

漢字・日本語処理技術の発展： 漢字コードの標準化

芝野 耕司
東京外国語大学
shibano@aa.tufs.ac.jp

1 我が国における漢字コードの標準化

我が国における文字コード、符号化文字集合 (Coded Character Set) の標準化は、もうほとんど「歴史」というに相応しい「とき」を経てきた。最初の文字コードの標準化についての「情報処理」での報告は、和田弘氏が1960年の「情報処理」Vol.1, No.2に電気試験所で1958年に発足したコード会の報告¹⁾に始まる。半世紀を経て、まさに、「歴史」、それも「科学史」の対象とする「とき」に至ったのであろう。

文字コードの検討から半世紀、最初のJIS漢字コード (JIS C 6226, 現JIS X 0208) の制定から四半世紀を経てはいるが、漢字コードの問題は今でも生きている問題である。1つには、日本語を用いるほとんどの人が、日々、漢字コードを用いていることによる。もう1つには、漢字コードの問題は国語と呼ばれる日本語の問題との強い関連から人々の関心を引きつけずにはおかないからである。

ここでは、「科学史」の問題としてではなく、我々が用いている文字コードの標準化がどのように行われ、今日、日常的に広く用いられているJIS漢字コードが最初どのようにして開発されたかについて開発思想を中心に述べよう。なお、JIS X 0208の技術内容に直接関連するものについては、JIS X 0208:1997の解説²⁾に、Unicodeとして知られるJIS X 0221³⁾、JIS X 0208:1997⁴⁾およびJIS X 0213:2000⁵⁾の開発・改正意図については、それぞれ参考文献を参照されたい。

2 情報処理学会漢字コード委員会

和田弘氏は、1971年の論文⁶⁾の中で、「わが国ではカナ文字のためのextension法を必要とするので、早くからこの (ISO/TC97/SC 2, 現ISO/IEC JTC 1/SC 2) 会議で発言しておりましたし、近くは漢字のために1 pageに代わって1冊とでもいふべき拡張法を提案して承認を得ておりま

す」と報告している。

この報告にもあるとおり、我が国における文字コードの標準化は、国際規格との整合性を基本とし、最初のISOコードであるISO R646の1967年の制定を受け、1年半後の1969年にカタカナを含む1バイトコードであるJIS C 6220 (現JIS X 0201) を制定した。漢字コードについても、早くも1969年12月に情報処理学会規格委員会に、国立国語研究所所長をはじめ、文部省で戦後の国語政策の中心を担われた林大氏を主査とする漢字コード委員会を設けた⁷⁾。すなわち、「国語」の担当者に漢字コードの検討を依頼した。

和田弘規格委員会委員長から漢字コード委員会に示された基本方針⁷⁾を原表記に忠実に要約すると次のようになる。

- (1) code extension procedureに準拠し、
 - (2) 2個の7単位符号で1つの文字あるいは記号を表し、
 - (3) 7単位符号の2～7列のうち7/15 (DELETE)を除いた $95 \times 95 = 9,025$ の文字を表し、
 - (4) JIS C 6220 (現JIS X 0201) で規定するローマ文字用符号およびカナ文字用符号を含むが、
 - (5) 現行の符号との関係は原則として考えないと指示している。
- また、漢字用制御符号については、今後の検討とし、escape sequenceについては仮に与えるとしている。
- (6) 漢字コード用文字記号、
 - (7) Sub setおよびSub setとfull setとの関係、
 - (8) key boardとcollating sequence (並び順)
- については、検討を依頼している。

この漢字のコード化の問題に対して、林大氏は、「体・躰・體を同一字の3体と見るか、3字各別のものと見るかが問題となる。……委員会では3主の字形それぞれに、すなわち3字として別々のコードを与えるという考えに傾いている。しかしやはり、……1点1画の差、……単なる筆法の差などについては、これらをそれぞれ別のものとして取り扱うのは、あまりにもむだな感じがする」

と記している。

また、林大氏は、「8,000字を選ぶか、4,000字、3,000字を選ぶかでは、選び方がちがってくる」とし、「8,000字の場合には、制限の色彩は非常に薄くなる。過去の文献を引用復刻するための文字が多くなるが、8,000字という範囲には、1,000字なり2,000字なりの実用的な漢字セットは、どのような目的のものでも、みな含まれることになろう」と書く。

そして、林大氏は、漢字の配列順序、すなわち collating sequence に大きな紙幅を割く。

この第1期漢字コード委員会は、1971年10月に6,100字からなる標準コード用漢字表(試案)を公表する⁸⁾。この第1期委員会の活動報告は、それぞれ情報処理1971年12月号、1972年7月号に掲載された。

この試案の選定方針は、(1)日下部重太郎「現代国語思潮続編」(昭8)(日下部表、6,478字)の付録表にある漢字でB1~8のいずれかに現れるものを選び、(2)日下部表にあってB1~8に現れないもの、日下部表になくてB1~8およびC1~4のいずれかに現れるものの中で、委員2人以上の意見の合致したものを加え、(3)全体として6,000字前後となることを目安とした。

また、異体字については、(4)当用漢字で、新旧字体の著しく相違するものは、両体を採用し、(5)それ以外の異体字は、(1)および(2)の手続きによって拾われる範囲で採用した、とする。

試案で参照したB1~8およびC1~4は、次の通り。なお、B1~8については、後述の林大氏の保存の委員会資料に基づき、それぞれの漢字数を付すこととする。

- B1 大西雅夫「日本基本漢字」(昭16) 3,000字、
- B2 新村出「広辞苑」第2版付録「通用漢字一覧」(昭44) 2,935字、
- B3 朝日新聞社「統一基準漢字書体表」(昭32) 4,000字、
- B4 全日本漢字配列協議会「常用漢字目録」(昭43) 4,000字、
- B5 日本活字鑄造株式会社「標準漢字目録」5,000字、
- B6 国会図書館用NDL70用コード表(昭45) 3,965字、
- B7 共同通信社「漢テレハンドブック」5,694字、
- B8 日経FAM・M型タイプ文字表 4,753字、
- C1 講談社国語辞典付録「漢字音訓総覧」(昭41)、

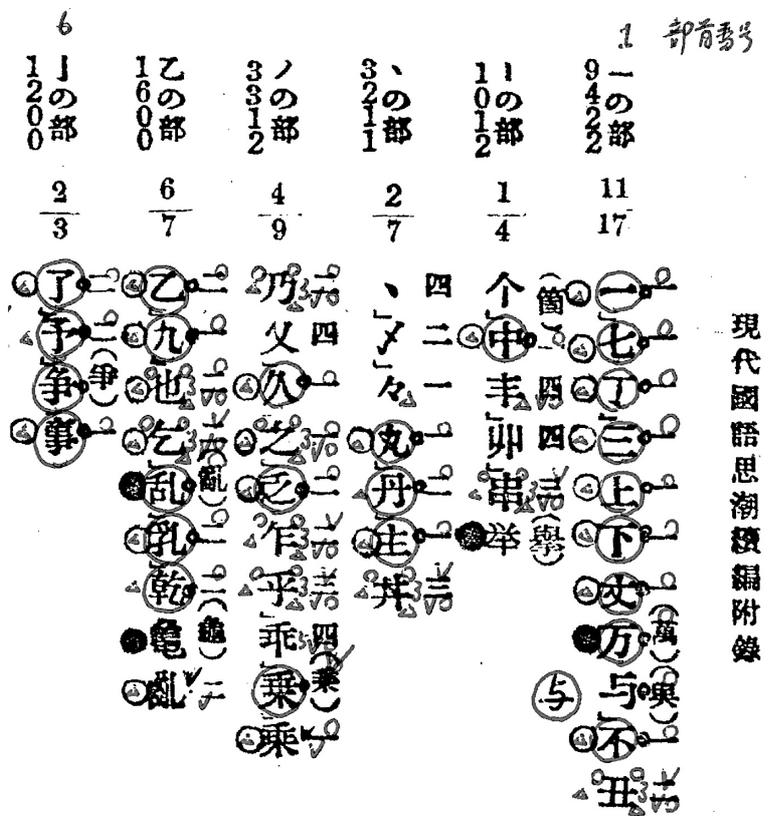


図-1 漢字コード委員会資料

C2 野村広氏「4万5千の姓氏に使われている文字の調査」(昭44)、

C3 国県郡市町村大字名および中学校名に用いられた漢字、

C4 国語研両調査に共通に現れた表外字。

このように当初、記号類を除き8,000字水準を想定して漢字の選定作業を行った委員会であるが、実際のデータから得られる情報から、大漢和辞典の漢字番号順に並べた(康熙字典順の)6,100字の試案ができることとなった。

日下部表、B1~8およびC1~4の参照文字表は、国語分野での常用字調査を基本とし、新聞、活版印刷、図書館、通信社など各分野での漢字利用を考慮に入れたほか、人名、地名にも配慮した漢字表となっていた。

この情報処理学会試案のもととなった資料を林大氏が保存されていたので、これを図-1に示す。このもとなっている漢字表が調査対象で最も文字数の多い日下部表であり、B1~8の漢字表にどの漢字が現れているかをそれぞれの文字の周りに赤字で書き入れたものであ

漢字	官報	辞書	方訓	字体	官報頻度	対応							
						情報処理学会	JIS	内閣府	JIS	大蔵省	国土庁	日本生命	国土行政
一	1	イチ	イチ	一	0	○	○	○	○	○	○	○	○
丁	1	チヨウ	チ	丁	5942	○	○	○	○	○	○	○	○
ナ	1	ナ	ナ	ナ	0	○	○	○	○	○	○	○	○
マン	1	マン	マン	万	8393	○	○	○	○	○	○	○	○
文	1	ブン	ブン	文	141	○	○	○	○	○	○	○	○
ミ	1	ミン	ミン	三	0	○	○	○	○	○	○	○	○
上	1	ジョウ	カエ	上	17795	○	○	○	○	○	○	○	○
下	1	ゲ	シダ	下	13104	○	○	○	○	○	○	○	○
不	1	フ	フ	不	3940	○	○	○	○	○	○	○	○
与	1	ヨ	アザル	与	2119	○	○	○	○	○	○	○	○
大		ジョウ	タケ	大	0								
匹		フウ	ウシ	匹	0								
四		カイ	コウ	四	1	○							
丑	3	チュウ	ウシ	丑	8	○	○						○
且	2	カン	カン	且	219	○							○
丕		ヒ	オキ	丕	3	○							
世	1	セイ	ヨ	世	1500	○	○	○	○	○	○	○	○
丘	1	キウ	カ	丘	197	○	○	○	○	○	○	○	○
丙	8	ヘイ	ヒエ	丙	358	○	○	○	○	○	○	○	○
丙		ヘイ	ヒエ	丙	0								
丞	3	ショウ	ジョウ	丞	7	○							○
兩	1	リョウ		兩	1066	○	○	○	○	○	○	○	○
並	1	ヘイ	ナラベル	並	3212	○	○	○	○	○	○	○	○
个		コン		个	0								
介		カ	カ	介	0	○							
中	1	チュウ	ナカ	中	19154	○	○	○	○	○	○	○	○
卯		カン	コウ	卯	0	○							
串		カン	クシ	串	99	○	○	○	○	○	○	○	○

図-3 対応分析結果

この検討では、情報処理学会試案を基本としたほか、1972年の官報の偶数日の使用頻度と7つの漢字表および調査結果を加え、対応分析⁹⁾を行い、2,817字の行政情報処理用基本漢字¹⁰⁾を選定した。

対応分析結果を図-3に示す。

この行政管理庁の調査研究は、偶数日に限定されているが、官報での漢字使用頻度を1年分にかけて調査したほか、行政で必要とする地名データを国土行政区画総覧から使用漢字一覧としてまとめた。また、日本生命収容の人名漢字をその調査に含めた。

官報頻度の調査は、漢字コード選定のために、漢字表を参照するだけでなく、実際の漢字利用を調査した1次資料を独自に調査したことは高く評価できる。国語学分野では、国立国語研究所の100万字調査が有名であるが、日本の近代活版のもととなったWilliam Gambleも100万字調査を行っており、戦後も大蔵省印刷局や毎日、朝日などの新聞社による活字頻度調査がある。こうした実際の頻度調査によって、工業としての活版印刷の基盤が形成され、工業規格としての漢字コードの基礎資料にもなった。国語学などの「学」ではなく、産業の基礎と

して、政府調達仕様の基盤としてのJISには、このような調査を持つことは必要不可欠であった。

行政管理庁の調査研究のもう1つの特徴は、地名・人名への配慮である。地名に関しては、最も充実した行政地名の集成である国土行政区画総覧で使用されている漢字を独自に調査し、国土行政区画総覧使用漢字として、その対応分析に納めた。また、人名に関しては、漢字情報処理を先行して行っていた生命保険会社の中から最大の生命保険会社である日本生命の収容漢字表を採用した。後述するようにIBM社の文字コードであるEBCDIC漢字コードも、ある保険会社の人名漢字をもとにしている。この点を考慮すると、NTTの電話帳がデータベース化される1980年代以前に得られる最良の人名漢字データをJISの基礎に据えることができたことになる。

この行政管理庁の動きに触発され、JIS原案委員会が設けられることになる。

4 JIS C 6226¹³⁾ (現JIS X 0208) の制定

1974年に森口繁一氏を委員長とし、3年計画でJIS漢字コードの開発を目指した漢字符号標準化調査研究委員会が日本情報処理開発センターに設けられることになる。

この調査研究では、対応分析結果に加えて、新たに12の漢字表を追加し、20の漢字表の漢字を新字源第63版をもとに整理した¹¹⁾。図-4に漢字整理番号表を示す。最終的には、37の漢字表の漢字を整理し、漢字集合の検討を行った。

この委員会によって78JISと呼ばれる最初のJIS漢字コードが開発された¹²⁾。

このJIS原案委員会での検討でも、林大氏を中心に西村恕彦氏らが中心となって検討が行われた。

この検討では、これまでの漢字集合だけの検討ではなく、非漢字集合を含め、文字集合を確定させ、符号化文字集合を完成させたことをまず特筆すべきであろう。また、漢字集合に関しては、“現代日本社会で通用する漢字全体の集約1万弱”とし、“基本度の高いと評価されている漢字集約3,000”を第1水準、さまざまな専門分野で用いられる漢字を第2水準と、漢字集合を2つの水準に分け、符号化する漢字の総数を94 × 94 = 8,836中の7,000 ± αと想定し、漢字集合の確定を目指した。

詳細は、JIS X 0208の解説²⁾を参照されたいが、頻度調査とクラスタ解析によって、第1、第2水準を決定したとしているが、実質的には、情報処理学会試案、行政



4 EBCDIC 漢字コード, シフト JIS, EUC と漢字コードの混乱

当時、世界のコンピュータのシェア 60%以上を占める巨人 IBM は、文字コードでも JIS/ISO/ASCII に対して、EBCDIC という独自の文字コードの世界を築いていた。日本語の世界においても、この EBCDIC に独自にカナ拡張および漢字拡張を行い、ある生命保険会社の頻度情報をもとに、漢字集合とその並び順を決めた EBCDIC 漢字コードを開発した。主要国産各社も英数字に関しては、IBM の EBCDIC に合わせるとともに、アーキテクチャも IBM に合わせ、M シリーズコンピュータの開発を行っていた。カナおよび漢字に関して、日本 IBM から日立などの国産各社に対して、EBCDIC 漢字コードの統一が打診されたが、この統一は実現されなかった。結果として、国内では EBCDIC とはいっても、各社独自の漢字拡張が行われることになる。

ISO/IEC 4873 で規定する 8 ビット符号の構造と規則など基本的な ISO/IEC 2022 の用い方は、エスケープシーケンスを多用するのではない。ISO/IEC 2022 の適合性の箇条でも、実質的には、最小のエスケープシーケンスを推奨している。すると、78JIS の規格票に規定するエスケープシーケンスを最小化したより効率的で自然な符号化が出現するのは、当然である。この代表的な EUC と呼ばれる文字コードである。

ISO/IEC 10646 (Unicode) では、当初 7 ビットの通信路で通信可能な UTF-7 を規定していたが、この UTF-7 では、ASCII がそのままでは利用できない。このため、Unix を含め、多くの処理系での国際化のコストが大きくなるという問題が生じた。このため、ISO では、UTF-7 の代わりに当初 FSS-UTF (File System Safe UTF) と呼ばれた UTF-8 を UTF-7 に代わって、規定した。

日本でも同様に、広く利用されている JIS カナを含む 8 ビットコードの空き領域を用いて漢字を利用可能としたいとの自然な要求が生まれ、シフト JIS がパソコンとほぼ同時期に誕生した。

すなわち、符号化方式としては、JIS、各社の EBCDIC 拡張、EUC、シフト JIS がネットワーク上に併存することとなった。

一方、1983 年に JIS 漢字コードが改正され、83JIS と呼ばれる規格が生まれた。83JIS では、常用漢字表の訓令・国字への対応を主とし、(1)非漢字の拡張、(2)漢字の拡張、(3)第 1 水準と第 2 水準の一部漢字の入れ替え、(4)

常用漢字などの政令漢字についての字形の変更、(5)字形の変更を行った。

83JIS への対応は、この 5 つのどれに対応するかが各社別々のものとなった。たとえば、日立は 83JIS 全部を受け入れ、NEC は (78 年に出された正誤表へは対応しない) 78JIS 初刷を基本とし、(4)の政令字形への変更だけを行った。また、NEC が 83JIS である日立製のプリンタをビジネス用に出したことによって、新旧 JIS 問題が発生する。

DOS/V のもととなり、その後のマイクロソフト漢字コードの基本となった OADG の漢字コードを決めた IBM のシフト JIS では、(1)および(2)だけを採用することとなった。

謝辞 最初の JIS 漢字コードの委員長であり、1997 年の最新の改正に関しても詳細に読まれ、修正意見をお寄せいただいた東京大学名誉教授森口繁一先生が 2002 年 10 月に亡くなりました。森口繁一先生の逝去を悼むとともに、森口繁一先生が JIS 漢字コードを含め、日本の情報処理技術に多大な貢献をされたことに謝辞を捧げます。

参考文献

- 1) 和田 弘, 高橋 茂: コード会のコードについて, 情報処理, Vol.1, No.2, p.107-109, 情報処理学会 (1960).
- 2) 芝野耕司編著: 増補改訂 JIS 漢字字典, 日本規格協会 (2002).
- 3) 芝野耕司: JIS X 0221 (ISO/IEC 10646) の目指すもの—文字コードと日本の国際対応—, 情報処理学会情報規格調査会 NEWSLETTER, No.40 (Dec. 1998), <http://www.itscj.ipsj.or.jp/topics/tp40.html>
- 4) 芝野耕司: JIS X 0208:1997 “7 ビット及び 8 ビットの 2 バイト情報交換用符号化漢字集合” の改正, 情報処理学会情報規格調査会 NEWSLETTER, No.35 (Sep. 1997), <http://www.itscj.ipsj.or.jp/topics/tp35.html>
- 5) 芝野耕司: JIS X 0213 (7 ビット及び 8 ビットの 2 バイト情報交換用符号化拡張漢字集合) の制定, 標準化ジャーナル, Vol.30, pp.3-7, 日本規格協会 (Mar. 2000).
- 6) 和田 弘: 情報処理と標準化, 情報処理, Vol.12, No.2, pp.63-67, 情報処理学会 (Feb. 1971).
- 7) 情報処理学会企画委員会編: 電子計算機の国際標準化 ISO の動きとわが国の歩み, 情報処理学会 (1971).
- 8) 情報処理学会漢字コード委員会: 標準コード用漢字表 (試案), 情報処理学会 (Oct. 1971).
- 9) 行政管理庁: 行政情報処理用標準漢字選定のための漢字の使用頻度および対応分析結果, 行政管理庁 (Mar. 1973).
- 10) 行政管理庁: 行政情報処理用基本漢字に対する符号付与に関する調査研究報告書 [付表], 行政管理庁管理局 (Mar. 1975).
- 11) 日本情報処理開発センター: 漢字整理番号表, 日本情報処理開発センター (Mar. 1975).
- 12) 日本情報処理開発センター: 情報交換のための漢字符号の標準化に関する調査研究報告書, 日本情報処理開発センター (Mar. 1975, Mar. 1976).
- 13) JIS C 6226-1978, “情報交換用漢字符号”, 日本規格協会 (1978).
(平成 14 年 10 月 28 日受付)