

漢字の入出力処理について

坂本義行 (電子技術総合研究所)

1. はじめに

日本語テキスト(漢字かな混り文)を計算機へ入力する方法については、すでに種々の方式が開発されている。ここで紹介するのは、対象となる利用者として研究者(素人)自身が小論文(数千文字から成る)等のテキストを直接入力する方法についてである。すなわち、大型計算機に接続されているTSS画像端末から会話形式により、簡単な規則で確実に打鍵入力が可能な方法である。

打鍵は、画像端末の英数字鍵盤によりローマ字記述で行なうが、字種の識別を行なうため、大小シフト・キー等数種の識別符号を挿入する。ローマ字かな変換処理では、種々の表記法を可能な限り許すこととし、例外的なものについては、パラメータの変更により表記する方法が用いられている。かな漢字変換処理に関しては、対象となるテキストの特性から、汎用的な辞書として恒久的に利用する固定辞書とテキストあるいは利用者固有の一時的に利用する浮動辞書の2種類を設けている。この辞書の検索方法は、会話形式により、同音異字については、画面に表示し、選択追加処理を行なう。同時に、「語音」の処理が可能である。

なお、当システムで作成されたテキストの編集処理は、独立に開発された漢字テキスト・エディタにより行なうことができる。さらに、テキストの横組編集出力を行なう漢字RUNOFFシステムを開発中である。

プログラムは、処理の簡便さからTSS FORTRANで記述されている。辞書は、Tree構造で表現されており、固定辞書の蓄積データとして、当用漢字音訓表の「字音」とその漢字符号が用いられている。漢字符号は、汎用性の面から整数10進4桁符号が割り当てられ、文字パターンは、大型計算機のディスク上にベクトル方式で、onlineランダム・ファイルとして約3,000種の文字が記憶されている。

2. 方法

2.1 入力の諸方式

漢字のかな混りテキストを計算機へ入力する方法としては、

- (1) 漢字コードを直接打鍵入力
- (2) かな(ローマ字を含む)漢字変換
- (3) 漢字認識(OCR)

の三種に大別される。

(1)の方式は、漢テレ等を用いるフルキー方式と呼ばれるもので、数千種の鍵盤から選択し打鍵するため、かなりの訓練を要する奥で操作性に問題がある。

(2)の方式は、(1)に比較し操作性では優れているが、変換を自動化するには、ぼう大な辞書と複雑な手続きを必要とし、現在までのところ完全な変換結果は得られていない。

(3)の方式は、試験的な段階であり、英数字、かな、一部の漢字や特殊字形については、処理システムが開発されているが、実用という面からは将来の向

題であろう。

以上いつれ的方式にも特徴を有しており、使用面から有用性が考えられる。

テキストの性質と容量によって、つぎの2種類について上記の(1),(2)を検討すると、

- (a) 大量のデータを熟練者に委託する。
- (b) 少量のデータを研究者(素人)が入力する。

(a)に対しては、(1)の方法を用いると、offlineで紙テープ等の媒体上に漢字符号を直接入力できる。この方式は、熟練者では鍵盤探索操作が安定し、打鍵速度が平均50字/分程度となる。漢字への変化処理による誤差を発生しない長、大量のデータの一括入力処理に適している。打鍵誤りの発生率は、30万字の実験で、0.01% (4/41件)、その内容は、表1表のような傾向を示し、二回の修正処理により完全なデータが得られた。

(b)の方法では、素人という長、操作が容易で訓練を要しない長がもっとも大きな問題となる。すなわち、(2)の方式で同音異字の処理を人間の判断に委ね、正しい漢字を選択する。また視認性をよくするたり、画像端末上に漢字パターンを表示する方法をとる。

第1表 打鍵誤りの内容

誤りの内訳	件数	比率(%)
欠落文字	90	21.7
" 機能記号	49	11.8
過剰打鍵	41	9.9
原文不鮮明によるもの	63	15.2
隣接キーを誤り打鍵	76	18.3
類似語	13	3.1
同音異義語	5	1.2
外字機能記号誤り打鍵	5	1.2
文字種誤り(ひらがなとカタカナ)	4	0.9
" (漢字とかな)	3	0.7
その他	65	15.7
計	414	

2.2 ローマ字表記とカナ変換

汎用性を目的としているため、特殊なカナ鍵盤を用いず通常の英数字鍵盤を有する画像端末からの入力を考えた。これは打鍵速度の長で、文字数の増加分は、鍵盤の操作性、訓練の必要性という面からカナによる入力と差がなれどおもしろい。

ローマ字表記には、ハボン、訓令、日本といった諸式があるが、外来語、外国の人名、地名に関しては、表記が一定していない。こゝでは可能な限り種々の表記が行なえる「ゆるやか」な規則とし、表記が重なるものについては、スイッチを設け、パラメータによる選択方法をとっている。漢字、カタカナ、ひらがなによる表記の異なりを表2表のように5個の仮引数 BOIN, SHIIN, KANA, JP1, JP2

表2表 選択表記規則

		BOIN	SHIIN	KANA	JP1	JP2	WI	WF	WO	TU	DU
ひらがな	a ₁	CBOIN	CSHIIN	0	0	0	お	ゑ	え	つ	つ"
	a ₂	"	"	0	1	0	ゐ	ゑ	を	と	と"
	a ₃	"	"	0	1	1	ゐ	ゑ	ゐ	と	と"
カタカナ	b ₁	SBOIN	SSHIIN	100	1	1	ウ	ウ	ウ	ト	ト"
	b ₂	"	"	100	1	0	ウ	ウ	ヲ	ト	ト"
漢字	c	"	"	100	0	0	井	工	ヲ	ツ	ツ"

で選択する。なお、表3表に漢字に付いた規則表を示した。その結果、表記に関する自由度が減少した。

表3表 漢字字音のローマ字表記

ア	イ	ウ	エ	オ	カ	キ	ク	ケ	コ	ガ	ギ	グ	ゲ	ゴ
a	i	u	e	o	ka	ki	ku	ke	ko	ga	gi	gu	ge	go
カ	キ	ク	ケ	コ	ギヤ	ギイ	ギユ	ギエ	ギョ	ザ	ジ	ズ	ゼ	ゾ
ka	ki	ku	ke	ko	gya	gyi	gyu	gye	gyo	Za	{zi ji}	zu	ze	zo
サ	シ	ス	セ	ソ	シャ	シイ	シュ	シェ	ショ	ダ	ディ	デュ	デ	ド
sa	{si shi}	su	se	so	{sha cya}	syi	{shu syu}	{she sye}	{sho syo}	da	di	du	de	do
タ	チ	ツ	テ	ト	チャ	ジイ	ジュ	ジエ	ジョ	バ	ビ	ブ	ベ	ボ
ta	{ti chi}	{tu tsu}	te	to	{cha ja}	zyi	{zyu ju}	{zye je}	{zho jo}	ba	bi	bu	be	bo
ナ	ニ	ヌ	ネ	ノ	チャ	チイ	チュ	チェ	チョ	パ	ピ	プ	ペ	ポ
na	ni	nu	ne	no	{cha cya}	cyi	{chu cyu}	{che cye}	{cho cyo}	pa	pi	pu	pe	po
ハ	ヒ	フ	ヘ	ホ	ヒヤ	ヒイ	ヒユ	ヒエ	ヒョ	ファ	フィ		フェ	フォ
ha	hi	{hu fu}	he	ho	hya	hyi	hyu	hye	hyo	fa	fi		fe	fo
マ	ミ	ム	メ	モ	ビヤ	ビイ	ビユ	ビエ	ビョ	ヴァ	ヴィ	ヴ	ヴェ	ヴォ
ma	mi	mu	me	mo	bya	byi	byu	bye	byo	va	vi	vu	ve	vo
ヤ		ユ	イエ	ヨ	ピヤ	ピイ	ピユ	ピエ	ピョ	ツァ	ツイ		ツェ	ツォ
ya		yu	ye	yo	pya	pyi	pyu	pye	pyo	tza	tzi		tse	tso
ラ	リ	ル	レ	ロ	リヤ	リイ	リュ	リエ	リョ	グァ	グイ	グウ	グェ	グォ
ra	ri	ru	re	ro	rya	ryi	ryu	rye	ryo	gua	gui	guu	gue	gao
ワ	ヰ		ヱ	ヰ	ヂヤ	ヂイ	ヂユ	ヂエ	ヂョ	グァ	グイ	グウ	グェ	グォ
wa	wi		we	wo	dya	dyi	dYu	dye	dyo	gwa	gwi	gWU	gwe	gwo
ン					ヂヤ	ヂイ	ヂユ	ヂエ	ヂョ	クワ				
n(')					tya	tyi	tyu	tye	tyo	qua				
ン					ニヤ	ニイ	ニユ	ニエ	ニョ					
m(')					nya	nyi	nyu	nye	nyo					
					ミヤ	ミイ	ミユ	ミエ	ミョ					
					mya	myi	myu	mye	myo					

{ } ; ' ' の表記の可

字種の識別を行なった後に、各々の字種による変換処理を行なうために、操作性を考慮して、数種の簡単な記号とシフトキーを定めていた。

(1) 漢字 - 鍵盤上に英文字の大小によるシフトキーが備えられているので、これを利用して漢字部分に付いて小文字で表現する。

(2) ひらがな - シフトキーにより大文字で表現する。

(3) 英数字, 特殊記号 - 「?」記号で両端を囲んだ文字列で表現する。ただし, 「?」自身は, !, ?, @, 英数字を除く記号で「?」を囲む。

(4) カタカナ - 「!」記号で両端を囲んだ文字列で表現する。

2.4 固定辞書と浮動辞書

かなから漢字への変換処理は, 辞書による自動変換と人用を介して変換する同音異字処理の部分からなっている。ここでは確実に変換することを主たる目的としており, 処理速度は考慮されていない。

かな漢字変換のための辞書として,

(i) 固定辞書

(ii) 浮動辞書

の二種類が設けられている。

(i) の辞書とは, 汎用性のある基本語について作成されている辞書で, 半永久的に検索のみに使用されるものを指す。一字訓読み辞書, 当用漢字音訓辞書等が構成されている。一字訓については, 野村氏⁽¹⁾の調査によれば, 第4表のように一字訓が訓読みの語にあって大多数を占める。また同音異字が音読み結合形に比較し, 極端に少ない値を示している。異なり語音⁽²⁾項目について, その語音と漢字符号を辞書として登録した。なお, 同音異字が最大5個のものが存在した。当用漢字音訓辞書は, 約3,500項目が登録されている。

(ii) 浮動辞書は, 本システムの特徴をなす部分で, 利用者(研究者)自身で作成し, 検索を行なう辞書を指す。

利用者が処理対象とするテキストの大きさは, 高く数千から数万文字から成るものと考えられる。そのテキスト内で出現する異なり語(漢字列)は, 実験から第5表に示すような値が得られている。すなわち, 二文字以上で構成されている語(漢字)での同音異義語は, ほとんど出現しない。一字の場合も訓読みとして出現し, 同音異義語として現れない。この裏から, テキスト内で始めて出現した語(漢字列)を辞書に一回登録することにより, 以下の処理では, 自動的に検索が行なわれ, 漢字への変換が行なわれる。

2.5 漢字テキストの編集処理

この処理は, 本システムとは独立に, EPICS (TOSBAC 5600 大型計算機システム) の TSS 環境下において, 漢字テキストを含むテキスト・ストリングを編集する目的で, 東京に外注して作成したものである。ここでは, 本システムと結合して使用するため, その特色を簡単に述べるにとどめ, 詳細は, その解説書にゆずる。

2.5.1 漢字符号と入力表現

漢字符号の割付けは, JIS 標準化の作業が現在行なわれている段階で, 尋

第4表 漢字音訓の用法分析 (数字は延べ度数, ()内の数字は%)

用法	音	訓	計
自立	12,280(1.9)	80,489(53.9)	92,769(11.0)
結合	514,822(81.9)	54,089(36.2)	568,911(73.1)
接辞的(前部分)	13,926(2.2)	2,293(1.5)	16,219(2.1)
接辞的(後部分)	87,829(14.0)	12,511(8.4)	100,340(12.9)
計	628,857(100.0)	149,382(100.0)	778,239(100.0)

(言語生活 1975.6 より)

第5表 テキスト内での出現漢字語の分析

データ	延べ語数	異なり語数	異なり漢字語数	λ (%)	μ (%)
1	1,241	525	201	42.3	16.2
2	1,662	721	281	43.4	16.9
3	2,618	1,100	574	42.0	21.9
4	2,637	885	415	33.5	15.7
5	3,208	1,208	607	37.6	18.9
6	3,425	1,116	513	32.5	14.9

λ =(異なり語数/延べ語数) $\times 100$

μ =(異なり漢字語数/延べ語数) $\times 100$

この問題を含んでいるわけが、一意に決定することは困難であるが、プログラム上でのデータ表現、一定の順序配列、約一万種の文字表現等の長から、10進法桁で表現した。ただし、JIS が決定されれば、この符号との変換を可能なものとする。

入力表現は、英数字のみからなるテキストで用いられている α -mode のデータを一定の条件のもとで、混在すること許すため、漢字データの区切りを示す、開始記号「'''」と終了記号「"""」で囲まれた数字列で表現する。

2.5.2 漢字データの出力

漢字データの出力形式には、つぎの4種が備えられている。

(i) α -mode による出力では、入力表現と同じもので、開始、終了記号を含む数字列からなる。

例 '''16821065054200700534'''

(ii) PRINT 出力では、(i) の場合の開始、終了記号が除かれ、漢字データの文字区分を明示するために、4桁の先頭の数字に「+」を重お書きする。

例 +682+065+542+070+534

(iii) DISPLAY 出力では、漢字パターンで表示される。

例 漢 字 デ ー タ

(iv) EXHIBIT 出力では、漢字符号での表示および漢字パターンで表示され、第1行を漢字符号列 (PRINT 表示)、第2行を漢字パターン (DISPLAY 表示) を並記する表示方法である。

例 +682+065+542+070+534
漢 字 デ ー タ

なお、1行は最大18個(漢字)表示可能であり、1画面当りの収納行数は、数字列で35行、漢字パターンで17行収納できる。

2.5.3 編集の機能

編集処理を行なうための命令は、機能を示す動詞の部分と、機能を修飾するオペランドとの組合せから成り、その形態は、表6に示す8種類であり、その動詞の種類は、表7に示すものが可能である。

表6 編集命令の形態

①	VERB
②	VERB m ; r
③	VERB m : st
④	VERB m : st ; r
⑤	VERB m : st , st
⑥	VERB m : st : st
⑦	VERB m : st , st ; r
⑧	VERB m : st ; r : st

注

m : s, L, null
(モード表示)

r : n, *(繰返回数)

st
(ストリングフィルタ)

表7 編集制御の種類

BACKUP or B	COPY	RUNOFF	NONVERIFY
FIND or F	CUT	DISPLAY	ASCERTAIN
PRINT or P	PASTE	CASE	NOASCERTAIN
DELETE or D	BUILD	STANDARD	CONFIRM
INSERT or I	LINE	VERIFY	NONCONFIRM
REPLACE or R	STRING	KANJI	NORMAL
EXHIBIT			

2.6 漢字テキストの表示処理

漢字テキストを一定の様式に従って出力する処理は、前節で述べた漢字テキスト・エディタ同様、EPICと結合してゆく漢字RUNOFFサブシステムが、現在開発中である(76年12月、稼働予定)。その概要を紹介すると、改行、改頁、負付け等のテキストの様式を指定するための制御語が多数設けられている。これを表8に示す。

表8 様式指定の制御語

①	• CHSIZE	n	②①	• PAGE	x, y, n
②	• CHPITCH	n	②②	• PAPERLENGTH	n
③	• LNPITCH	n	②③	• PARAGRAPH	
④	• BEGINPAGE	n	②④	• POINT	n
⑤	• BOTDOMMARGIN	n	②⑤	• REFERENCE	n
⑥	• BREAK		②⑥	• SCOREUNDER	n
⑦	• CENTER	n	②⑦	• SINGLESPEACE	
⑧	• COMMENT		②⑧	• SPACE	n
⑨	• DOUBLESPACE		②⑨	• SUBHEADER	x, n
⑩	• FILL		②⑩	• SUBFOOTING	x, n
⑪	• FOOTING	x, n	②⑪	• SUBPARAGRAPH	n
⑫	• HEADER	x, n	②⑫	• TOPMARGIN	n
⑬	• INDENT	n	②⑬	• UNIDENT	n
⑭	• LEFTIDENT	n	②⑭	• TABULATE	n, ..., n
⑮	• LINELENGTH	n	②⑮	• NOTAB	
⑯	• MARGIN	t, b, l, r	②⑯	• BOLDFACE	n
⑰	• MULTISPEACE	n	②⑰	• HALF	
⑱	• NODENT		②⑱	• FULL	
⑲	• NOFILL		②⑳	• JUSTIFY	

これ等の制御語をテキスト中に挿入するのは、2.5で述べた漢字テキスト・エディタを用いて行なうことが出来る。また、これ等の制御語を9,000番台の数字によって表現することも出来る。各制御語の説明は省略する。

この様式化されたテキストをファイルに蓄積または印字出力するために、REFORM命令があり、73の3種類の形式がある。

- i) REFORM ファイル名1, ファイル名2, PRINT
- ii) REFORM ファイル名1, ファイル名2
- iii) REFORM ファイル名1, PRINT / DISPLAY

制御語の中で、とくに JUSTIFY は禁則処理を行なう命令で、横組における、左端の2文字、左端の1文字に関する禁則規則を、フルピッチ、ハーフピッチの文字の組合せによる処理を行なうことが可能である。

漢字のフォントパターンは、その符号配列として、特殊記号、数字、英字、ギリシア文字、ロシア文字、ひらがな、カタカナの順に 0 ~ 599 の範囲に割付け、601 から上位に、当用漢字、人名漢字、補正漢字を部首別で割付けられている。漢字パターンは、文字を直線に分解し(最大62)、端末へ伝送するベクトル方式がとられており、ドットによる直接伝送方式より、伝送速度の向上がはかられている。

3. 実験と結果

3.1 辞書の蓄積と検索

固定辞書となる一字訓読み、当用漢字音訓のデータは、その字音をローマ字表記し、漢文字符もかっこに囲む、つぎのような表現で入力される。

例: kan (1177, 787, 1776, 1608, ..., 1781)

なお、同音異字は「,」で区切り連続して入力される。

辞書の内部では、ローマ字からカタカナに変換され、文字を単位節とよぶ Tree 構造で、漢文字符は Table の形で蓄積されている。このため、辞書の修正は容易である。

浮動辞書への未登録語の蓄積は、固定辞書を検索するために、字音または字音と「+」記号の形で入力する。

例: 一字訓 aida / (CR)
音訓表 kan+ / (CR) (CR) 改行キー

同音異字の選択は、画面の下部に表示されたものの中から送んで、その番号を打鍵する。蓄積は、語を単位として処理されたため、「+」記号を加える。

例 kan+ / (CR) . ji+= / (CR) ; 漢字

この実例を才士団に示した。選択された語は、テキストに収められるとともに、浮動辞書に登録される。

複合語の表現は、

例 漢字 漢字入力 漢字入力処理

に対して、

(((漢字) 入力) 処理)

の形で3語として登録される。この3語の検索は、

例; kanji kanji, nyuuryoku kanji, nyuuryoku, shori

の形で行われる。
つぎに、語音補正処理は、

例 michi + /CR , hon + /CR got = nihongo / CR ; 日本語

の形で打鍵する。
検索全体の手順を下記図に示す。

(電子技術総合研究所)

1. はじめに 日本語テキスト

日本語テキスト(漢字かな混り文)を計算機へ入力する方法については、すでに種々の方式が開発されている。ここで紹

01 02 03 04 05 06 07 08 09 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30

-01 介灰回会快戒改怪海悔皆界械開階絵塊解懐壊街貝

日本
=go+=nihongo

01 02 03 04 05 06 07 08 09 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30

-10 五互午吳後悟娛御基語誤護期

図1 画面に於る選択表示の例

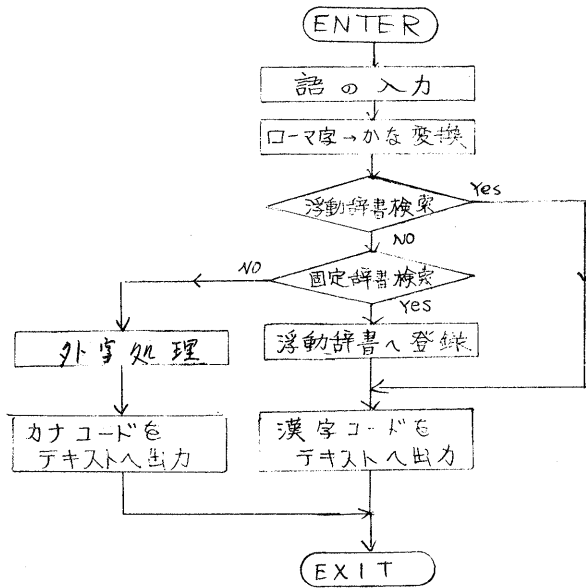


図2

検索の手順

3.2 テキストの蓄積と編集

1個のテキストを蓄積する単位をファイルと呼ぶ、1ファイルは複数のレコードで構成され、1レコードは、最大文字数を指定することにより構成される。蓄積時におけるレコード単位の修正は、バッファに一時的に記憶されており、次のレコードが入力された時点で、ファイルに蓄積される。ファイル単位の修正は漢字テキストエディタにより処理を行なう。この出力例を図3に示す。

テキストの様式化は、漢字RUNOFFサブシステムを用いて、画像端末あるいは、他の出力装置により、印字出力を行なう。

3.3. プログラムと結果

主プログラムは、TSS FORTRAN で約1,500行から成っており、漢字パターンの表示は、ユーティリティ・プログラム PLOT-10 を主プログラムから呼び出しにより行なう。一時刻、音割、浮動辞書、テキスト・ファイルは、ディスク上の online-file に蓄積され、検索処理開始前にコアに読み込まれる、その容量は約4KB (1K=36bit) である。

固定辞書の検索速度は、EPICS の TSS 利用者の利用度によく影響される、とくに、同音異字が多い場合、この文字の漢字パターンを端末へ伝送、表示するのに、数十秒から数分を要する。さらに、画像端末として、蓄積型が用いられているため、画面の一部消去が不可能であり、固定辞書を検索するたびに画面全部を書きかえを行なう長で多くの時間を要している。

4. あとがき

今後の問題として、固定辞書の拡張、パターン伝送の高速化、画像の部分制御、種々の様式が可能にハードコピー等の改良とともに、新しい方式として、長谷川大⁽³⁾の英字におけるかな漢字変換方式の並置を検討中である。

応用として、既存の大量データの編集、活用処理⁽⁴⁾、Concordance⁽⁵⁾の諸プロ

て □ 漢 字 の 入 出 力 処 理 に つ い
 術 総 合 研 究 所) □ □ 坂 本 義 行 は じ め に 電 子 技
 語 テ キ ス ト □ □ 日 本 語 テ キ ス ト (漢 字 日 本 本
 な 混 り 文) を 計 算 機 へ の 入 力 方 式 が 開 発 さ れ っ
 い て は 、 と じ 紹 介 す る の は 、 対 象 と な っ て
 □ 利 用 者 と し て □ □ □ □ □ □ □ □ □ □ □ □ □ □

end of file

才3回 漢字テキストエディタの出力 (DISPLAY)

プログラと結合し、日本語テキスト分析を完了する。

最後に、本プログラムの作成を助けて下さった、電気通信大学の佐藤雅之君、漢字パターンを提供下さった人間機械研究室、共通ソフト開発で検討会に参加して下さった諸君方と武東楚の担当者の方々に感謝いたします。

参考文献

- 1) 野村雅昭; 新聞の文章につかわれた漢字, P35, 言語生活, 1975.6, 筑摩書房.
- 2) 林四郎; 漢字の役割, P25~26, 言語生活, 1975.6, 筑摩書房.
- 3) 長谷川貞夫; 英文字号による漢字を含む普通文字の印刷, 日本特殊教育学会才13回大会発表論文集, 1975年9月.
- 4) 坂本義行; 日本語の活用処理, CL5-2, 計算言語学研究会資料5, 1976.5, 情報処理学会.
- 5) 坂本義行; 岡本哲也; 日本語のコンコ-ダンス, CL2-1, 計算言語学研究会資料2, 1975.7, 情報処理学会.