

## かなべた文の単語分割アルゴリズムの一方式

沼田泰之、藤田克彦、林 大川、山内佐敏  
(（株）リコー・技術本部)

### 1はじめに

日本語の入力作業において、最も大きな問題の一つが多数の漢字の存在であったことは言うまでもない。特に、専任のオペレータではない一般の人間にとっては、数千にもおよぶ文字の中から希望の文字を探して入力することは、事実上不可能に近いことであった。

かな漢字変換技術はそのような背景のもとに産み出され、これにより、一般の人間にも日本語の入力が従来よりはるかに容易に行なえるようになったのである。

しかし、現在の日本語ワード・プロセッサで主流をなしている方式は、単語ないし文節といった文法上の単位、概念を意識しながら入力作業を進めなければならないという問題点を残している。何が理想的な日本語入力法であるか、の議論は別にしても、より人手の介入の少なくてすむシステムのほうが望ましいという点については異論はないであろう。すでに、そのような観点から、べた書きされたかな文字列に対するかな漢字変換方式<sup>1~3</sup>の提案もいくつかなされている。

こうした中で、我々は入力に追随して逐次変換結果の出力ができるべた文のかな漢字変換方式を開発した。この方式は、日本語文章中において漢字表記される単語の多くがいわゆる漢語であるこ

とに着目し、それら漢語の読みであるかな文字列の性質を利用した、高速かつ精度の良いかな漢字変換方式である。本稿では、特にその中の単語区切りの方式について述べる。

### 2日本語文章の特徴

図1に、日本語の単語分類を示す。

単語	自立語	外来語	漢語		
			カタカナ語		
和語（やまとことば）			和語（やまとことば）		

図1 日本語の単語の分類

まず、単語は自立語と付属語とに分類できる。ここで、付属語は一般にひらがな表記される。自立語は、さらに和語と外来語に分類できる。和語には、動詞、形容詞などが多いが、これらは必ずしも漢字表記される必要はない。一方、外来語には漢語とカタカナ語とがあるが、かな漢字変換は、かな表記されたものを、適切に漢字に変換するこ

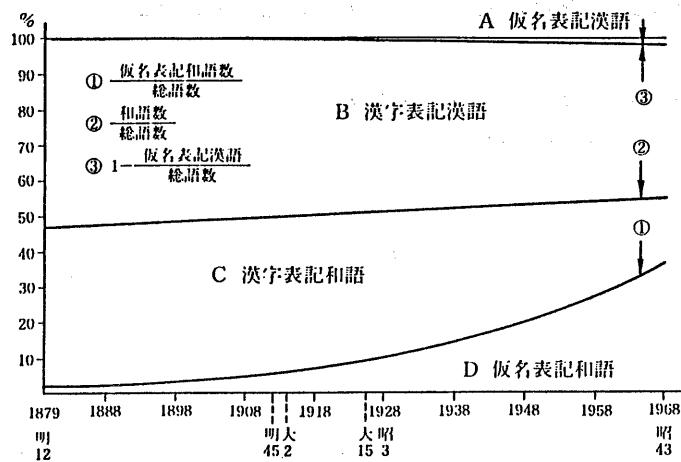


図2 語種別の語表記率の変遷<sup>4</sup>

とが目的であるから、カタカナ語はカタカナで入力することを原則とすれば、処理の対象としては漢語の方に注目すればよい。

次に、図2が示すのは、語種別にみた漢字表記率の変遷である。全体として漢字による表記の割合は低下している傾向は明らかだが、漢字で表記されている単語に占める漢語の比率は、高くなっていることがわかる。

筆者らの事務文に対する調査においても、漢字表記された単語の8割弱が漢語であるという結果が得られた。このことは、漢字で表記されるものの多くが漢語であることを示している。

以上の事から、かな漢字変換の対象として、漢語が極めて重要であることがわかる。

### 3 漢語を特徴づける漢字音

かな漢字変換において重要な対象であるこの漢語を特徴づけるものは、何か。これは、明らかに漢字音（音訓でいう音）である。

この漢字音には、かなで表記すると1字のもの、2字のもの、3字のものがある。

<sup>5</sup>新明解国語辞典には、造語成分の読みとして漢

字音が272種が記載されているが、そのうち2字以上のものが224種で、1字のものが48種ある。さらに2字以上のものについて、促音化したり半濁音化したもの考慮すると、現実には上記の他に113種が漢字音として存在しうることになる。

つまり2字、3字の漢字音が337種、1字の漢字音が48種ということである。

漢語の読みは、これら漢字音によつて構成されているわけである。

ここで、以下の説明の便宜上、改めて漢字音を定義しておくことにする。

漢字音——上記漢字音のうち、かなで2字以上である337種。

そして、それ以外のもの、すなわち上記337種の漢字音の構成要素になつてないものは、一般のかなとみなすことにする。このように定義すると

例 開始 (かいし)

においては、「かい」が漢字音、「し」がかなとすることになる。

表1に、上の定義による漢字音の例を示す。

じゅう	じゅく	じゅつ	じゅっ	じゅん
じょう	じょく	じょっ		
すい	すう	すん		
ずい	ずう	ズン		
せい	せき	せつ	せつ	せん
ぜい	ぜき	ぜつ	ぜつ	せん
そう	そく	そつ	そつ	そん

表1 漢字音の例（上記定義による）

#### 4 漢字音による単語の分類

以下の議論のため、漢字音とかなとの組合せによる単語の分類を行なう。

以下に示したのが、単語の形態的分類である。

##### PATTERN 1 漢字音+漢字音

例 開校 (カイ+コウ)、設定 (セッ+ティ)

##### PATTERN 2 漢字音+かな

例 要素 (ヨウ+そ)、新規 (シン+き)

##### PATTERN 3 かな+漢字音

例 期間 (き+カン)、事実 (じ+ジツ)

##### PATTERN 4 かな+かな (+……)

例 緑 (み+ビ+リ)、川 (か+わ)

##### PATTERN 5 かな

例 胃 (い)、絵 (え)

##### PATTERN 6 漢字音

例 法 (ホウ)、感 (カン)

##### PATTERN 7 その他

例 私 (わ+タク+し)

考え (カン+が+え)

注 カタカナは漢字音を示す。

この分類にしたがって、筆者らが作成した辞書中の単語の比率を算出したものが表1である。

分類	語数	比率 (%)
PATTERN 1	1 2 3 2 4	32.1
PATTERN 2	3 3 5 7	8.8
PATTERN 3	3 3 4 0	8.7
PATTERN 4	1 4 3 9 6	37.5
PATTERN 5	6 6 8	1.7
PATTERN 6	1 3 0 5	3.4
PATTERN 7	3 1 1 5	8.1
合計	3 8 5 0 5	100.0

表1 単語の分類別比率

#### 5 漢字音の利用

漢字音に関しては、次に示すような経験的事実がある。

- 1 漢字音は漢語の構成要素として、多く用いられる。
- 2 漢字音が2個連続して漢語を構成することが多い。
- 3 漢字音と直後のかなで漢語を構成することが多い。
- 4 かなと直後の漢字音とで漢語を構成することが多い。

これらの漢字音の特徴を利用して、あらかじめ、辞書検索の前段階において、単語の切目を推測することを考える。

そのためには、入力される区切りの対象となるかな文字列に対しても、漢字音による特徴づけを行なう必要がある。

##### 区切りの対象例文

じしょのけんさくにあたり

のような入力について、漢字音を単位として区切りを行なってみると次のようになる。

##### 処理結果

じ／ショ／の／ケン／サク／に／あ／た／り

ここで、カタカナで表記されている部分が上述の漢字音である。これによれば、漢語相当部分が推定できると考えられる。

そこで、区切り対象となる入力文字列についても、その先頭からのかな文字の並びによって、次のように分類を行なっておくことにする。

##### [1] TYPE 1 漢字音+漢字音

### 例 ケン／シュツ

[2] TYPE 2 漢字音+かな

例 カン／じ

[3] TYPE 3 かな+漢字音

例 か／ノウ

[4] TYPE 4 TYPE 1, 2, 3以外

例 こ／の

注 カタカナは、漢字音を示す。

今後は、区切り対象文字列の属性を、上記TYPE 1, TYPE 2, TYPE 3, TYPE 4と略称する。

これら文字列は、単語との関連でいえば、それぞれ次のような特徴を有する。

#### TYPE 1

PATTERN 1かPATTERN 6の単語が存在している、つまり抽出される可能性が高い。

#### TYPE 2

PATTERN 2かPATTERN 6の単語が存在している可能性が高い。

#### TYPE 3

PATTERN 3かPATTERN 5の単語が存在している可能性が高い。

#### TYPE 4

PATTERN 4かPATTERN 5の単語が存在している可能性が高い。

このことを用いると、区切り対象文字列に対する辞書検索を限定することができる。

## 6 付属語と漢字音

表3は、自立語と、それに後続する付属語の接続の頻度を、事務文227例を対象として調査し

NO	自立語 の品詞		出現頻度	%累積 % 出現率	
				出現率	出現率
1	名詞	の	1217	16.7	16.7
2	名詞	に	891	12.2	28.9
3	名詞	を	628	8.6	37.5
4	サ名	の	623	8.6	46.1
5	サ名	を	420	5.8	51.9
6	サ名	に	409	5.6	57.5
:	:	:	:	:	:
14	名詞	から	109	1.5	80.2
15	サ名	が	105	1.4	81.6
:	:	:	:	:	:
23	サ名	する	55	0.7	90.3
24	サ名	から	52	0.7	91.0
:	:	:	:	:	:
107	サ名	ずつ	1	0.0	100

表3 自立語と付属語との接続頻度

た結果である。

自立語と付属語の組合せ107種中、上位14位までで80%を占め、23位で90%を越えている。

この結果より、次のことが明らかになる。

名詞の直後には「の」「に」などの助詞が現れやすい、また、サ変名詞の直後にはサ変動詞の、「さ」「し」が現れやすいことがわかる。

つまり、直前に抽出され、認識された単語の品詞が名詞で、かつ、それに続く文字列の先頭が「の」「に」などであれば、それは助詞である可能性が相当高く、同様にサ変名詞直後の「さ」「し」などはサ変動詞である可能性が高いといえる。加えて、経験的事実より、名詞の直後に用言が接続するケースは稀である。

しかしながら、複合語の存在を考慮すると、

「の」や「し」については、複合語の後半部を構成する名詞（あるいは接辞）の先頭文字である可能性と、付属語の可能性とを並列に検討する必要がある。

ここで、判断の材料として漢字音を利用することができる。

たとえば、直前に名詞が抽出された状態で、そこでの文字列の先頭文字が上記の付属語相当文字の類であったとする。漢字音による前処理を行なった結果、その文字列がTYPE 1あるいはTYPE 2であれば、その先頭文字は漢語の一部である可能性が高く、複合語のケースと考えられるようになる。

一方、TYPE 4 の場合は、少なくとも漢語の一部ではないといえる。また、TYPE 3 の場合は、可能性としては、付属語と漢語の先頭文字としての両方が考えられる。ここでは、複合語の出現より付属語の出現のほうが頻度的に大きいと考えて、付属語と考えることにした。

以下に、この考え方に基づいた単語区切りアルゴリズムを示す。

## 7 単語区切りアルゴリズム

### 7.1 アルゴリズムの概要

図3は、単語区切り対象文字列から単語区切りを1つ見出すまでのアルゴリズムをブロック図として示したものである。

単語区切りのプロセスは、大別して次の2つとなる。

- (a) 一般単語の読みの抽出（被検索文字列の設定－辞書検索－最適候補選択）
- (b) 特徴付属語の抽出

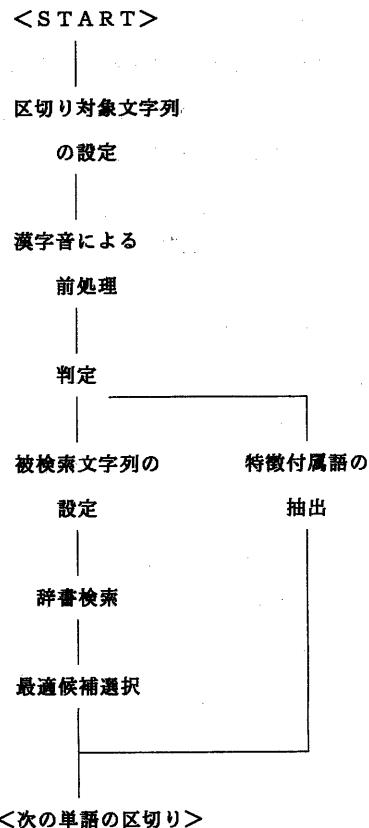


図3 アルゴリズム概要

ただし、(b) から抽出される単語は、後述する一部の付属語に限られる。また、ここでは、通常の辞書検索は行なわれない。

図中の【判定】処理において、次の3項目がチェックされて、すべての条件を満たすときのみ、(b) のプロセスが適用される。

#### チェック項目

- 1) 漢字音による前処理の結果、そこでの区切り対象文字列の属性がTYPE 3, TYPE 4 である。
- 2) 直前に抽出された単語の読みの品詞が名

詞がサ変名詞、形容動詞語幹のいずれかである。

3) 区切り対象文字列の先頭文字が、設定した付属語（サ変動詞語尾、形容動詞語尾等を含む）と一致する。

以下に、個々の処理を説明する。

## 7. 2 一般単語の読みの抽出

### 7. 2. 1 被検索文字列の設定と辞書検索

かなでべた書きされた区切り対象文字列に対し、漢字音による区切りを実行する。その先頭からの漢字音とかなの並び（TYPE）により、辞書検索の対象となる被検索文字列を限定する。

これにより辞書へのアクセス回数を減少させることができることができる。

### 7. 2. 2 最適候補の選択

上記の被検索文字列に対する辞書検索において、複数の単語の読みが得られる場合がある。解析を続行するためには、単語の読みの評価を行ない、一つに候補を絞らなければならない。

従来、評価法としては最長一致法（読みの長さを基準とする）、最尤評価法（読みの長さと使用頻度を組み合わせて評価する）などが知られている。

本方式では、品詞間の接続の重みを重視し、従来よりもやや細かい品詞分類を採用した。評価は、この品詞間の接続の重みと読み単位の評価値との和の形となっている。

なお、辞書検索で候補となる単語の読みが見出されない場合は、辞書検索の限定を解除して、検索を実行し、それでも候補が見出されない場合は、直前の解析に戻る、いわゆるバック・トラックを採用している。もちろん無制限なバック・トラッ

クを防ぐために評価値の累積による制御を行なっている。

## 7. 3 特徴付属語の抽出

漢字音をさらに有効に利用するために、上述のように一部の付属語の抽出において、区切り対象文字列のTYPEによる処理を行なっている。

## 8 単語区切りの具体例

以下に、具体例に即して処理の過程を説明する。

### 入力文字列の例

- (1) さぎょうのうりつがいい
- (2) さぎょうのうきがきまる
- (3) さぎょうのこうかがひくい
- (4) さぎょうのきばがおおきい

いずれの文例においても、文頭から「さぎょう」という単語の読みが抽出されたとすると、残る文字列は次のようになる。

- (1) のうりつがいい
- (2) のうきがきまる
- (3) のこうかがひくい
- (4) のきばがおおきい

「さぎょう」という単語の読みは名詞またはサ変名詞であり、その直後に位置する「の」は助詞の「の」である可能性が高いことになる。しかし、その反面、名詞連続、すなわち複合語を形成する名詞の一部である可能性もある。従来の漢字音を用いないアルゴリズムでは、辞書検索を実行するまでは、いずれの「の」であるか、まったく判断の方法がなかった。しかし、漢字音が漢語を構成する成分であるという特徴を利用しTYPEによ

る処理の切替を行なうことにより、辞書を検索する前にかなり妥当な判別ができる。その様子を次に示す。

(1) ノウ／リツ／が／い／い TYPE 1

(2) ノウ／き／が／き／ま／る TYPE 2

(3) の／コウ／か／が／ひ／く／い

TYPE 3

(4) の／き／ぼ／が／お／お／き／い

TYPE 4

これらに対する妥当な区切りは次のようになる。

(1) のうりつ（名詞）／が

(2) のうき（名詞）／が

(3) の／こうか（名詞またはサ名）／が

(4) の／きぼ（名詞）／が

つまり、これらは漢字音による前処理の結果、漢字音の一部であれば、漢語の先頭文字であり、逆に漢字音の一部でなければ付属語であるという考え方たに即した例となっている。

この処理の特徴は、被検索文字列の設定、辞書検索、最適候補選択といった一連の処理を不要としながらも一部の頻出する付属語を的確に抽出することができる点にある。その結果、後続の文字列の処理へ即座にと移ることができる。

また、この方式は最長一致法などで問題となる「読みの長い単語が常に優先される」傾向への解決策の一つともなっている。

## 7 本方式の性能

市販の文例集にある事務文21文例を対象に、汎用機上で行なった本方式による単語区切りの実験結果では、単語数5167語に対し区切りを誤

った単語数は202で、単語区切り率は96.1%であった。

ここで、単語区切り率とは、同音異表記語への誤りは無視した、単語区切りの成功率のことであり、次に示す式で定義している。この数値は、同音語の学習が可能な単語の率もある。

$$\text{単語区切り率} = (\text{正しく区切られた単語の数} / \text{総単語数}) * 100$$

次に、上記の漢字音処理を行なわない場合と行った場合との効率の差を比較するために行なった辞書検索回数の調査によれば、単語の候補数は、漢字音処理を行なった場合、行なわない場合よりも約40%削減できた。

## 9 おわりに

単語区切り率、辞書検索回数から考えて、本方式は十分な性能を有すると考えられる。

しかし、漢字音による処理により、かえって誤解釈を生じる例もみられたので、今後はそれらに対して辞書で個別に対応するなどの方法により、さらに変換性能を向上させる方向へ向けて、試行するつもりである。

## 参考文献

- 1 牧野ほか：べた書き文の分かち書きと仮名漢字変換  
情報処理，v o l 2 0 , N o . 4 (1 9 7 9 年)
- 2 内田ほか：自由入力形式のカナ漢字変換  
情報処理学会N L研資料 2 7 - 3 (1 9 8 1 年)
- 3 吉村ほか：最長一致法と文節数最小法について  
情報処理学会A I研資料 2 4 - 1 (1 9 8 2 年)
- 4 森岡健二：「近代語の生成—明治期語彙編」  
(明治書院、1 9 6 9 年)  
ただし、本文中の図は、林大監修「図説日本語」(角川書店、  
1 9 8 2 年)より再録した。
- 5 金田一ほか：「新明解国語辞典(第一版)」(三省堂、1 9 7 2 年)
- 6 安田賀計編：「企業経営文例全書」(ぎょうせい、1 9 7 8 年)