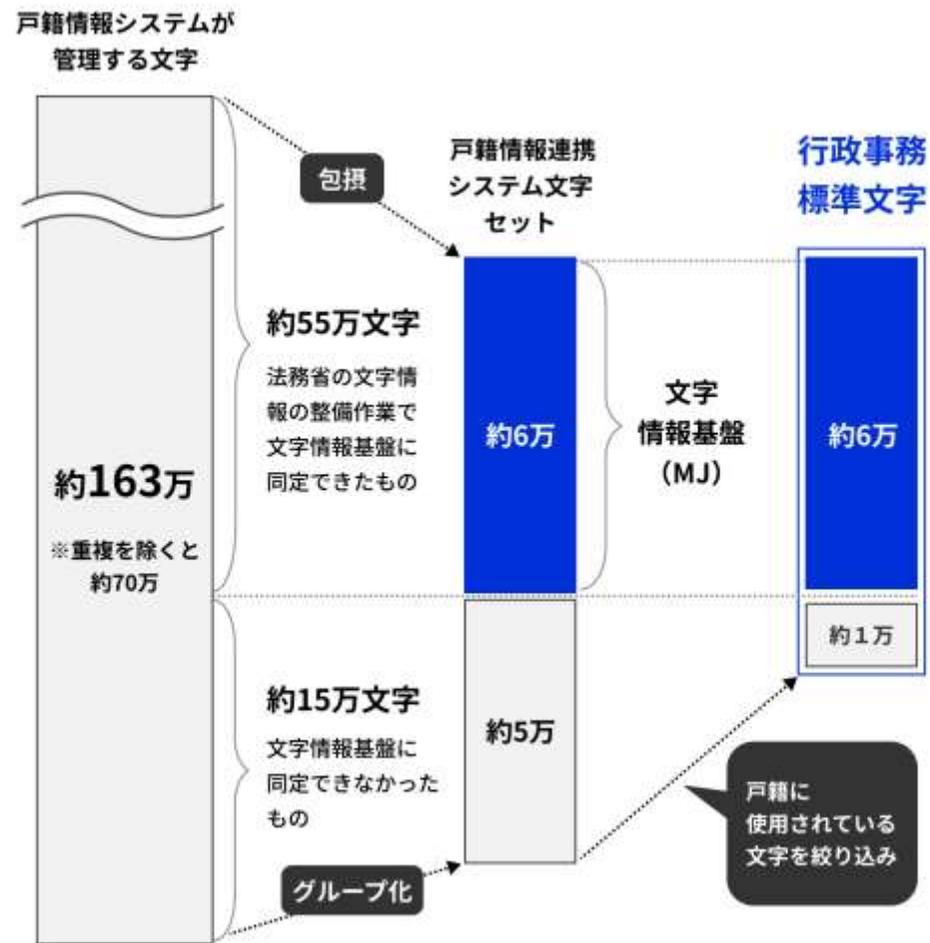
常用漢字 (2136字)

法令、公用文書、新聞、雑誌、放送など、一般の社会生活において、現代の国語を書き表す場合の漢字使用の目安を示す。



行政事務標準文字

行政文字標準文字とは



外字による様々な課題の解消

- **外字作成コスト**

- 地方公共団体ごとや、一つの地方公共団体の中でも手続を行う部署ごとに文字データを作成・維持しています。

- **地方公共団体の職員・住民の負担大**

- 転入時に住民票の写しを即日発行できません。地方公共団体の職員が外字を作成したのち、住民は後日住民票の写しを受領します。
- 地方公共団体等の公的機関では、手続きの際、氏名の入力誤りが発生しないようするための確認作業等で大きな負担が発生します。
- 地方公共団体ごとにコンピューターで管理する文字が異なることで、効率的な行政サービスの実施および大規模な災害発生時の迅速な対応等の妨げとなります。

- **システム選択時の制約**

- システム更新の際には、外字が支障となって、同じベンダーに依存し続けるベンダーロックインの状態になりやすく、文字の作成や維持・管理にコストがかかる要因にもなります。

- **システム間での情報連携を阻害**

- 住民サービスに用いられているコンピューターのデータが別々のシステムで管理されているため、外字を使用している場合は、異なる地方公共団体間のデータ連携だけでなく、同じ地方公共団体の中の組織であっても、部署が異なればシステム間のデータ連携の際に文字化け等が発生し、情報の連携を阻害する要因になっています。

一般社団法人 文字情報技術促進協議会

<https://moji.or.jp>

全ての文字を全ての環境へ

- フォント技術
- 標準化
- 技術支援
- 政策

※Unicode Consortium alliance organization

<https://moji.or.jp>

文字情報技術促進協議会

ホーム 各活動について 文字情報基盤 Unicode IVS対応製品 セミナー・関連資料 協議会について

UNICODE IVS

Unicode IVSは、外字を使わずにいた字を表現できる国際標準規格です。(Unicode、ISO-10646規格の一部)

[詳しくはこちら](#)

文字情報技術促進協議会が「文字情報基盤」の成果物を、情報処理推進機構から移管

～ 相互運用性のさらなる拡大へ ～

[詳しくはこちら](#)

- Unicode IVS**
Unicode IVSは、外字を使用せずに異体字を表示、印刷、データとしての送受信を可能にする仕組みです。国際標準であり、高い相互運用性が保証されています。
- IVS 対応製**
既に多くのフォント、アプリケーション、オペレーティングシステムが、Unicode IVSに対応しています。今からでも、人名、地名に正しい漢字を使用できるシステムへ。
- 文字情報基盤**
これまで外字でしか表現できなかった人名、地名を国際標準規格に沿って使用できるようにする取り組みです。Unicode IVSに対応したシステムは、文字情報基盤に準拠したシステム構築が可能です。
- 人名・地名**
Unicode IVSに対応した製品を使用することで、正しい人名・地名を表示、印刷することができます。マイナンバーを見据えた情報システム構築をサポートします。

令和4年度体制 (2022/05)

- 令和3年度 理事会・社員総会が開催されました。(2022/05/25)
- 文字情報技術促進協議会『年次特別講演会』の申し込み受付を開始しました。(2022/05)
- 小林会長ブログ「UCS水平拡張とブリックレビュー」を掲載しました。(2022/02/08)
- 小林会長ブログ「M文字固有形名とUCS符号位置との対応」を掲載しました。(2021/06/03)
- 文字情報基盤検索システムを公開しました。(2021/03/29)
- 文字情報技術促進協議会『年次特別講演会』の講演スライドを掲載しました。(2021/03)
- 株式会社モリサフ 高井幸代様「フォント擬人化CODE」

文字が描く、デジタル社会の未来図

～ 一般社団法人文字情報技術促進協議会 15周年シンポジウム ～

文字・フォントの最新動向から政府の外字問題解消の取り組みまで、文字・フォントに関する今をご紹介します

日時： 12月12日（金）
13:00 ～ 17:00

場所： J P タワーホール&カンファレンス
東京都千代田区丸の内二丁目7番2号
KITTE 4, 5階

参加費： 無料



内容	講演者
開会あいさつ	一般社団法人文字情報技術促進協議会 会長 小林 龍生
基調講演 (1)	デジある庁 デジタル社会共通機能グループ 統括官 楠 正憲
基調講演 (2)	Unicode Consortium Technical Director Ken Lunde
セッション1	大学共同利用機関法人 人間文化研究機構 国立国語研究所 教授 高田 智和
セッション2	一般社団法人文字情報技術促進協議会 理事 / 日本電気株式会社 袴田 博之
セッション3	Monotype 開発部 津田 昭 モリサワ デザイン部門長 岡 繁樹 株式会社イワタ 代表取締役社長 水野昭
パネルディスカッション	
閉会あいさつ	一般社団法人文字情報技術促進協議会 事務局長 田丸 健三郎



文字情報技術促進協議会

