

# 多言語を統一する情報交換用4バイトコードの研究 - 漢字の属性情報に対する符合化法の提案 -

国立国語研究所  
齋藤秀紀

## 1. はじめに

JIS C6226-1978 は、1978 年に日本語処理のための情報交換用漢字符号として制定された。(1987 年 JIS X0208 に改称)。しかし、JIS X0208 は、漢字符号に連続番号を使用しているため、文字の追加・更新に対する初期状態や履歴を管理することが困難な状態にある。また、1993 年に国際標準化機構が発表した国際符号化文字集合 (ISO/IEC 10646-1, JIS X0221-1995) は、各国の漢字識別方法を規定していないため、実装される文字フォントが包摂基準による「字形統合」の範囲が「ソースコード分離規則」を併用するのが曖昧である。さらに、データに各国語が混在しているデータは、符号からはどの国の漢字であるのが判別できない。

本論文では、調査・研究や漢字字典に記載されている漢字の属性情報(属性)を漢字データベースで規定し、これらを形式的に符号化する方法を提案する。符号化実験は、サーバに置いた JIS C6226-1978 の 6349 字に 40 種の属性情報を附加した漢字データベースを基本に、以下に示す 16 種の属性情報を対象に行った。

(1) JIS X0208-1983 で改定された情報、(2) 総画、(3) 部首部分を省いた画数、(4) 教育漢字識別符号、(5-6) 学習漢字の識別符号 2 種、(7-9) 当用漢字の識別符号と各補正符号 2 種、(10) 常用漢字識別符号、(11-14) 人名漢字識別符号 4 種、(15) 国名、(16) 読み。

なお、中国語は、GB 2312 を基準に読みはワークステーション(日本電気株式会社製 EWS4800/32OEX)付属の漢字入力システム用電子辞書から引用した。

## 2. 漢字データベースと転置ファイル

文献・資料の電子化と再現処理には、明確に規定された符号と文字の関係をデータ作成者と利用者の双方が共通に使用できる

環境をもつことが重要である。また、符号と属性情報を安定状態で運用するためには、情報の追加や更新に対して符号と文字の初期状態が維持されていることが必要になる。そのためには、符号が、装置で使用する内部符号や漢字データベースの編成方式から独立していることが要求される。

本論文では、漢字データベースの索引として設けた転置ファイルから形式的に属性情報を符号化するため、漢字データベースで規定した属性情報の符号化(以下、属性符号)に対して3層モデルを導入する。符号長は、4バイトコードとし、構造要素、要素を結合した構造、構造間の関係に対応させ、構造の要素には符号に要求する機能の単位を当てた(図1)。

第1層は、論理符号を漢字データベースの編成方法から独立させるため、表1に示した転置リストの構成要素である属性情報と漢字データベースへの接続情報の2項関係を論理符号とした。属性情報の位置は、漢字データベースに記載されている見出しからの位置を項目番号とし、同一属性をもつ場合は前項目番号からの相対位置(最大94種)の二つを使った。相対位置は、データへの属性情報を埋め込む処理用に使い、1バイトで表した。また、漢字符号は、3バイト(830,584字)で表現した。

第2層は、属性情報1バイト論理符号の3バイトを結合したものと拡張 UNIX コードの2バイトを結合した符号系である。拡張 UNIX コードは、クライアントで使用する漢字符号であり、論理符号はサーバで使用する符号である。また、データに埋め込む属性情報は、漢字符号部分の拡張 UNIX コードと論理符号をデータ符号と共用させ、属性を表す1バイトを使った。

第3層は、装置で使用する内部符号である。既存の漢字符号との併用を図るため G3 領域を使った。

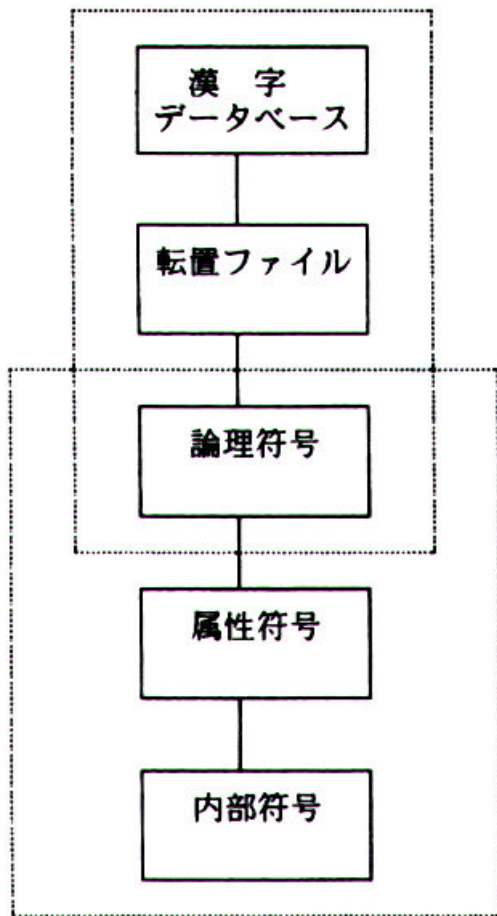


図1 論理符号と属性符号に対する3層構造

表1 転置リストの一例

音読み	個数	同音の漢字
ア	13	亜啞娃阿雅蛙窪亞埜猗両鋸闕
アイ	20	娃哀愛挨喝乃哇噫埃曖欸曖矮穢 藹阨隘霽霽鞋
アク	14	啞阿悪握渥厄埜幄惡扼輒阨鸞齧
アチ	1	藹
アツ	10	亜圧鞆亞壓扎藹軋遏闕

- - 以下略 - -

訓読み 個数 同訓の漢字

アア	1	乎
アア	11	悪於于吁咨嗚嗟噫惡猗羌
アイ	7	始逢間際相藍胥
アイダ	1	間
アイト	2	対對
アウ	13	逢会遇合遭值會翕觀覲逅邁避

- - 以下略 - -

画数 個数 同一画数をもつ漢字

01	5	ノ一乙、
02	29	勺丁儿十冂冂力口乃；ヒ口 刀冂冂ムセト人又二九メ八几
03	59	宀女士也勺子下巾才寸又夕上 与干巳己凡子士万山久乞刃大 及三弓小亡丸《口井川九么广 廴丈夕口兀于工弋中夕シ彳个 女尸千巳刃
04	109	刈文分比父少日天手匹夫方氏 夫无支女升死片尺仏什不仄厄 受毋气孔亢欠从仍式月歹公仆 仇犬午互元五幻尸勾彳支仇毛 木牛勿卞尤勾丐卅止友予彳斤 爻今区天兮六凶之双中火彡丑 弔屮井冗允尹巴円廿引日爪扎 切心云化戈斗水反牙壬太王丹 屯介句内仁

- - 以下略 - -

### 3. 属性情報の符号化の方法

選択された属性情報の管理は、属性情報の見出し漢字からの位置を示す項目番号と項目番号からの相対位置の二つを使った。項目番号は、属性情報に設けた分離記号'、'(カンマ)を数え検出した。複数の属性値をもつ場合は、項目番号からの相対位置による検索で得た。また、漢字データベースから属性情報を読み出す処理は、制御レコードの項目番号とデータの相対位置を使った。項目番号は、ヘッダーの一部として使用する制御レコードに記録し、該当するデータの管理に使用した。

4バイトコードの漢文字符号は、3バイトで漢字を表し、4バイト目を47個の異体字登録と、以下に続く情報が「属性情報なし」か「属性情報付き」かを判別する機能にを使った。漢文字符号と属性情報の接続状態は、以下の組み合わせとなる。

#### (1) 属性情報なし2バイト漢文字符号

全てのデータを拡張 UNIX コードで表現する。

#### (2) 属性情報なし4バイト漢文字符号

1バイト目から3バイト目に論理符号の3バイト部分を埋め、4バイト目に「属性情報なし4バイト漢文字符号を表す'A1'から'EC'を入れる。

### (3)属性情報付き 2 バイト漢字符号

1 バイト目と 2 バイト目に拡張 UNIX コードを入れ、3 バイト目と 4 バイト目に「属性情報に対する接続情報」をす '21D0'を入れる。

### (4)属性情報付き 4 バイト漢字符号

1 バイト目から 3 バイト目に論理符号の 3 バイト部分を埋め、4 バイト目に「属性情報付き 4 バイト漢字符号を示す'D1'から'FE'を入れる。

## 4 . プログラム機能と実験の概要

実験装置は、ワークステーション 1 台にクライアント・サーバ・モデルを実装した。符号化実験は、漢字データベースで規定した属性の符号化と位置情報から属性情報を再現する二つの処理を行った。本実験では、クライアント側漢字データベースに登録されていない漢字がある場合は、サーバの漢字データベースから補填した。日本語と中国語の識別は入力文字がカナの場合に日本語とし、ピンインの場合は中国語とした。属性情報のデータへの埋め込みは、プログラム起動時に画面表示された 16 種のリストから選択した最大 8 種の属性情報を表す位置情報を指定順に結合した。クライアント側の漢字符号は、拡張 UNIX コードを使い、サーバでは 4 バイトコードを使用した。

実験用プログラムは、サーバ(4BServer)用とクライアント(4BText)用の二つを作成した。4BServer では、漢字データベースの検索と管理を行った。4BText は、(1)ファイル入力、(2)ファイル出力、(3)日本語や中国語に対する読み漢字変換、(4)データと属性情報の表示、(5)読み情報の修正処理、(6)新旧漢字の選択表示機処理、の各機能をメニュー - 選択方式で実行する方法をとった。

図 2 の左画面は、テストデータの内容である。右画面に、JISC6226 で規定した区点番号、総画情報、部首を省いた画数情報(内画)、テストデータで使われた漢字の読みと国名を表示した。国名は、日本語を、'J'で表示し、中国語は'C'を使った。画面右側に示した項目名「4 バイト漢字符号」と「2 バイト漢字符号」はサーバの論理符

号とクライアントで使った符号である。画面中央の副画面は、漢字データベースに日本語と中国語を入力する場合に同じ読みをもつ漢字の候補を示したものである。読みの入力は、副画面の変換方法指定ボックスから「音」、「訓」または「ピンイン」を指定し、次に、読みと漢字を選択する方法をとった。

## 5 . おわりに

本研究で提案した符号化法を実用化するためには、属性符号を 1 バイトから 2 バイトに増加させることや間接指定法を標準化し、分類用レコードの作成を省略することが必要になる。また、属性と漢字の 2 項関係を符号化基準に使用する方法は、ハンゲルや漢字と異体字についても符号化できる可能性をもつものであり、多言語を異体字と同じ枠組みで符号化するための重要な課題となる。これは、今後、漢字の異体字と属性情報を一つの符号化法で処理する実験で確認する予定である。

[付記] 本研究は、平成 9 年度・文部省科学研究費(創成的基礎研究費)(「国際社会における日本語についての総合的研究」研究代表者 水谷修, 課題番号 09NP0701)の分担として交付を受けた。

## 参考文献

- [1] 斎藤秀紀 漢字情報と文例情報を結合した日本語データベースの構築・情報処理学会, 人文科学とコンピュータ研究会, Vol.96.No.42, pp.35.-40(1996).
- [2] 斎藤秀紀 4 バイトコードと対応文字の部分集合に対する利用者規定の方法, 情報処理学会第 50 回全国大会講演資料集 pp.201-202(1995).
- [3] 斎藤秀紀 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, 第 19 巻 5 号, pp.223-233(1994).
- [4] 斎藤秀紀 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.9, pp.1119-1126(1994).
- [5] 斎藤秀紀 漢字の属性情報に対する符

号化法の提案，計量国語学（審査中）

The screenshot shows a software window titled "4leat" with a menu bar containing "ファイル入力", "ファイル出力", "キーボード入力", "コード表示", "読み設定", "字体変更", and "終了". The main text area contains Japanese text, including "月日は百代の過客にして、行きかふ年も又見出し文字".

Overlaid on the table is a "キーボード入力" dialog box with the following content:

検索  
 月, ゲツ/ガツ/, ツキ/, yue4/yue/  
 朝, キノ, ツキノ,  
 突, トツ/, ツ(ク)/ツ(キ)/, tui/tui/  
 付, フイ, アタ(エル)/ツ(ケル)/ツ(ク)/ツキ/ツケ/, fui/fui/  
 拜, ハイ/, オカ/ツキ/, hui/  
 変換方法  
 v 音読み + 訓読み v ビンイン  
 入力 つけ 候補数: 6  
 OK CANCEL

The table in the background contains the following data:

月日は百代の過客にして、行きかふ年も又見出し文字																				
A1	B7	C6	A4	C9	C2	A4	B2	B5	A4	M4	M4	A1	B9	M4	M4	C7	M4	C8	2	バイトコード
A1	EE	FC	CF	B4	E5	CE	E1	D2	C8	B7	C6	A2	D4	A0	AB	D5	AF	E2	F4	
31	30	3A	21	40	29			47		28	24	4	バイトコード							
85	D3	DF	CA	AE	A7			D1		C4	D2									
40	78	5F	27	34	30			28		3F	45									
A1	A1	A1	A1	A1	A1			A1		A1	A1									
23	38	41	34	18	21			25		39	43	JIS	区点番号							
78	92	20	69	65	50			52		15	84									
										6	6	2	読角							
										0	3	0	内典							
										ユ	ト	マ	読み							
										(	シ	タ								
										ク										
										)										
J	J	J	J	J	J			J		J	J	国名								

図2 属性情報を付加したデータ