

# 日本語の活用処理

坂本義行 (電子技術総合研究所)

## 1. はしがき

自然言語の自動処理において、その処理の最小の意味単位の決定が必要であろう。この言語単位として、屈折言語(欧米語)では、「語」という概念があるが、日本語では、これにあたる形態上明確な定義がない。詞、辞、文節という単位は、その認定が不十分で、かつ分ち書きの問題も含めて日本語の計算機による機械処理の大きな問題になっておられると思われる。

ここで語の単位を英語の単語と対比させ、英和辞書の訳語項目に着目して、分析を行なった。これについては、すでに報告を行なった。<sup>(1)</sup>この結果から日本語の文節と呼ばれる単位に近い性質を有していると思われる。

この単位の認定の自動化を目標として、助詞あるいは字種連鎖に着目した自動分割処理、分割された単位を文節と定義し、文節を単位とする辞書の情報構造の設定、活用、接続空間の構造ならびに、その処理プログラムを作成したので、ここで報告する。

## 2. 分ち書きの手法

日本語の分ち書きには多くの問題を含んでいる。すなわち自動的な分ち書き処理 (automatic segmentation) の評価自体、問題がある。

自動分ち書きの手法については、種々の方法が提案されているが、大量な辞書等を用いた複雑な処理をもってしても完全なものはいずれも得られていない。むしろ人間が処理する単位とは別に、計算機による機械的処理に適する分ち書き方法を見出すべきであろう。ここでは、複雑な辞書を用いず、単純な手順による以下の(1)の2種類の手法を試みた。なお分割された単位としては、文節を目標とした。

### 2.1. 字種による分割

日本語の文書は、字種の混合文で記述されており、文字連系の字種が変化する部分により分割を行なうことは、機械的な処理として、十分有効な手段であることが知られている。

<字種> ::= <欧文字> | <片仮名> | <平仮名> | <漢字> |  
<特殊文字> | <未登録漢字>

<欧文字> ::= A | B | C ... a | b | c ... A | B | B ... a | p | r ...

<片仮名> ::= ア | イ | ウ ... ア | ズ | ヲ | ヱ

<平仮名> ::= あ | い | う ... め | り | ゑ

<漢字> ::= 々 | 亜 | 阿 ... (約 8,200字)

<数字> ::= 1 | 2 | 3 ... (1 | 2 | 3 ... 1 | 0 | 1 | 3) ... I | II | III ...

<特殊文字> ::= + | - | / | , | 。 ...

<未登録漢字> ::= || <片仮名> ||

日本語の混合格文中で漢字の果方役割は、意味の基底となる内容を提示し、漢字列に後続する平仮名列は、漢字の示す意味、概念の文中での役割を規定する場合が多い。また片仮名および欧文字で表現された文字列は、漢字列と同等の意味を有するものもある。したがって字種による分割を拡張し、漢字列または、それと同等の文字列に後続する平仮名列は結合し、これ以外の異字種の間で分割を行なう。

### 2.2 助詞による分割

日本語に於いては助詞が、欧米文における屈折がよむ語順に依存してゐる役割を担つてゐる。以下に助詞による分割方法について述べる。

(イ) 表2.1に示す格助詞による分割 (以下表は参考文献(2)による)

表2.1 格助詞表

が, を, に, へ, て, と, から, まで, の, と, や, とか, だの, より
---

(ロ) 表2.2に示す全助詞による分割

表2.2 助詞表

格助詞	が, を, に, へ, と, から, まで, の, と, や, とか, だの, より
副助詞	は, も,こそ, へ, こ, すら, だけ, ばかり, のみ, など, なん, まで
接続助詞	て, たり, たり, は, と, ても, ても, たゞ, たゞ, ながら, が, けれど, し, から, ので, のに
終助詞	か, ね, ず, よ, わ, な, ぜ, ぞ

(ハ) 表2.3に示す出現頻度上位の助詞による分割

表2.3 出現頻度上位10位までの助詞表

の, を, に, は, が, て, と, まで, (の), も
---------------------------------

(注) 表2.3表に採用した助詞は「電子計算機による新開の語彙調査Ⅱ」(国立国語研究所編)の助詞の頻度統計より頻度の高い助詞上位10位までを採用した。ただし「の」は用法の違いにより、1位と9位に存在するので、実際には9種について処理を行つた。

### 2.3 分割実験と考察

2.3.1 字種による分割では、比較的文節の単位で分割が行なわれるが、平仮名列の中に「詞」を含む場合が起り、すなわち、分割の必要条件を満たしてゐると考えられる。誤りの特徴として、(a) 漢字または平仮名のみからなる副詞。(b) 平仮名列中の助詞。(c) 平仮名・漢字列からなる詞。

2.3.2 格助詞による分割は、表中の助詞にフッては比較的正しく認定されているが、構文上重要で、かつ出現度の高い「は」、「も」が含まれていないため、分割の単位が大きくなり、最小単位に区切られていると云いがたり。(付録Iを参照)

2.3.3 全種類の助詞による分割では、全ての助詞が認定されるが、助詞として、非常に多くの平仮名列が含まれている結果、極端に誤りが増加する。とくに平仮名で構成される「詞」、副詞、接続詞等がほとんど分割され、「詞」としての形をとらぬなり。(付録IIを参照)

2.3.4 出現頻度による分割では、この実験結果から頻度の高い助詞は、構文上重要な役割を担っている助詞として正しく認定され、頻度の低い助詞は、他の文節の部分文字列である傾向がみられた。また、2.3.2、2.3.3 に比較して、比較的意味単位としてまとまった(文節)分割がなされた。

### 2.4 分割処理手続

連系(string)のパターン・マッチングが処理の主体となっている英から、SYMBOLSをアノテーション言語として用いた。この言語は文字のパターン・マッチングが主たる機能であり、複雑なパターンを簡単に構成できるように、各種の組込関数やプリミティブ・パターンを用意しておき、簡単なスタート・エンドで複雑なパターン・マッチングが行える英からこれを用いた。

例: 助詞表による文字列照合

KAKUJOSHI = '8+ : 9V'! '8V : 9L'! '8M : 8+!' '9- : 8U'! '8' "' : 9W'!  
                   か   ら    ま    て    と    か    た    の    ま    り  
 + '9+!' '83'! '8/'! '8Q'! '9L'! '8M'! '8U'! '8#'!  
           か    ま    た    へ    て    と    の    や

上記の表を定義すると、文(BUN)中の文字列との照合は、

BUN            KAKUJOSHI    :S(JOSHIARI)F(NASHI)

で簡単に行なえる。

### 3. 文節と辞書構造

日本語では、文と語の中間に、文節という単位がある。この文節の定義をBackus notationで表現したものを示す表を示す。

日本語の文節辞書の構造は、見出し、品詞、派生、活用、接続の5個の情報空間が身1回のように構成されているものとする。

1) 見出し - 派生による接辞や活用による語形変化の語尾を切り離した残りの文字列を基本の見出しとする。

第3表 文節の構造

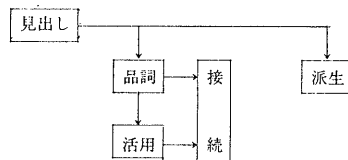
<文節>::=<詞> <文節> <辞>
<詞>::=<第一種の詞> <第二種の詞>
<第一種の詞>::=<体言> <用言>
<第二種の詞>::=<副詞> <連体詞> <接続詞>
<体言>::=<名詞> <代名詞>
<用言>::=<動詞> <形容詞> <形容動詞>
<辞>::=<助動詞> <助詞>

2) 品詞 - 各見出しに対し、1個以上の品詞が存在する。

3) 派生 - 見出しが、その形態が同一か、あるいは一部変化(接辞等が附加)して、他品詞となる場合、その変化形と派生による新品詞(派生品詞)を与える。

4) 活用 - 用言と助動詞に語形変化によって分類された活用型と助動詞、助詞への接続関係が与えられる。

5) 接続 - 活用する語は、後続の助動詞、助詞との結合で語形変化がなされ接続する。



第1図 文節辞書の情報空間

#### 4. 活用と接続

合成 - 用言、助動詞と助動詞、助詞との結合における語尾変化の表を表4.1表に示した。動詞に於いては、例外的な変化を有するものについては、類別別に項目として示す。さらに、助動詞、助詞と用言、助動詞との接続関係を表を表4.2表に示した。

分析 - 文節内の分析は次節で述べるように文節の左端より処理を行うため、合成に用いられた表を分解し、助詞、助動詞のみ、活用語尾、用言、助動詞の語幹への処理用に表4.3表と動詞の一部語幹、助動詞語幹を表を表4.4表に示した。

#### 5. 文節の分析手続

分析の対象となる文節は、2節で処理された出力について考える。  
すなわち、

〈文節〉 ::= 〈詞〉 | 〈文節〉〈詞〉 | 〈文節〉〈辞〉

ただし 〈辞〉 は空を認める

分析の処理手順の概略図を表5.1図に、その詳細図を表5.2図に示す。

各辞書の構造は、  
5.1 助詞辞書

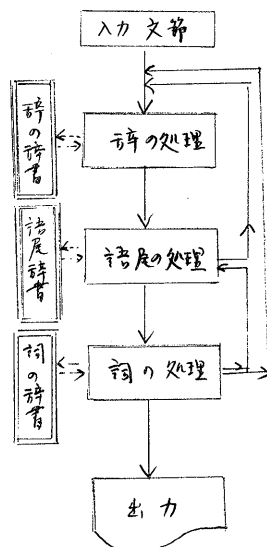
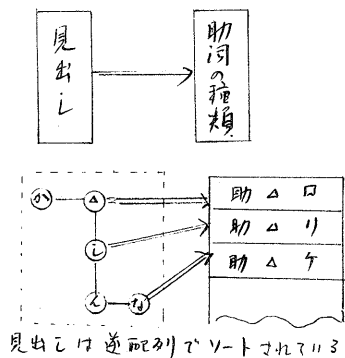
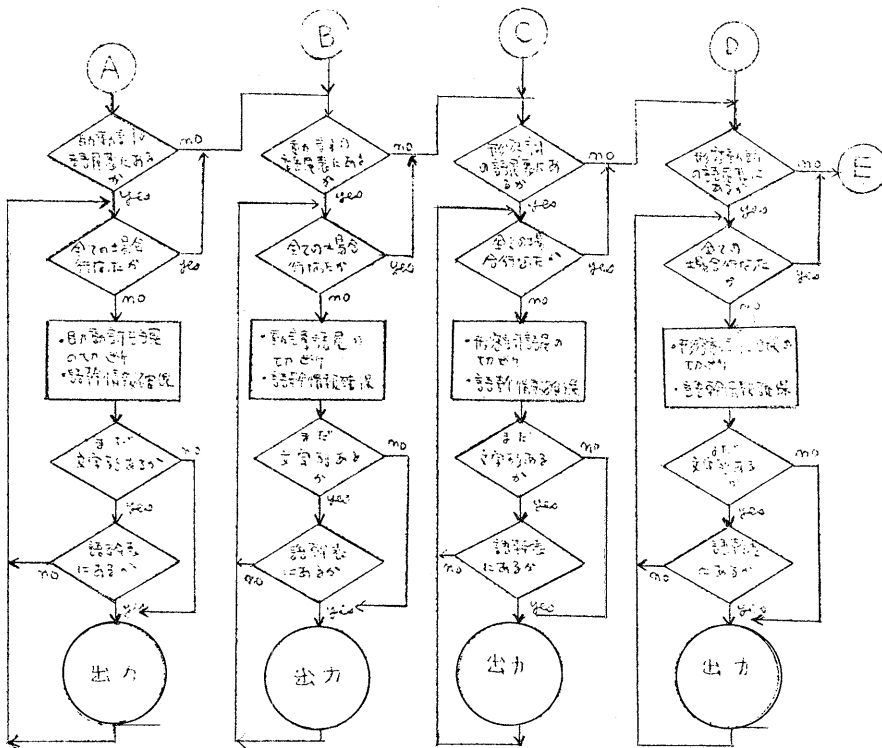
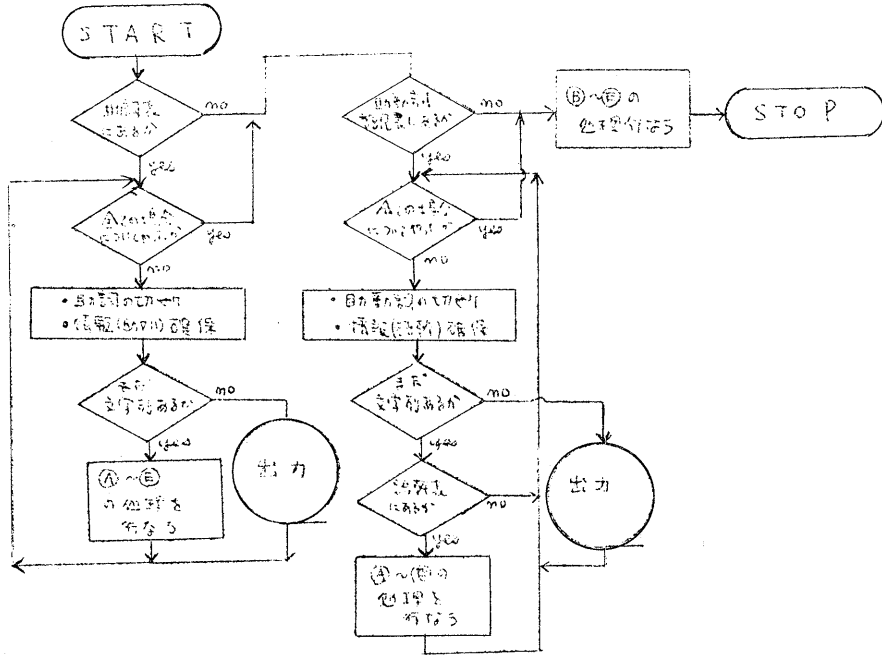


表5.1図 文節の分析手続

才 5.2 図 文節の分析手順図













4.4 表A 動詞の語幹表

type	ロ	ハ	ニ	ホ	ラ	ネ	ヘ	ト	ル	タ	レ	ツ
語幹と	在(る)	来(る)	きこえ(る)	で(る)	思(う)	ござ(る)	増え(る)	見え(る)	思出(す)	生(む)	漬(む)	怒(る)
語尾	う	かね	たり(る)	限(る)	おも(う)		ふえ(る)	知(る)	おもい(だ)	産(む)	住(む)	おこ(る)
表中の番号	1	2	3	4	14	13	7	8	9	10	11	12

オ	ク	ヤ-1	マ-1	マ-2	チ	リ	ヌ	ヲ	ワ	カ	ヨ	ソ	ナ	
在(る)	おっし	来	行(く)	(ゆく)	5段	5段	5段	5段	5段	5段	5段	5段	5段	
有(る)	ヤ(る)		往		カ	ガ	サ	タ	ナ	バ	マ	ラ	ワ	
あ(る)														
表中の番号	15	16	17	18	34	20	21	22	23	24	25	26	27	28

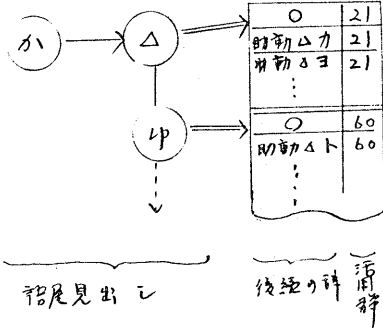
ム	ウ	キ	ノ	ヤ-2
サ変	(する)	サ変	サ変	(くる)
-1		-2	-3	
29	30	31	32	33

注: 表の上段に入る番号  
→ 表の下段に入る番号

4.4 表B 助動詞の語幹表

type	イ	ニ	ホ	チ	ヌ	ヲ	ワ	カ	ヨ	イ	ロ	ハ	ヘ	ヨ	カ	リ	ル	ト
基本型	よう	まい	ござ(います)	な(い)	で(す)	ま(す)	らし(い)	させ(る)	られ(る)	う	じや	な	た	れ(る)	せ(る)	(だ)	(ぬ)	(たい)
表中の番号	2	5	6	9	11	13	14	15	17	1	3	4	7	18	16	10	12	8

5.2 活用語尾辞書



語尾見出し

後述の辞  
活用辞

5.3 活用の種類群

種類群	活用型
1	イ
2	イ ロ / 1 =
3	イ ロ / 1 = ホ ヘ ト
...	...

5.5 処理結果

参考文献

- (1) 坂本義行: 『日本語の辞書』, CL研究会資料, 1975-3-15.
- (2) 鈴木重行 『日本語文法・形態論』, 東京書房。

活用接続処理出力例

入力文字列	文字列	用言語幹	活用型	語尾	助動詞	助詞
{ あるのに }	{ あ }	{ 5段ラ }	{ 動ソ }	{ る }	{ }	{ のに助ソ }
○ { あるのに }	{ あ }	{ 動オ }	{ る }	{ }	{ のに助ソ }	{ }
{ 見えるけれども }	{ 見え }	{ 5段ラ }	{ 動ソ }	{ る }	{ けれども助へ }	{ }
○ { 見えるけれども }	{ 見え }	{ 動ト }	{ る }	{ }	{ けれども助へ }	{ }
{ 見えるけれども }	{ }	{ }	{ }	{ }	{ も助ア }	{ }
○ { 思う }	{ 思 }	{ 動ラ }	{ う }	{ }	{ }	{ }
{ 注目して }	{ 注目 }	{ サ変2 }	{ 動イ }	{ し }	{ て助ヲ }	{ }
○ { 注目して }	{ 注目 }	{ (する) }	{ 動ウ }	{ し }	{ て助ヲ }	{ }
{ 注目して }	{ 注目 }	{ 5段サ }	{ 動ア }	{ し }	{ て助ヲ }	{ }
{ 注目して }	{ 注目 }	{ サ変1 }	{ 動ム }	{ し }	{ て助ヲ }	{ }
○ { を産む }	{ を }	{ 産 }	{ 動タ }	{ む }	{ }	{ }
{ であると }	{ であ }	{ 5段ラ }	{ 動ソ }	{ る }	{ と助オ }	{ }
○ { であると }	{ で }	{ あ }	{ 動オ }	{ る }	{ と助オ }	{ }