

字形検字方式による漢字入力の一方法

磯道義典 (広島大学)

新谷 順子 (日立ソフトウェア
エンジニアリング(株))

1. はじめに

現在、さまざまな漢字入力法が考案され、⁽¹⁾⁽²⁾⁽³⁾⁽⁴⁾⁽⁵⁾ 次々実用化されているが、次の様な利点を持つ字形検字方式による漢字入力の一方法を考案したので、報告する。

(1) 入力装置が小型である。

(2) オペレータの負担が軽い。

(3) 入力できる漢字数が比較的多い。

本論文では、まず字形検字方式を選択するに至った理由を述べ、その例を示した上で、コード作成の過程を説明する。

2. 字形検字方式

漢字を計算機で処理する場合、その入力装置として考えられるものに、まず、普通のTSS端末に備わっているキーボードがある。これは、一般に漢字入力に用いられているタブレットとかOCRといった装置に比べて、小型で安価であるという利点がある。

常用キーを用いる場合はさらにどのように漢字を指定するかという問題が残る。一番手っ取り早い方法としては、テキストをすべて仮名文字化して入力し、漢字へ変換するのは計算機まかせという、仮名漢字変換方式がある。しかしこの方法は、日本語の特徴である同音意義語の処理のため、構文解析や意味解析にまで立ち入ったかなり複雑な変換プログラムが必要となる。

逆に計算機内での処理が最も簡単にすむのが、漢字に対してコードを割り当てそのコードを入力するというコード入力方式である。この方法はコード

を暗記しなくては使えないという取っ付きの悪さがあるが、一度コードを覚えてしまうとかなり高速の入力が可能になる。

他に、コード入力方式におけるコードを暗記する手間を省いた、対話型入力方式がある。コード入力方式のように漢字とそのコードが1対1に対応しておらず、コードに対して漢字が一意に定まらない時の最終決定は、改めてオペレータが入力するようになってくる。この方法は初心者には使いやすいが、熟練してもあまり速度が上がらないことが多い。

この3つの方式のうち、本研究ではコード入力方式を考えた。何故なら、コードを暗記するというオペレータにとって一見不利な条件も長い目で見れば、入力速度が向上し負担が軽くなることにつながるからである。

さて、コード入力法における最大の問題は、漢字にどのようなコードを割り当てることができるかである。この漢字のコード化にも、ただ順に数字コードを割り当てるという実用にむかないものから、人間工学的にかなり考慮されているラインナップ方式まで、いろいろの方法がある。本研究では入力対象として約4000字種の文字を前提としていたので、比較的多い字種を入力でき、また速度の点においても期待できる字形検字方式をとった。

これは漢字をその構造に基づいてコード化するものである。コードを全部丸暗記しようとする、人間が記憶できるコードの数には限界があるため、入力できる漢字数が限られてしまう。この点、字形検字方式では、コード化の規則さえ覚えれば視覚的にコードを

導き出すことができるので、比較的多い字数の入力が可能となる。また訓練すれば反射的にコードが出るようになり、入力操作の際原稿だけを見て打つタッチ打鍵に近づけることができる。

3. 字形検字方式の例

字形検字方式では、漢字をその構造に基づいてコード化するわけだが、その方法にも2通りある。一つは、字根（漢字を構成している基本的な要素）とともにそれらの結合操作を指定する方法である。もう一つは、字根を並べてコード化する時に漢字の特定の位置から字根をひろってくる方法である。前者は、コードの長さが不定で、複雑な漢字ほどコードが長くなってしまふという欠点があるので、本研究では後者を採用した。

後者の方法において、字根の数と字根をひろってくる位置は研究によって異なり、それによって使用するキーの数や打鍵数が異ってくる。例として、四角号碼法と三角コード法がある。

四角号碼法とは、漢字がだいたいの正方形で四つの角を有しているのに着眼した王雲五が考案したものである。〔1928年〕まず10種類の基本字根を定義し、0から9までのアラビア数字で表現した。次に各漢字の四角の字根を左上、右下、左下、右下各角の順にひろっていき、4位の10進数で漢字をコード化した。四角号碼法の基本字形とコーディングの例を図1に示す。この方法は異なる字で同コードになるものが多く、すべての漢字を入力するのは困難である。

これに対して、胡（L.Hu）らは三角コード法を提唱した。〔1979年〕胡はまず300個の基本符号を選び、99のグループに分け01から99ま

で10進数2桁のコードを与えた。各グループ内の基本符号は皆類似した形を有し、しかも一桁目が四角号碼に関連している。次に、コード化の規則としてはZ形の角コード化原則を決定した。それは左から右、上から下、外から中へといった順序に従って23つの基本符号をとらえ、そのコードを配列するものである。すなわち6桁の10進数で漢字をコード化することができる。この方法では同コードで異った文字の場合が少なくなる。

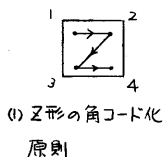
コード	0	1	2	3	4	5	6	7	8	9
基本字根	上	一	丨	十	キ	ロ	フ	ハ	小	
			ノ	丨	又					
			ノ	丨						

(1) 基本字根とそのコード

$\begin{matrix} 2 & 3 \\ 2 & 3 \end{matrix}$ 伏 ₃	= 2323	$\begin{matrix} 0 & 8 & 6 & 1 \\ 6 & 1 \end{matrix}$ 説 ₁	= 0861
$\begin{matrix} 3 & 1 \\ 1 & 4 \end{matrix}$ 汗 ₄	= 3114	$\begin{matrix} 9 & 7 \\ 2 & 5 \end{matrix}$ 輝 ₅	= 9725

(2) コーディング例

〔図1〕四角号碼法



1	馬 = 馬	35 00 00
12	什 = 丨十	22 40 00
13	法 = 丨士 丨	33 41 74
132	郝 = 丨 丨 小	41 15 91

(2) コーディング例

〔図2〕三角コード法

4. コード化の手順

本研究では次の手順で漢字のコード化を行った。まず漢字を構成する基本字形を26のグループに分け、A～Zの英字1字で表わした。そして各漢字から最高3つの基本字形を選び、書く順序に従ってこれらのコードを配列した。同コードとなる漢字を区別するためには数字、記号を補助的に用いた。その結果漢字約3000字を42キーの2～3打鍵で入力することができた。

なお、仮名、数字、記号、英字などはすべて2打打ちとし、オ1打目はキーボード上の文字に従い、オ2打目は記号を用いた。

本コードの特徴を次にあげる。

(1) 漢字の視覚的形を重んじている。

基本字形を26個のグループに分ける際、筆画の形によって視覚的に分類した。

(2) 良く使うものは短く、あまり使われないものは長いコードとなる。

漢字は2～3打鍵で入力するが、2打鍵となるのは基本字形1つで構成される漢字、すなわち比較的出現頻度の高いものとなっている。

(3) 言語一般において頻繁に使うものは定形からはずれずよいか、使われないものは簡単な規則にのっとり形作られている。例えばbe動詞などは規則にのっとりないが、あまり表われない動詞では3人称は-S、過去過去分詞は-ed(-en)となる。この点をコード化でも重視している。

コード化の規則は、基本字形1つで

構成される漢字、2つで構成される漢字、3つ以上で構成される漢字のそれぞれで異なる。基本字形の数が少ないものほど、本来の英字コード以外に数字や記号を補助的に用いており、定形からはずれている。

5. 基本字形のグループ分けとそのコード

漢字はいくつかの基本字形から成る。例えば「胃」は田と月から成り、「漢」はシとサと果から成る。このとき、田、月、シ、サ、果をそれぞれ基本字形という。

基本字形はその筆順のオ1画目～オ3画目の形によってグループに分けた。そのため主に左上の角あたりの形が類似した基本字形が同一グループに入る。

具体的な過程を示す。まず漢字の筆画はその形によって次の5種類に分類できる。

「丶」、「一」、「丨」、「ノ」、「フ」

このことから基本字形をオ1画目の形によって5つのグループに分けることができる。さらにオ2画目、オ3画目の形によって最終的に26個のグループに分けた。表1に各グループの筆画の形とそのグループに入る基本字形の例、そしてグループに与えられた英字1文字のコードを示す。

特にオ1画目が「一」(よこ)、オ2画目が「丨」(たて)のグループ(E～I)、オ1画目が「丨」(たて)、オ2画目が「フ」(かぎ)のグループ(M～P)に入る基本字形はいろいろな漢字に多く使用されているため、漢字コードを作成した時同コードが多くなるのを避ける意味で細かく区分した。

各グループに与えた英字1文字のこ

ードは、漢字コードとして打つ時の指の動きを規定する。よって頻度の高いグループにはキーボード上の打ちやすいキーの英字を割り当てるというような配慮が必要である。しかし頻度を決定するのは入力テキストにおける漢字自体の頻度であり、これはテキストの

内容や種類によっても大きく変動する。本研究では主にコードの覚えやすさに重点をおいたので、ここでは打ちやすさは考慮に入れずただ羅列的にアルファベットを割り当てただけである。しかし入力速度をあげるには、この基本字形のコードをいかに割り当てるかが問題となる。

なお、同コードとなる漢字を区別するために補助的なコードを定めおく。これは基本字形1~2個で構成される漢字に用いられる。それは、前述の四角号碼と同様な部分字形を9種定義し、1~9の数字1字で表現したものである。(表2)

すでに述べた英字1字のコードは基本字形の左上の角あたりの形を示すものである。よってその不備を補うために、右下の角に注目してその字形によって四角号碼的な数字コードを与えてやる。

[表1] 基本字形のグループとそのコード

コード	キ1画	キ2画	キ3画	基本字形の例
A	・	・		ンシツオ斗...
B		一		一言言云云立产产产文亦主衣方广
C		177		リ内州/ノム並並羊必半首火米炎/之良ネ心
D	一	二		平来/一二千开キ弋井夫未天末 三半夫
E		1	一	十土土走产弋赤
F				丁工下王正耳牙井上コ/キ圭/弋寸
G			1	廿廿廿廿廿廿世世世臣
H			/	不才弋本
I			7	行日東東宙五互巨瓦雨再内内可更事車
J		/		厂ナ厂犬尤太犬不戌成百百而夕豕石
K		7		七工工工ラ戸互与万夷事
L	1	7		1小リ/ト非巨上止凸馬片/川/小小
M		7		口几口/丹/月且用目/内内内
N			7	口中虫史串央巨回
O				日里甲申卑門馬果
P			17	巾田内曲曲凹册由皮山出/水内
Q	/	・		ノハハ米人入倉金メ谷父乎米牛
R		一		ノ箭午牛缶矢失朱生年作竹片又夕
S				千生我升糸毛手筆筆重欠牛
T		1		ノ隹介丘丘氏身母白白尔自血鳥向夕禹
U		17		行夕/カ夕乃九丸夕夕夕久久儿几及勿也凡
V	7	・		乙くマ羽 又又又又又 777
W		一		月ヨ尹尹良尸尺弓己巳韦韦民
X		1		了子承尸巴尸也丑正丑尸口
Y		/		刀刃女力
Z		7		母系夕么白西《

[表2] 補助コード

コード	部分字形
1	一し
2	1」ノ
3	、し
4	+又
5	キ+
6	口
7	7」7
8	人ハ
9	小木

6. 漢字のコード化の規則

基本字形1つで構成される漢字(199字)は、原則として2打ちで入力する。1打目は、基本字形のコード(A~Z)であり、2打目に漢字の右下角の部分字形のコード(1~9)である。

例) 文 = B4
言 = B6

以上で同定しきれないものは3打目に記号(□, ;, ., /)を任意に与えた。

例) 亡 = B | □
立 = B | ;
主 = B | ;

コードの1部を表3に示す。

基本字形2つで構成される漢字(878字)は3打ちで入力する。1打目と2打目は、2つの基本字形のコード(A~Z)を書く順序に従って並べる。3打目は原則として□(空白)だが、同コードの漢字がある場合には、漢字の右下角の部分字形のコード(1~9)を与える。

例) 童 = B O □
竜 = B O 1
音 = B O 6
課 = B O 9

ただしそれでも同定できなければ3打目に漢字の右下角でなく、2つの基本字形のうちどちらかの右下角に注目して数字コードを与える。

例) 冗 = M U □

[表3] 基本字形1つで構成される漢字とそのコード(1部分)

		2nd stroke								
		1	2	3	4	5	6	7	8	9
1st stroke	A				斗					
	B	亡, 五, 王,		衣	元		言	方		亦
	C	並, 乙,	羊, 州,	良	必	半	首	心	火	米
	D	一, 二, 三,	死, 干,	-	井, 寿,				夫, 天,	采, 未, 末,
	E	土, 士,		走	十					赤
	F	土, 玉, 五,	丁, 耳, 牙,	下	寸					
	G	世, 臣,						廿, 廿,		
	H		才		本					不
	I	亞, 瓦, 五,	可		兒, 更, 再,	兼, 康,	西, 酉,	雨, 南, 再,		東, 東,
	J	尤	不	戊, 成,	文		百	而	大, 太, 犬,	
	K	七, 互,			与		戸	石	走	
	L	上, 止,	川	氷, 止,			凸	片		小
	M	且			丹		目	内, 月, 用,		
	N		中	虫	史	串	口		央	

(注) 漢字の右下角の記号が3打目にある。記号が書かれていない漢字は3打目は必要でない。

見 = MU 1
肌 = MU 7

基本字形3つ以上で構成される漢字
(1888字)は、漢字から基本字形
3つを選び出し、書く順序に従ってそ
のコードを配列する。よって英字3文

コード表の1部を表4に示す。

[表4] 基本字形2つで構成される漢字とそのコード(1部分)

		2nd stroke								
		A	B	C	D	E	F	G	H	I
1st stroke	A		泣, 洋, 脊, 洲, 洗, 浪, 瑪, 河, 考, 玃, 柱	头, 洗, 浪, 瑪, 河, 考, 玃, 柱	玃, 河, 考, 玃, 柱	玃, 河, 考, 玃, 柱	玃, 河, 考, 玃, 柱	玃, 河, 考, 玃, 柱	玃, 河, 考, 玃, 柱	玃, 河, 考, 玃, 柱
	B	造	訪, 註	心, 忘, 忘, 許	庄, 計, 辛, 封, 序, 訂	庄, 計, 辛, 封, 序, 訂	庄, 計, 辛, 封, 序, 訂	庄, 計, 辛, 封, 序, 訂	庄, 計, 辛, 封, 序, 訂	庄, 計, 辛, 封, 序, 訂
	C	料, 為	粒	祥, 笑, 迷, 春, 道	字, 笑, 社	字, 笑, 社	字, 笑, 社	字, 笑, 社	字, 笑, 社	字, 笑, 社
	D	迂			秦, 秦, 耕	秦, 秦, 耕	秦, 秦, 耕	秦, 秦, 耕	秦, 秦, 耕	秦, 秦, 耕
	E	过	坊, 裁	志, 南	坪	圭	寺	堪	裁	載
	F			恥	扶, 得, 扶	扶, 得, 扶	式, 打	弄		拒, 捕, 珂
	G		芳	芝, 心	芳, 芳	基	葉, 茸	華	某	莖, 奇
	H		柱	株, 述	株, 述	杜, 戎	村, 程, 奎, 柑, 戎	村, 程, 奎, 柑, 戎	村, 程, 奎, 柑, 戎	柄, 捷, 栖
	I			患, 患, 速, 速, 速, 速, 速, 速	軒, 重, 栗	車, 車, 耐	車, 車, 耐	車, 車, 耐	車, 車, 耐	輔
	J			及, 退, 逐, 威, 殘, 辰	压, 在	耐, 左, 至	耐, 左, 至	耐, 左, 至	耐, 左, 至	奇
	K		房		至	匡	匡	匡	匡	匡
	L	馬		業						
	M		肪		肝		肘			同, 單
	N	團	蛇, 跡, 野	思, 患	味, 團, 跨, 吐, 叶	呈, 團, 團	呈, 團, 團	呈, 團, 團	呈, 團, 團	呈, 團, 團
	O	黑		開, 開, 惡, 開, 昧, 旦, 早		是, 開, 旺, 華	是, 開, 旺, 華	是, 開, 旺, 華	是, 開, 旺, 華	是, 開, 旺, 華
	P	边		思, 畔, 缺		听, 畏	听, 畏	听, 畏	听, 畏	听, 畏
	Q		忙	乞, 必, 悉	余, 錢, 鑄, 針	全, 飯, 鉦	全, 飯, 鉦	全, 飯, 鉦	全, 飯, 鉦	全, 飯, 鉦
	R	無	筮	造	筮, 筮, 壯	筮	筮	筮	筮	筮
	S	科		秋, 秋	秤					
	T	焦, 鳥, 鳥	位, 住, 依	伴, 近, 逸, 仁, 缺	仕, 代, 集	代, 付, 射	代, 付, 射	代, 付, 射	代, 付, 射	代, 付, 射
	U	冬	往	冬, 冬, 忽, 行		征, 征	征, 征	征, 征	征, 征	征, 征
	V			脫		弁	桑	司		
	W	尽		忌, 退		既	展			
	X		防	隊						陳, 陣, 阿
	Y		妨	妨, 忍, 忍, 妹			姬			
	Z		紅, 綉	災, 逃		紅	紺			練

注) 漢字の右下の口(空白)の下に数字の 3rd stroke.

字のコードとなる。4つ以上の基本字形を持つ漢字では、どの3つを選ぶかという問題がある。主に次の3つのいずれかでコード化した。

(1) 最初の3つの基本字形をとる。

例)

1	2
3	

 徴 = イ山王
 U P F

(2) 3打目に最後の基本字形をとる。

例)

1	2
3	

 徴 = イ山心
 U P C

(3) 途中をとばす

例)

1	2	3

 戯 = 業ギ
 L L H

以上で同定できないものは3打目に字形とは無関係な、他の漢字で使われていないコードを選んで与えた。このような例外は63字あった。

例) 請 = B F Z

cf. 靖 = B F M

なお3打目には主に X ~ Z を与えておいた。表5に基本字形3つ以上で構成される漢字とそのコードを一部分あげておく。

7. おわりに

本コードによって実際に日本語テキストが入力できるシステムを作成した。使用機器は HITAC M-180 (広島大学情報処理センター) とグラフィック端末テクノロクス 4010 である。

試行の結果、入力速度としては1分間に17~18文字程度であった。これは製作者自身の場合で、コード作成期間約2週間、練習期間約1ヶ月の場合である。まだ漢字を見て反射的にコードが出るという段階までいたったことがないので、訓練しただけではもっと速くなると思われる。なお全くの初心者には試していいないが、英文タイプやカナタイプに熟練した人であればコードの規則さえ覚えればすぐ利用できる。

終わりに、字形データとして、電子技術総合研究所の FONT 4000 を使わせて頂いた。ここに謝意を表す。

<参考文献>

- 1) 田中二郎, 山田尚勇: タッチ打鍵による日本文入力法の研究, 125 pp., 東京大学理学部情報科学科
- 2) TECHNICAL REPORT 78-01
山田尚勇: 日本語テキスト入力法の人間工学的比較, 32 pp., 東京大学理学部情報科学科 TECHNICAL REPORT 78-06
- 3) 木澤誠: 漢字の入力, 電気学会雑誌, 97巻2号, pp.90-92 (昭52-2)
- 4) 渡辺定久: 漢字入力装置, 電気通信学会誌, vol.63, No.7, pp.707-712 (1980年7月)
- 5) 増田功: 日本語文字読み取り装置, 電気通信学会誌, vol.63, No.7, pp.719-723 (1980年7月)
- 6) 楊維楨: コンピュータによる漢字処理をめぐる問題, システムと制御, vol.24, No.6, pp.372-378 (1980)
- 7) 渡辺茂: 漢字と図形, 245 pp., NHK ブックス 264 (昭51)
- 8) 読売新聞社: 日本語の現場 オール集, 210 pp. (昭51)

9) 文部省：外国人のための漢字辞典，大蔵省印刷局（昭41）
* * *

10) Donald E. Knuth：基本算法／情報構造，pp.87-148，米田信夫・笈捷彦共訳，サイエンス社（昭53）

[表5] 基本字形3つ以上の構成される漢字とそのコード（1部分）

A

1st stroke		3rd stroke																										
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	
2nd stroke	A					涌	滂																					泳
	B																											
	C																											
	D																											
	E																											
	F																											
	G																											
	H																											
	I																											
	J																											
	K																											
	L																											
	M																											
	N																											
	O																											
	P																											
	Q																											
	R																											
	S																											
	T																											
	U																											
	V																											
	W																											
	X																											
	Y																											
	Z																											

B

1st stroke		3rd stroke																									
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
2nd stroke	A																										
	B																										
	C																										
	D																										
	E																										
	F																										
	G																										
	H																										
	I																										
	J																										
	K																										
	L																										
	M																										
	N																										
	O																										
	P																										
	Q																										
	R																										
	S																										
	T																										
	U																										
	V																										
	W																										
	X																										
	Y																										
	Z																										