# 失われた文字コード

安岡 孝一(やすおか こういち)

## ■ はじめに

「失われた文字コード」とタイトルに付けたものの、実は、文字コードそのものが失われているわけではない。失われているのは、ある文字コードを読むための環境であって、その文字コードで書かれたファイルそれ自体は、古い CD-R や磁気ディスクやその他の記憶媒体の中に、ひっそりとしかし厳然と存在していたりする。そのような「失われた文字コード」で書かれたファイルが、何かのはずみで我々の前に姿を現す。実際、読めない文字コードで書かれたファイルほど、始末の悪いものは無い。しかも、完全に読めないならまだし

図 1 GB 2312 の 16 区 (シフト GB の 889F ~ 88FC)

% 败拜稗斑班搬扳般颁板版扮拌伴瓣半办绊邦帮

<sup>80</sup> 梆榜膀绑棒磅蚌镑傍谤苟胞包褒剥

図2 JIS X 0208の16区(シフト JISの889F~88FC)

 $16\, \overline{\boxtimes}$  01 02 03 04 05 06 07 08 09 10 11 12 13 14 15 16 17 18 19

9 里壁灶門农友沃加建关四闸芯连连旭革户黔

20 梓圧斡扱宛姐虻飴絢綾鮎或粟袷安庵按暗案闇

<sup>40</sup> 鞍杏以伊位依偉囲夷委威尉惟意慰易椅為畏異

% 移維緯胃萎衣謂違遺医井亥域育郁磯一壱溢逸

🖁 稲茨芋鰯允印咽員因姻引飲淫胤蔭

も、平仮名や片仮名など一部が読めるのに他が文字化けしていたりすると、何としても残りも読みたくなってしまうのが人のサガである。

本稿では、そのような「失われた文字コード」 のうち、本誌の読者の手元に残っていそうな文字 コードとして、特に「シフト GB」と「文字鏡コード」 を取り上げる。また、これらに加えて「住基コード」に対しても、警鐘を鳴らしておくことにする。

## ■ シフト GB

GB 2312<sup>[1]</sup> に収録された簡体字を、日本語版 MS-DOS や日本語版 Windows 3.x で使うために 開発された文字コードである。本来 GB 2312 は、  $94 \times 94$  の文字表を A1A1  $\sim$  FEFE に割り当てて

用いるが、日本語版 MS-DOS では FD や FE のコードは使えなかった。そこで、 JIS X 0208 をシフト JIS に変換するやり 方をまねて、GB 2312 の 94 × 94 の文字表を  $47 \times 188$  に折り曲げ、1 バイト目に  $81 \sim 9F \cdot E0 \sim EF$  を、2 バイト目に  $40 \sim 7E \cdot 80 \sim FC$  を使うことで、 いわばシフト JIS の上に GB 2312 を載っけてしまうのが、シフト GB の基本的アイデアである。

当然、同じ文字コード上にシフト GB とシフト JIS が重なることになり、同じ「889F」という文字コード(図 1・2 の16-01 にあたる)であっても、シフト GB では「啊」を、シフト JIS では「亜」を、それぞれ意味することになる。たとえば

「白雪皑皑」という文字列は、シフト GB では「88D5 9948 88A6 88A6」という文字コードになるが、これをシフト JIS だと思って表示すると「易僣始始」となってしまうわけだ。逆に言えば、「易僣始始」とシフト JIS で書かれたファイルを、シフト GB のフォントで表示すれば、見事「白雪皑皑」になるという仕掛けだった。これにより、フォント切り換えの可能なワープロソフト上では、日本語と中国語を混在して表示できたのである。

#### ★ Windows 98 以降のシフト GB

ところが Microsoft は、Windows 98 において全てのフォントを Unicode 化することを決定、同時に MS-Word の内部コードを Unicode 化したため、それまでのシフト GB は使えなくなってしまった。シフト JIS の文字コードのうち、CP932 において Unicode との対応を規定されている部分しか、使えなくなってしまったのである。たとえば、CP932 は 9873 ~ 989E(図 3 の 47-52 ~ 47-94)に、文字を収録していない。一方シフト GB は、この部分に「显」~「晓」の 43 字(図 4 の 47-52 ~ 47-94)を収録していた。この結果、Windows 98

 $47 \, \boxtimes \, \, \, \text{ or } \, \, \text{$ 

- 量 蓮連錬呂魯櫓炉賂路露労婁廊弄朗楼榔浪漏
- 20 牢狼篭老聾蝋郎六麓禄肋録論倭和話歪賄脇惑
- 40 枠鷲亙亘鰐詫藁蕨椀湾碗腕

図3 JIS X 0208 の 47 区(シフト JIS の 9840  $\sim$  987E  $\cdot$  9880  $\sim$  989E)

 $47 \boxtimes \ \ \text{01} \ \ \text{02} \ \ \text{03} \ \ \text{04} \ \ \text{05} \ \ \text{06} \ \ \text{07} \ \ \text{08} \ \ \text{09} \ \ \text{10} \ \ \text{11} \ \ \text{12} \ \ \text{13} \ \ \text{14} \ \ \text{15} \ \ \text{16} \ \ \text{17} \ \ \text{18} \ \ \text{19}$ 

- <sup>11</sup> 稀息希悉膝夕惜熄烯溪汐犀檄袭席习媳喜铣
- 20 洗系隙戏细瞎虾匣霞辖暇峡侠狭下厦夏吓掀锨
- % 先仙鲜纤咸贤衔舷闲涎弦嫌显险现献县腺馅羡
- % 宪陷限线相厢镶香箱襄湘乡翔祥详想响享项巷
- 80 橡像向象萧硝霄削哮嚣销消宵淆晓
- 図4 GB 2312 の 47 区(シフト GB の 9840 ~ 987E・9880 ~ 989E)

 $13 \, \boxtimes \, \, \text{ or } \, \text{$ 

01

显险现献县腺馅羡宪

20 陷

限线相厢镶香箱襄

- 40 湘乡翔祥详想响享
- 项巷橡像向象萧硝霄削哮 器销消宵淆
- 図5 Unicode 対応シフトGBの13区(8740~877E・8780~879E)

表 1 シフト GB の Unicode 対応

オリジナルのシフト GB		Unicode 対応シフト GB		
GB 2312	シフト GB	移動先	NEC 拡張文字版シフト GB	IBM 拡張文字版シフト GB
$47-52 \sim 47-61$	9873 ~ 987C	13-11 ∼ 13-20	874A ~ 8753	874A ∼ 8753
$47-62 \sim 47-77$	987D ∼ 987E ⋅ 9880 ∼ 988D	$13-32 \sim 13-47$	875F ∼ 876E	875F ∼ 876E
47-78 ~ 47-80	988E ~ 9890	$13-63 \sim 13-65$	877E⋅8780 ~ 8781	877E⋅8780 ~ 8781
47-81	9891	13-66	8782	FA59
47-82	9892	13-67	8783	8783
47-83	9893	13-68	8784	FA5A
$47-84 \sim 47-88$	9894 ~ 9898	$13-69 \sim 13-73$	8785 ~ 8789	8785 ~ 8789
47-89 ~ 47-93	9899 ~ 989D	$13-75 \sim 13-79$	878B ~ 878F	878B ∼ 878F
47-94	989E	13-83	8793	8793
84-01 ~ 84-94	easf $\sim$ easc	89-01 ~ 89-94	ED40 $\sim$ ED7E $\cdot$ ED80 $\sim$ ED9E	FASC $\sim$ FA7E $\cdot$ FA80 $\sim$ FABA
$85-01 \sim 85-94$	EB40 $\sim$ EB7E $\cdot$ EB80 $\sim$ EB9E	$90-01 \sim 90-94$	ED9F $\sim$ EDFC	fabb $\sim$ fafc $\cdot$ fb40 $\sim$ fb5b
86-01 ~ 86-94	EB9F $\sim$ EBFC	91-01 ~ 91-94	EE40 $\sim$ EE7E $\cdot$ EE80 $\sim$ EE9E	FB5C $\sim$ FB7E $\cdot$ FB80 $\sim$ FBBA
$87-01 \sim 87-78$	EC40 $\sim$ EC7E $\cdot$ EC80 $\sim$ EC8E	$92-01 \sim 92-78$	EE9F $\sim$ EEEC	FBBB $\sim$ FBFC $\cdot$ FC40 $\sim$ FC4B
87-79	EC8F	13-84	8794	8794
87-80	EC90	13-88	8798	8798
87-81 ~ 87-90	EC91 ∼ EC9A	92-81 ~ 92-90	EEEF $\sim$ EEF8	fa40 $\sim$ fa49
87-91	EC9B	13-89	8799	8799
87-92 ~ 87-94	EC9C $\sim$ EC9E	$92-92 \sim 92-94$	eefa $\sim$ eefc	FA55 ∼ FA57

や Windows 2000 では、シフト GB の「显」~「晓」 を表示できなくなってしまった。

この問題を回避するため、シフト GB を使用していた各社は、これらの文字を 13 区に移動させることにした(図 5)。CP932 の 13 区にあたる部分(8740~877E・8780~879E)は、NEC特殊文字と呼ばれる独自拡張文字を収録していたが、この部分をこれまでのシフト GB は使っていなかったからだ。また、GB 2312 の 84 区以降に関しても、NEC 拡張文字のある 89 区以降に移動させた。移動の詳細を表 1 に示す。なお、NEC 拡張文字(ED40~EEFC)に対しては、全く同じ文字が多少順序を変えて IBM 拡張文字(FA40~FC4B)にも収録されており、ファイル入出力の際に相互に自動変換される。このような自動変換がおこなわれた場合の文字コードを、表 1 に「IBM 拡張文字版シフト GB」として示しておく。

#### ♣ シフト GB のサルベージ

前世紀末には一世を風靡したシフト GB も、今世紀に入ってからは衰退の一途をたどっている。開発の一翼を担ったオムロンも、2010 年5月で『楽々中国語』(cWnn)のサポートを完全に打ち切ってしまった。現時点でシフト GBをサポートしている製品は、わずかに高電社の『ChineseWriter』のみ<sup>[3]</sup>となってしまっている。では、何かの拍子にシフト GB のファイルを発掘してしまった場合、『ChineseWriter』を購入して読むしかないのだろうか。

実は、シフト GB のサルベージに関しては、針谷壮一氏が公開している「中国語コンバータ [4]」というフリーウェアが、現時点において最も強力なツールである。各社固有のシフト GB テキストを GB 2312 の A1A1 ~ FEFE に変換してくれると同時に、シフト GB・シフト JIS 混在の rtf 形式も変換可能というスグレモノだ。すなわち、doc 形式のシフト GB が発掘された場合でも、MS-Word の「リッチテキスト形式」でそのまま保存しなおすことができれば、「中国語コンバータ」によって、フォント名も含め変換可能 [5] となっている。

ただし、サルベージ変換後のテキストやrtfは、現在は読み書きが出来ても、いずれ遠い将来には、また「失われた文字コード」となりかねない。これらのテキストやrtfは、プリントアウトして紙の形にしておくか、せめてフォント埋め込みのPDFにしておく必要があるだろう。

## ■ 文字鏡コード

『今昔文字鏡』で使われている10進6桁の文 字鏡番号に対し、000001~005640、005641 ~ 011280, 011281 ~ 016920, ·····, & 5640 個ずつに区切り、各 5640 個を JIS X 0208 の 16~45・48~77 区に割り当てた文字コー ドである。文字コード上は完全にダブった形とな るが、000001 ~ 005640 は「Mojikyo M101」 フォントに、005641 ~ 011280 は「Mojikyo M102」フォントに、011281 ~ 016920 は 「Mojikyo M103」フォントに、……、それぞれ 収録することにより、フォント切り換えによって 全ての文字を扱うことができる、というのがウリ であった。逆に言えば、文字鏡フォントが実装さ れていない環境下では、文字鏡コードで書かれた 「一」「声」「憺」……は、全て「亜」に文字化け することになる (図6~8および図2)。

文字鏡番号の000001 ~ 049964 は基本的に『大漢和辞典』の検字番号を踏襲しているが、それ以降の番号は文字鏡研究会が独自に付けたものである。管理を担当する文字鏡フォントセンターは、2000 年 9 月の使用許諾書において、文字鏡番号と他の文字集合との対応表および、その対応表を用いて文字の変換をおこなうソフトウェア等の作成・配布を、明確に禁止している。

## ☆ 文字鏡コードの現況

文字鏡コードで書かれた MS-Word や Excel あるいは一太郎などのファイルは、文字鏡フォントが実装されていない環境では文字化けが起こるが、文字鏡フォントセンターは 2008 年 3 月以降、文字鏡フォントの無償ダウンロードを打ち切っている。文字鏡番号の付与をおこなってい

たはずの文字鏡研究会は、ほぼ同時期に活動を凍結し、そのWWWサイト<sup>[6]</sup>も「再構築中」の状態が続いている。また、現時点で文字鏡コードをサポートしている製品は、エーアイ・ネットの『今昔文字鏡』のみである。ただし、各文字鏡番号に対応する字形に関しては、ISO/IEC 10036 のグリフ登録サイト<sup>[7]</sup>で閲覧可能だ。ISO/IEC 10036 のグリフ番号 10000001~10576000が文字鏡番号 000001~576000に対応しており、128 ドットのgif 画像が入手できる。

このような状況において、文字鏡コードで書かれたファイルが発掘された場合、どうやってサルベージすればいいのか。残念ながら文字鏡フォントセンターが、そのようなサルベージをおこなうソフトウェアを明確に禁止してきたため、現時点では全く妙案がない。化けてしまった一字一字に対して、その字のJIS X 0208 における区点番号を調べて、フォント名と共に文字鏡番号に変換 [8] し、さらに 100000000 を加えて、ISO/IEC 10036 のグリフ登録サイトから gifを入手 [9] するしかない。もし、フォント名がわからない場合は、とりあえず

「Mojikyo M101」の文字鏡番号を計算し、その グリフ番号に順に 5640 を加えていって当該グリ フを入手、というのを繰り返すことになる。こん なのを手作業でおこなうのは、全く馬鹿げている としか言いようがないのだが、ソフトウェアを作 成できない以上、どうしようもない。

## ■ 住基コード

住民基本台帳ネットワークにおける統一文字 コードとして、2001年2月に検討版が作成され、 2001年10月に制定、2002年8月より運用開始された文字コードである。Unicode の基本多言 語面(U+0000~U+FFFF)を独自拡張した2バイトの文字コードであり、0000~9FFF および

袋 丞丟弫丽至北両正丣宪丽两听並並太恆考茵丽

₩ 排死些毘売畾腊 | 以 个丫中中乳 ≠ 丰 中 丫 中

図 6 「Mojikyo M101」の 16 区 (文字鏡番号 000001 ~ 000094)

80 夏曼酸氫夏奧耜夏矮毀慶夏麥復變

図7 「Mojikyo M102」の16区(文字鏡番号005641~005734)

 $16 \, \boxtimes \ \ \text{O1} \ \ \text{O2} \ \ \text{O3} \ \ \text{O4} \ \ \text{O5} \ \ \text{O6} \ \ \text{O7} \ \ \text{O8} \ \ \text{O9} \ \ \text{10} \ \ \text{11} \ \ \text{12} \ \ \text{13} \ \ \text{14} \ \ \text{15} \ \ \text{16} \ \ \text{17} \ \ \text{18} \ \ \text{19}$ 

20 慢鮲噶燭憹憺憻懏홨憽懆噤憾懅懅熬憿懀嶫懁

<sup>40</sup> 懂歎懄馀懅懆貇豤懈艀噟敤懊懋懌懍憤懎儶懗

59 里念[新]逐[涿]朱芯龙[肝芯]忘念[失芯]羊[禾]貝[回]阿|艸

図8 「Mojikyo M103」の16区(文字鏡番号011281~011374)

E000 ~ FFFF に関しては Unicode と互換である (図 9)。ただし、ハングルにあたる領域(AC00 ~ D7A3)に関しては、ハングルではなく独自の 拡張漢字や変体仮名などを、勝手に割り当てる 形で設計された(図 10)。これに加え制定版では、JIS X 0213 の第  $3\cdot 4$  水準漢字のうち漏れていた 漢字を、AAA1 ~ ABCF に追加している(図 11)。

住基コードの仕様は一般には非公開となっており、現時点で唯一のフォントである「KAJO J 明朝」を購入して、文字コードを調べるしかない。ただし、2009年3月発行の『汎用電子情報交換環境整備プログラム成果報告書別冊』(日本規格協会)には、住基コードの漢字部分が収録されており、漢字に関しては一応の調べがつくようになった。

0 1 2 3 4 5 6 7 8 9 A B C D E F

4E0x 一丁万七上丁 万丈三上下丌不与丏

4E1x 丐丑丒 且丕世丗丘丙业 丞丢

4E2x 両丢丣两 並 | 山个丫丬中丮丰

図 9 住基コードの 4E00 ~ 4E3F

0 1 2 3 4 5 6 7 8 9 A B C D E F

ADIx 式甲盯听亟亟茄亜砝爾丰 暢々主棄

AD2x 乂乂メ人凢乕乵乳亂豫夾亥交亮亭亮

AD3x 牽曺齊襄齊廖輔齋齋齋廖團齎亹變个

図 10 住基コードの AD00 ~ AD3F

0 1 2 3 4 5 6 7 8 9 A B C D E F

AAAX 一丈斥自へ個嗶集嘿價儼吳浴冂剱

AABX 劉边勵斗卓云员硴喜嗎噌唔噔里土块

AACx 垆华烬绘堅捇垭据望壤增寿莫螤姆妣

図 11 住基コードの AAA1 ~ AADF

#### ☆ 住基コードの崩壊

2008 年 12 月 に Tai Viet<sup>[10]</sup> が U+AA80 ~ U+AADF に収録されたことから、住基コードは急速に崩壊が進んでいる。というのも、たとえばU+AABE は「TAI VIET VOWEL AM」という非前進文字なので、通常の Unicode 処理では直前の文字に重ねて表示される。この結果、住基コードの AABE にあたる「土」は、直前の文字と重なって表示されてしまう [11] のである。

今後、Unicode の U+AAA1  $\sim$  U+ABCF に 新たな文字が追加されるたび、住基コードの崩壊が加速していくのは間違いない。本来、これらの独自拡張漢字は、U+20000  $\sim$  U+2FFFF の拡張漢字面にちゃんと収録するか、さもなくば IVS[12] を使って見分けるべきなのだ。元々の設計が完全に誤っ

ている上に、住基コードの仕様を公開しないという方針を堅持し続けているため、 もはや誰からも救いの手を差しのべようがない、というのが現状である。

#### ■ おわりに

シフトGB、文字鏡コード、住基コード、の3つの文字コードについて概要を述べ、これらの文字コードがどのようにして「失われていった」のかを解説した。これらの文字コードに共通する特徴は、特定のベンダやグループによってややもすれば閉鎖的に使われており、標準化からは程遠い立場にあったことである。言い換えれば、標準化をおこたった文字コードは、いずれ「失われる」ということだろう。

もし、読者諸氏が今現在、仕様が一般 に公開されていない文字コードを使って いるのなら、悪いことは言わない、さっ さと別の文字コードに乗り換えるべきだ。 閉じた仕様の文字コードは、いずれ誰に もサポートできなくなるのだから。その

点を最後に指摘して、本稿をしめくくることにする。

#### 注

- [1] 中華人民共和国の国家標準局が1980年に発布した漢字コード。94×94の文字表の中に、非漢字682字、第1級漢字3755字、第2級漢字3008字、あわせて7445字を収録していた。現在はGB18030に、とって代わられている。
- [2] Microsoft の独自拡張シフト JIS。
- [3] 高電社版のシフト GB である「CW コード」では、1 ~8区の非漢字部分は GB 2312ではなく JIS X 0208 である。すなわち、非漢字部分はシフト JIS で、漢字 部分だけが Unicode 対応シフト GB だと言える。
- [4] http://www5b.biglobe.ne.jp/~harigaya/chcnv.html
- [5] 『ChineseWriter』の「CW-GB Mincho」や「GB 中国

明朝」、『楽々中国語』の「OM 中国 GB 明朝」、あるいは『NiHao2』の「NH 簡体宋」などのフォントを、Microsoft Windows の「SimSun」に変換し、同時に文字コード変換をおこなう。ゴシックなども同様。

- [6] http://www.mojikyo.org/
- [7] http://10036ra.org/glyph-index.html
- [8] 16~45区の場合は、区番号から16を引いて94倍し、点番号を加える。48~77区の場合は、区番号から18を引いて94倍し、点番号を加える。さらに、フォント名の番号部分(「Mojikyo Mxxx」のxxxの部分)から101を引いて5640倍した値を加えれば、文字鏡番号が得られる。
- [9] グリフ番号 10050001 の gif が http://**10036ra.org/** retrieve.php?gid=**10050001** で入手可能。
- [10] 黒タイ語や白タイ語に用いられる文字。
- [11] U+2123D に収録されている本物の「土」を使えば、 このような重なり表示の問題は発生しないが、住基コ ードでは U+10000 以降を使うことができない。
- [12] Ideographic Variation Sequence。Unicode において複数の字体を見分けるために付加する文字コード。たとえば「丈」と「丈」は、それぞれ「U+4E08 U+E0100」と「U+4E08 U+E0101」という IVS で見分けられる。