

インターネットにおける学術漢字の 符号化に関する基礎的研究

斎藤秀紀・柳沢好昭・横山詔一

国立国語研究所
情報資料研究部
日本語教育センター

技術教育を効率よく行うためには、日本語資料に対する提供手段の開発と分野別で使用される漢字の特長を知ることが重要である。本稿では、最初に構造化4バイトコードが大規模の漢字データベースと専門分野別文字集を一つの枠組みで符号化できることを示した。次に、学術漢字の教材化に適用可能な出現頻度・親密度データベースを紹介した。最後に、海外の日本語学習者が情報を収集し、その体系を発見することに有用であり、かつ、学習者自身が収集した情報から簡便に教材作成が行えるプログラムとリソースのサンプルについて述べた。

Basic Research on the Encoding of Kanji for the Transmission of Technical Information Via Internet

Hidenori Saito, Yoshiaki Yanagisawa, Shoichi Yokoyama

The National Language Research Institute
Department of Data Orientation
Center for Teaching of Japanese as a Second Language

In order to improve the efficiency of education in technical subjects, it is essential to develop the means to transmit technical information in the Japanese language, and to examine the characteristics of Kanji used specialized fields. In this study, we show that it is possible to use a structured four-byte code to encode both a large scale Kanji database as well as sets of characters used in specialized fields, within a single framework. Next, we describe a database which includes indices of frequency and familiarity of characters which could be used in adapting literature using specialized Kanji for educational purposes. Finally, we describe a software system for allowing people studying Japanese overseas to accumulate and organize their personal observations from a variety of media, looking for rules and systematic regularities which connect them.

I 学術漢字の符号化の方法

1. 1 はじめに

本研究では、最初に JIS X0208 および ISO/IEC10646 の問題点を示し、大漢和辞典の検字番号を基本に漢字符号に求める機能を 3 種の構造をもつ 4 バイトコードで表現できることを示した。次に、4 バイトコード（整数部 3 バイトと小数部 1 バイト）の小数部に異体字と中国・台湾・日本・韓国の情報交換用文字を配当する 1 字体 1 符号化法が、文字集合の全体から部分集合を抽出し、2 次符号化する方法に拡張できることを述べた。さらに、文字集合の部分は、専門分野別に漢字を符号化できることを示した。

情報交換用漢字符号 JIS C6226 は、(1987 年 JIS X0208 に名称変更) 1978 年に日本工業規格として制定された。現在 JIS X0208 は、利用方法の多様化とともに、多国語表現や古典・漢籍などにも使用されている。しかし、大規模の漢和辞書を始め古典・漢籍を電子化するためには、文字種の不足が指摘されていた。また、1983 年の規格改正では、水準間で 22 組みの字形入れ替えと 294 字の字形変更があり、新旧規格の間で互換性を崩した。問題には、符号間への文字の追加機能や変更情報の局所化など外界の変化に対する調整機能や、文字と符号に対する規範性に対する考え方の曖昧さ、異体字の整理に関するものがある(斎藤:1985, 1993)。

一方、国際標準化機構は、世界共通の情報交換用符号の利用を目的に 1993 年に ISO/IEC10646-1 を発表した。しかし、ISO/IEC10646-1 は、中国・台湾・日本・韓国で使用されている国内規格を 2 バイトコードの 20,992 の枠内におさめようとしたため新たな問題が指摘されている。問題は、各国の文字政策の変更と ISO/IEC10646-1 への反映方法や、漢字の読み・画数・部首などの属性情報の規定に関するもの、既存の漢字符号からの移行や、各国語の混在と個別処理を指定に係わるもの、符号化規則に対する「ソースコード分離」や「字形の統合」の併用に係わるものがある(斎藤:1994a, 1994b)。

1. 2 4 バイトコードの構造化

4 バイトコードの符号化領域は、JIS X0208 (G0 領域) や拡張 UNIX コード (G1 領域) との併用を考慮し、G3 領域を使用する (図 1)。G0 から G3 の各符号の識別は、2 の 8 ビット目を使用した。10 進数で表現された検字番号の内部符号への変換は、94 進数と 16 進数の二重変換を行い (以下 16 進数 94 進数変換と呼称) 16 進数 '21' から '7E' に調整した。符号化できる字数は、整数部 830、584 字と小数部 94 字である (式 1、2)。4 バイトコードの構造は、整数部が (1) 漢字辞書やコードブックに付けた 10 進数 5 桁の検字番号を 3 バイトで表す構造 (図 2-1)、(2) 既存の 2 バイト系漢字符号を切り替え符号で統轄する構造 (図 2-2)、(3) 各桁の 2 の 8 ビット目が '01' である 2 バイトコードを 2 個使い内部符号とする構造 (図 2-3)、(4) 2 バイトコードと 4 バイトコードの各桁 2 の 8 ビット目を識別用とする 4 種である。小数部は、整数部に 1 バイトを付加した構造を基本に、(5) 異体字と各国語の情報交換用文字を配当する構造である (図 2-A)。

$$\text{整数部} = \text{HEX}(\text{検字番号} \bmod 94) + \text{重み} '21' \quad (\text{式} 1)$$

$$\text{小数部} = \text{HEX}(\text{小数部挿入位置番号}) + \text{重み} '21' \quad (\text{式} 2)$$

HEX: 10 進数値を 16 進数に変換する関数

MOD: 10 進数表現された検字番号の剰余を求める関数

検字番号: 10 進数 5 桁で初期値が '0' の数値

G 0	G 3	表外字コード領域
I-00	I-01	
G 2	G 1	内字コード領域
I-10	I-11	

G 3
G 0、G 1、G 2

図1 2バイトコードと4バイトコードの符号化領域

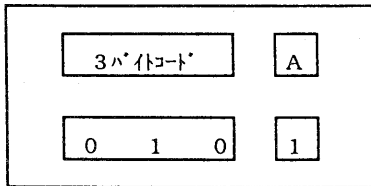


図2-1

検字番号を4バイトコード化した構造
符号の識別ビット列

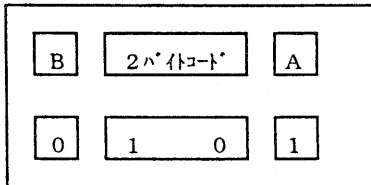


図2-2

JIS X 0 2 0 8を基本とする符号の構造
A: 文字追加用符号
B: 各領域の識別符号

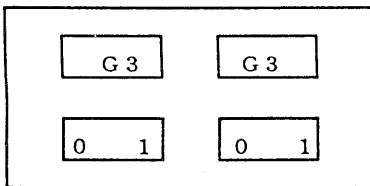


図2-3

G 3を基本とする内部符号の構造

図2 4バイトコードの構造

1. 3 1字体1符号化法

漢文字符の1字体1符号は、4バイトコードの整数部3バイトに見出しに相当する字体（特定の文字で代用）を対応させ、中国・台湾・日本・韓国の情報交換用文字集合や異体字を小数部1バイトに配当する方法をいう（斎藤:1994a）。この符号化法は、変更の多い情報を見出しから分離し、見出しと漢文字符の関係を安定させる効果がある。また、旧規格の漢文字符で作成したデータの変更情報の履歴を小数部に累積させることによって継続使用を可能にする。4バイトコード小数部への文字配当は、符号化の対象になる文字集合を部首順に配列したのち小数部‘1’から‘4’に中国・台湾・日本・韓国語の情報交換用文字を登録する（図3）。‘5’から‘94’には、新しい文字の追加機能を基本に見出しに対応する異体字の全リストを登録する。

見出し	情報交換用文字				異体字登録領域										
0	1	2	3	4	5	6	7	8	9	10	11	12	13	----	94
劍	G	T	J	K	劍	劍	劍	劍	劍	劍	劍	劍	劍	劍	劍

図3 小数部による文字指定の方法（中国：G，台湾：T，日本：J，韓国：K）

1. 4 基本漢字の選定と符号化の方法

4バイトコードで表現できる文字集合は、多様な文字集合を符号化できるが実務で使用するためには柔軟性に欠ける。文字集合や符号は、目的にあわせ自由に設定できることが重要である(斎藤:1995)。分野別や専門別に文字集合を規定することは、専門的な読みや漢字を一般利用者に理解させる必要がなくなり、利用者の負担軽減と外国人に対する専門分野別の教育への利用が期待できる。利用者規定の方法は、図2-1に対応する文字集合を規範とし、図2-2の構造の枠組みのなかで2バイト符号を与える。異なる符号間での情報交換は、規範となる4バイトコード対応の文字集合を中間符号に使用する。

現在、情報科学で使用している日本語から基本漢字を選定する作業の一貫として、日本語・中国語・英語で使用される用語を入力している。この資料に対して、以下に述べる分析を試み、日本語学習リソースの開発を目指している。

II 学術漢字の認知的特性の解析に向けて

認知心理学では、語彙の習得において学習材料の出現頻度や親密度(familiarity)が学習効率に大きな影響を及ぼすことが知られている。親密度とは、学習者が当該の言語刺激にどの程度見慣れているかを示す心理的指標を言い、評定尺度法で測定される。

日本語学習者のための語彙学習教材を開発するには、語彙の出現頻度や親密度の資料が重要な役割を果たすと考えられる。以下、日本語学習の分野に適用が期待できそうな出現頻度・親密度データベースを紹介し、学術漢字の教材化にどのように応用できるか簡単に述べる。

<出現頻度データベース>

漢字および漢字単語の出現頻度については、国立国語研究所の語彙調査が諸学界から高い評価を受けてきた。しかし、それらは約30年以上も前に実施されたものであり、時代的に現状にあっていないのではないかという声を耳にするようになってきた。国立国語研究所の語彙調査データは1990年代にも果たして通用するのだろうか。

この疑問に答えるため、野崎・横山・磯本(1995)、野崎・磯本・横山(1995)は1993年の朝日新聞記事CD-ROMテキスト・データベースを利用して漢字および漢字単語の出現頻度を調査中である。現在、1993年の朝日新聞記事CD-ROMに格納された約11万件の記事に登場する漢字1文字の出現頻度について集計を終えた。次の段階では、国立国語研究所が実施した1966年の新聞語彙調査データとの相関を分析し、漢字の出現頻度に時代差がどのような影響をもたらしているのか明らかにしていく。

<親密度データベース>

親密度については、これまで大規模なデータベースの整備には手が付けられていなかった。しかし、Amano, Kondo & Kakehi (1995)により約6万の見出し語を有するデータベースが構築されつつあり、状況は好転した。このデータベースは、文字データベース、単語データベース、音声ファイルデータベースの3要素から成り、それらが相互に参照可能である。文字データベースにはJIS X0208に規定された約6800文字に対する親密度が納められている。単語データベースには単語を文字で被験者に呈示したときの親密度と、音声で呈示したときの親密度のデータのほか、アクセント、音韻表記、文字表記などの情報が約6万語のそれぞれについて含まれている。音声ファイルデータベースには、単語データベースの見出し語に対応する音声情報が格納されている。

<学術漢字の教材化に向けて>

以上で述べた出現頻度データベースと親密度データベースを利用すれば、情報科学分野における学術漢字の認知心理学的特性をかなりの確に把握することが可能となる。その知見に基づいて教材開発に着手すれば、合理的な日本語教育教材を手にすることができると期待される。

Ⅲ 日本語学習リソース開発

従来の学習観では、外国語習得は学習コースに参加しクラスワークを行うもの、教師によるコースデザインがあり、それに沿って学習を進めるのが一般的である。しかし、学習により新たな知識や技能を獲得するという学習効果が海外の様々な学習者すべてに同じように、どのような学習においても、いつでも現れるわけではない。学習者の個人差やニーズに対処するには、できる限り多くの個人差の次元（学習活動：学習者の認知的・情意的前提能力、学習課題、授業活動、到達水準、学習速度、情意的成果等、個人内要因：知能、認知スタイル、性格、情意的認知、動機付け、原因帰属、効力感、統制感、やる気、関心度、集中度、適性、交友、生活意識、情報交換網、自己理解等）について考えることが前提となる。

一方、現地国の日本語指導者、あるいは仕事や研究や将来のために学習する一般成人は、自分の学習や研修について検討し、全体計画を立て自律的に学習や研修を行う能力を備える。法則性や体系化の習得は、一方的に伝えるより学習者側が発見する方が定着度と情意面において効果的であり、いろいろな情報に接し、検索、複眼的に整理、分析、組み合わせる新しい発見を生み出すと考える。

学習を、（１）空間的知覚能力、記憶力、精神速度、推理能力、運動能力等の基本的能力からとらえる活動、（２）収集した言語行動情報に存在する規則や体系の結論を学習者自身が発見する過程に参加する活動、（３）日常の情報の能動的処理能力を尊重し、複数の人間が協働して解答を創り出す活動、を前提に、学習支援システムをコンサルトと人的・物的・社会的リソースと分け、視聴覚的アピール度、付加情報のインデックス化、開架式、集積・再利用と自己開発の簡便性を持つリソースと、その交流について試行する。その際、コンピュータ・ソフトは、操作の簡便性、視覚的アピールの高さ、ハードへの依存度の低さ、アクセスの容易さを要件とする。

本研究では、各種のメディアを合成・編集するツール、「Intellegent Pad」（開発者：北海道大学工学部田中譲教授、以下、IP）を用いてリソースを開発する。IPは、すべての事物、絵やテキストや演算などのプログラムあるいはアプリケーションプログラムをパッドで表すものである。そのパッドは、それぞれ個別の機能を持ち、CRT上で紙面のごとく表され、マウスで移動、複写、結合・貼り合わせ、拡大縮小や固有の機能を操作ができ、コンパイルなどの操作なしで機能合成や複合が行える。OSやハードの垣根をも越える。このツールを使い、サンプル作成を行う。

作業１：日本語学習者の耳の上に小型CCDカメラを装着し、本人が希望する行動場面を録画し、そのデータをキャプチャーする。リンク・パッド、動画再生パッド、ビデオ／サウンド・パッド、開発ツール・パッドなどを使い、文字化資料、語彙や画像などのデータベースをリンクしたマルチメディア型素材を学習者自身が主体的に作成し、蓄積する。

作業２：学習者と指導者が作成物を視聴し、think-aloudによる内観活動を行い、素材の追加データとする。

作業３：上記素材の送受信を行う。

【構成図】

録音録画データ	ビデオ/サウンド・パッド	* 学習者自身が録画したもの (CVD-500 と CCD-MC1 (SONY))
文字化データ	テキスト・パッド	* 学習者が主体的に書き起こしたもの
語彙データ	辞書・パッド	* 既存の語彙データベースの活用
文化情報データ	イメージ/ムービー/CD・パッド	* 非言語行動・文化事情等のデータ
付加情報データ	テープ・パッド	* 指導者等が補足的に付け加えた練習問題等
追加情報データ	テープ・パッド	* think-aloud 等のデータ
ベース・パッド		
カーネル (Intellegent Pad)		
OS (Windows3.1/ 漢字Talk7.0/System7.0)		Intellegent Pad for Mac 1.0、漢字Talk7.5
ハードウェア (IBM/ATorMac)		PowerMac 8100/80AV (Apple)、1GB (HDD)、32MB (RAM)

参考文献

- 1) 斎藤秀紀(1985)「漢字コードの拡張法に対する試案」『研究報告集6』国立国語研究所報告83, pp. 57-103
- 2) 斎藤秀紀(1993)「漢字コードのメタコード化の方法」『情報処理学会第46回全国大会論文講演集(1), pp. 23-24』
- 3) 斎藤秀紀(1994a)「1字体に1符号を対応させる漢字符号化の方法」『計量国語学会誌第19巻5号, pp. 223-233』
- 4) 斎藤秀紀(1994b)「大漢和辞典の検字番号に基づく構造化4バイトコードの提案」『情報処理学会論文誌 Vol. 35, No. 6, pp. 1119-1121』
- 5) 斎藤秀紀(1995)「4バイトコード対応文字の部分文字集合に対する利用者規定の方法」『情報処理学会第50回全国大会論文講演集(3), pp. 201-202』
- 6) 財)札幌エレクトロニクスセンター(1994)「STEP NEWS vol. 83」
- 7) 中野照海編著(1979)「教育学講座第6巻: 教育工学」学習研究社
- 8) 織田守矢ほか(1990)「学習環境の構築」コロナ社
- 9) 今栄国晴編著(1993)「教育の情報化と認知科学: 教育の方法と技術の革新」福村出版
- 10) アルフレッド・トマティス(1994)「人間はみな語学の天才である」アルク
- 11) 後藤忠彦(1986)「教育とコンピュータ5: コンピュータと教育情報システム」東京書籍

引用文献

- 1) 野崎浩成・横山詔一・磯本征雄(1995)「KWIC作成プログラムの開発」『教育システム情育報学会第20回大会発表論文集, pp. 211-214』
- 2) 野崎浩成・磯本征雄・横山詔一(1995)「日本語教育のための語彙調査」『日本教育工学会第11回大会発表論文集』
- 3) Amano S., Kondo T., & Kakehi K. (1995) 'Modality dependency of familiarity ratings of Japanese words' *Perception & Psychophysics*, 57(5), pp. 598-603