

# 日本語形態素解析の基本設計

坂本 義行  
(電子技術総合研究所)

## 概要

本研究は、日英機械翻訳システム<sup>1)</sup>における科学技術文献を対象とした日本語文の形態素解析を行い、分析結果として日本語構文解析等に必要情報を辞書システム<sup>2)</sup>より検索し、出力することとを目的としている。解析は、辞書駆動型による文字列間の接続関係と文の左端からチェックし、文を単位として、可能な解を全て出力する Multiple-path 方式をとっている。本稿では サンプル文として JICST 文献抄録文について実験を行っている。そこに出現した文の表現形式を細部にわたり分析し、実用的な解析システムを構築するプロセスの開発をめざしている。

## 1. はじめに

日本語の形態素解析はその目的、用途によって、その方法も、自から異なったものとなるであろう。本研究では、日英機械翻訳システムとして必要な処理形態、言語的な興味といった目的も括って実用という面から、文献に現われる表記も可能な限り辞書に登録することにより解析の可能性を確かめるとともに、処理ソフトウェアと辞書とを明確に区別し、処理手順の簡略化、拡張の容易さをめざした。

日本語構文解析<sup>3)4)</sup>に必要とされる情報は、すべてこの形態素解析レベルで辞書を引くことにより付与する方式をとっている。

処理手順は、大きく次の3つからなる。

- 1) 入力テキストの定形化処理
- 2) 形態素の分割
- 3) 形態素解析の出力処理

## 2. 入力テキストの定形化処理

翻訳の入力文である JICST の抄録文から日本語形態素解析の入力文となる標準形式をつくりだすために以下のようなフォーマットを行う。

- 1) 抄録番号、記事番号といった書誌的事項の区別
- 2) 標題、本文、式、図などの区別
- 3) 印刷書式等の内容に直接関係しない部分の区別

### 2.1 システム中での処理の流れ (図1)

#### 2.2 フォーマット情報

フォーマット情報の取扱いを以下の4種に分ける。

A1. 総合システムが管理するもの (翻訳フラグをNILとする)

例. @p@e, @AR@

A2. 総合システム<sup>5)</sup>がすべて管理するもの。フォーマット情報のうち首尾照応するもので、中間に翻訳対象を含むもの。

本研究は国の科学技術振興調整費による「日英技術文献の通報システムに関する研究」の1部として行ったものである。

研究協力者として、木村睦子 (計量計画研究所)、山本稔 (東洋情報システム) の諸氏、ならびに本研究の委託を受けて、その遂行のために、自然言語処理の専門家から成る言語処理システム作業分科会、さらに処理システム作業グループ (京大)、辞書作業グループ (JICST) を組織し、その審議、指導のもとに研究をすすめている。

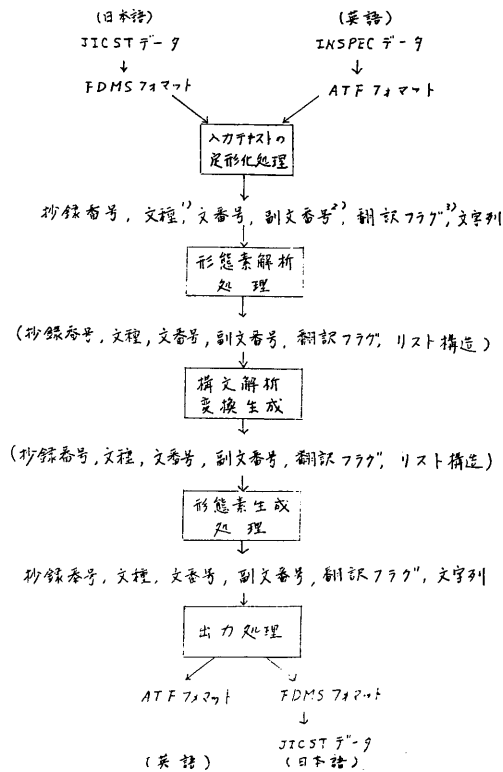


図1 翻訳システムのデータ形式

- 1) 文字列 (文) の種類。例: 標題, 本文, レイアウト。
- 2) 同一文番号内での枚番号。一文中にフォーマット情報が付いて複数のレコードに分かれた場合
- 3) 翻訳を必要とするかどうかのフラグ  
 T: 翻訳が必要 (通常の文)  
 NIL: 翻訳が不要 (主としてページ様式に関するフォーマット情報など)

フォーマット情報と翻訳対象文を分けてレコードにする。

例: @HS=H12.3@抄録の機械翻訳システム@HE@

B: 形態素解析の時点で単語の性質に反映するもの (単語情報)

例: @UN@システム@UE@  
 ただし、複数の単語にまたがって指定されている場合は「C」で扱う。

例: @UL@ある種の単語@UE@  
 また、単語の一部だけが指定されている場合は、その単語全体に対して指定したものとみなす。

C: 句読点と同様に一語扱いとして、文法処理に引き渡すもの。

### 3. 形態素の分割

漢字かな混じりの日本文を文節、語あるいはさらに小さな単位に分割する処理方法として、その手順をプログラム化して行う方法と辞書の検索を中心としたプログラムにより行う、大きく分けて2種類の方法がある。前者は「小まわりのきく」方法であり、後者は「辞書」に依存するためプログラムは単純化されるといった特徴を有している。本研究では翻訳に適合する、すなわち英語に訳すに最適であり、改良、拡張が容易な後者の方法である辞書駆動型の解析手順を選んだ。

テキスト名: 0121111 83.07.01

項目番号	参照番号	
0.0.0	1. 1. 0	@PG=50.66.1.5.9.4@BAR=5.5@HNS=H12.3@
0.0.0	1. 3. 0	抄録の機械翻訳システム@HE@このシステムは、JICST
0.0.0	1. 4. 0	の抄録を計算機により翻訳するものである。本システムの入力は、たとえば、富士通のFDMSのようなフォーマット機能をもったシステムに対する入力であってもよい。@NLB@このシステム
0.0.0	1. 5. 0	の特長は、次のとおりである。
0.0.0	1. 6. 0	

項目番号	参照番号	
0.0.0	1. 1. 0	E82060001, レイアウト, 0, 10, NIL, "@PG=50.66.1.5.9.4@"
0.0.0	1. 2. 0	E82060001, レイアウト, 0, 20, NIL, "@BAR=5.5@"
0.0.0	1. 3. 0	E82060001, レイアウト, 0, 30, NIL, "@HNS=H12.3@"
0.0.0	1. 4. 0	E82060001, 機械, 0, 40, T, "抄録の機械翻訳システム"
0.0.0	1. 5. 0	E82060001, レイアウト, 0, 50, NIL, "@HE@"
0.0.0	1. 6. 0	E82060001, 本文, 1, 10, T, "このシステムは、JICSTの抄録を計算機により翻訳するものである。"
0.0.0	1. 7. 0	E82060001, 本文, 2, 10, T, "本システムの入力は、たとえば、富士通のFDMSのようなフォーマット機能をもったシステムに対する入力であってもよい。"
0.0.0	1. 8. 0	E82060001, レイアウト, 2, 20, NIL, "@NLB@"
0.0.0	1. 9. 0	E82060001, 本文, 3, 10, T, "このシステムの特長は、次のとおりである。"

図2. 正規化処理の事例

### 3.1 分割の特徴

- 1) 構成要素 - フォルダ、辞書、接続表等テーブル類
- 2) 処理方式 - Multiple-path 方式として、1)の言語情報により形態的に接続可能なあらゆる単位語列を作成する。
- 3) 処理単位 - 設定した「文」を単位にその中の文字列がすべて接続可能な場合のみを出力する。

### 3.2 辞書引込の方法

辞書には自立語(複合語を含む)、付属語、接辞等を見出しし、形態素解析に必要な情報として、図1に示すような見出し情報、形態素情報が与えられる。

見出しは、用言については、終止形で表わされ、語尾字数によって語幹を判定する。異形語はすべて見出しと同等に扱う。分割の可否を判定するに

表1. 形態素辞書フォーマット

見出し語	語尾字数	
	漢字部	
見出し情報	語基	
	読み	
異形語	誕生語	
	関連語	
形態品詞	名 副 名 動 形 形 動 副 連 体	
	接 助 助 助 接 頭 接 尾	
動詞活用型	五 上 一 下 一 変 変 変 変	
	カ カ ガ サ タ	
形態活用型	ナ バ マ ラ	
	ダ ナ	
助動詞活用型	形 形 助 助 特	
	格 接 副 並 接 準	
細分類	接尾	体 助 形 形 助
	前接	
情報	接接	
	情報	

は、後接情報(後続語との接続関係と示<sup>2)</sup>)と前接情報(先行語との接続関係と示<sup>3)</sup>)との結合を接続表(表3)により判別する。すなわち、先行語の後接情報と後続語の前接情報とがうまく整合するかどうか。ブランク、1、2等の記号で示された表によって判別する。ブランク: 接続不可

- 1: 接続可能だがそこで切つてはいけない。
- 2: 接続可能で、かつ単語の切れ目となる。

表3-1 後接情報

コード	内 容	コード	内 容	コード	内 容
1	名詞	26	格助詞、並助詞、と、か	50	五段動詞語幹
2	副詞的・名詞(含数詞)	27	接続助詞	51	〃
3	副詞	28	接頭辞	52	〃
4	連体詞	29	五段動詞未然形(～ぬ)	53	上一段動詞語幹
5	接続詞	30	〃 (～う)	54	〃
6	数字	31	〃 接尾(清音)・接音便	55	〃
7	上・下・一 力変・サ変連用	32	〃 (濁音)	56	〃
8	五段動詞(サ行以外)連用	33	動詞・形容詞仮定形	57	〃
9	動詞終止形	34	上・下・一 力変未然形	58	〃
10	動詞・形容詞連体形	35	サ変未然形	59	〃
11	動詞命令形	36	〃	60	〃
12	上・下・一 語幹末語尾	37	〃	61	下一段動詞語幹
13	サ変語幹	38	形・形容詞未然形	62	〃
14	〃 終止形	39	〃 連用形	63	〃
15	形容詞連用形	40	サ変語幹	64	〃
16	〃 終止形	41	形容詞語幹	65	〃
17	形容詞連用形	42	形容動詞(サ型)語幹	66	〃
18	〃 終止形	43	〃 (サ+型)	67	〃
19	〃 連体形	44	五段動詞語幹	68	〃
20	〃	45	〃	69	〃
21	〃	46	〃	70	〃
22	〃	47	〃	71	〃
23	格助詞	48	〃	0	終止状態
24	〃	49	〃		
25	〃	50	〃		

表3-2 前接情報

コード	内 容	コード	内 容	コード	内 容
1	名詞	26	の、に、を	51	五段動詞語尾
2	名詞以外の自立語、接頭辞	27	格助詞	52	〃
3	形式名詞	28	〃	53	上一段動詞語尾
4	補助動詞(主に接続助詞)	29	〃	54	〃
5	助動詞	30	〃	55	〃
6	助動詞	31	〃	56	〃
7	〃	32	格助詞、と、に、を、並助詞	57	〃
8	〃	33	と(引用)	58	〃
9	〃	34	助詞	59	〃
10	〃	35	係助詞	60	〃
11	〃	36	副助詞	61	下一段動詞語尾
12	〃	37	接尾辞	62	〃
13	〃	38	〃 (体言接続)	63	〃
14	〃	39	サ変語幹・接続助動詞	64	〃
15	〃	40	サ変語尾	65	〃
16	〃	41	サ変語尾	66	〃
17	〃	42	形容詞語尾	67	〃
18	〃	43	形容動詞	68	〃
19	〃	44	五段動詞語尾	69	〃
20	接続助詞	45	〃	70	〃
21	助動詞	46	〃	71	〃
22	接頭辞	47	〃	0	終止状態(サ接点)
23	〃	48	〃		
24	〃	49	〃		
25	〃	50	〃		

活用語尾処理は、活用語尾表を用い、その前接情報により判定する。なお語幹を持たない語は、語尾導致により判定し活用処理を行う。(辞書台帳には特殊活用型として記載されており、辞書生成時に見出し等の情報が独立して記憶される。)

### 3.3 分割の手順

#### 3.3.1 辞書による単語抽出

処理は辞書引きから始まる。文字列を左端から一字づつ辞書を引き、1回の操作で最長一致する文字列まですべてをスタックする。スタック中の最長一致の文字列に対して、そのつぎの文字列について辞書を引き接続ナエツプを行い可能な文字列をスタックする。活用する語には、活用語尾処理を施す。以上の操作をくり返す。接続不可となったら、以下の処理を実行する。

#### 3.3.2 特殊記号処理

##### ① 特殊記号の判別方法の特徴

- ・文脈(前後の文字の字種や、文字の並びのパターン)によって判別される記号が多い。
- ・確定的な判別法はほとんどなく、実際の文例をもとにした経験的な手法が多い。

##### ② 特殊記号の取り扱いの特徴

- ・前の文字と同一の字種として扱う。
- ・後の文字と同一の字種として扱う。
- ・区切り符、特殊符号や「つなぎ」の符号等のように独立の符号として取り扱う。

また、翻訳対象の広がりと共に、判別方法の変更の追加等が頻繁に行われ、その取り扱い方法も拡張されること予想される。

以上の裏から、ソフトウェアは機能の追加、変更柔軟に対応できるようなテーブル駆動型の構造をとった。

##### ③ 特殊記号処理の手順

特殊記号処理は、辞書で引けなくなった文字列に適用され、特殊記号

処理テーブルにより判別を行う。特殊記号テーブルは、特殊記号、処理適用条件と実行処理からなる。

#### 3.3.3 未知語の処理

辞書になく、特殊記号処理によっても単語抽出ができなかった場合で、かつ1つ前のスタックにも、文完成スタックにも解がないとき、字種によるスキップ処理を行い、未知語としてスタックし、単語ポインタを変更する。スキップ位置の決定方法として、未知語の先頭文字が、

漢字、

カタカナ、

英字、

ギリシャ文字

指定した特殊記号、であった場合にその字種が連続する最後の文字の次の文字までスキップする。

#### 3.3.4 バックトラック探索

##### 1) 1文完成内でのバックトラック

単語抽出に失敗したとき、それ以前にある単語スタックを戻し、単語リストを書き換え、単語ポインタを変更し、3.3.1からの処理を繰り返す。

##### 2) 文完成後のバックトラック探索

文末までの処理が完了した後、マルケパスを行うため、残余の単語スタックを探索し、スタックが空になるまで処理を続行する。

### 3.4 形態素解析の出力処理

#### 3.4.1 辞書項目の付与

分割処理が終了し、文が完成した場合に辞書項目の付与を行う。

##### 1) 見出し情報

見出しの形(活用語は、その見出しの形)、派生語、関連語等表に記述されている情報

##### 2) 形態素情報

活用語には、その活用形を示す

##### 3) 構文・意味情報

##### 4) 格支配情報

5) 共起情報

6) 備考

以上のすべての情報をLISP-S式にて図3のように付与する。

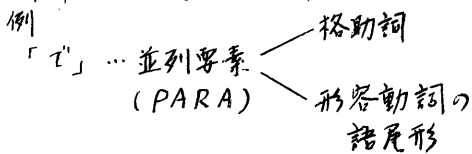
3.4.2 日本語構文解析への出力

1) 複数の文の出力

分割が異なる解が得られたとき、文を単位として複数の文を出力する。

2) 複数の形態素情報をもつ(PARA)

単位として同一の分割が行われたが複数の形態素情報が得られ、形態素解析では判別不能の場合



3) 連語

連続した文字列間に形態素的な結合関係がある。



表2-3 形態素分析継続表(1)

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	
2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43		
3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43			
4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43				
5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43					
6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43						
7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43							
8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43								
9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43									
10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43										
11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43											
12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43												
13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43													
14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43														
15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43															
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																
17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																	
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																		
19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																			
20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																				
21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																					
22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																						
23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																							
24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																								
25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																									
26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																										
27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																											
28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																												
29	30	31	32	33	34	35	36	37	38	39	40	41	42	43																													
30	31	32	33	34	35	36	37	38	39	40	41	42	43																														
31	32	33	34	35	36	37	38	39	40	41	42	43																															
32	33	34	35	36	37	38	39	40	41	42	43																																
33	34	35	36	37	38	39	40	41	42	43																																	
34	35	36	37	38	39	40	41	42	43																																		
35	36	37	38	39	40	41	42	43																																			
36	37	38	39	40	41	42	43																																				
37	38	39	40	41	42	43																																					
38	39	40	41	42	43																																						
39	40	41	42	43																																							
40	41	42	43																																								
41	42	43																																									
42	43																																										
43																																											

注 \*K~\*I ; 対角線上のみ

4) 複合語

辞書に見出しとして登録されているが語基情報をもっている



4. おわりに

本稿では、5ヶ年度行った基本設計に基づき、本年度開発に着手した段階での報告が、完成されたものではない。今後、実験をすすめて、分割レベルの問題、構文解析へのインターフェイス等、多くの改良、拡張を必要とすると思われる。マルチパス方式の収束性の問題がある。それらについては、漸次報告する予定である。

参考文献

- 1) 長尾真, 「科学技術系機械翻訳プログラムの概要」。
  - 2) 坂本義行, 「格構造を中心とした用言と付属語」。
  - 3) 辻井潤一, 「日本語構文解析」。
  - 4) 中村順一, 「文法記述用ソフトウェア」。
  - 5) 矢田光治, 長尾真, 「機械翻訳総合システムの基本設計」
- 以上のいずれも、本会で同時報告されたものである。

\*\*\* INPUT DATA \*\*\* => (E R 2 0 0 0 0 0 3 文 3 3 1 T システム工学、情報工学、ソフトウェア工学など新しい技術の急  
速な成長により、I E E および電気工学教育の全パタンでの重点が移動した。)

(E R 2 0 0 0 0 0 3 文 3 3 1 T  
(\$形態素分析)

(\$辞書情報  
(\$見出し情報 (\$見出し語 "システム工学")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 51629))  
1)  
(\$区切り符 ".")  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "情報工学")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 52763))  
1)  
(\$区切り符 ".")  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "ソフトウェア工学")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 59726))  
1)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "など") (\$読み "など"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 36)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (意味区分 "など") (格助詞との代弁性 "1"))  
(レコード番号 394))  
2)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "新しい")  
(\$語尾次数 1)  
(\$漢字部 1)  
(\$読み "あたらしい")  
(\$形態素情報 (\$形態品詞形) (\$後接情報 41))  
(\$原文-意味情報 (\$格助詞 (形容詞)  
(\$漢字部分類 "性質・状態")  
(\$格パターン "A 1")  
(\$辞書属性 "尺度")  
(\$名詞性程度 "一"))  
(\$格支配情報 (格パターン "A 1")  
(\$パターン "A 1")  
(\$表現格 "が" "深層格" "S U B" "名詞意味コード" "必須性" "1")  
(レコード番号 322))  
10)  
(\$辞書情報 (\$見出し情報 (\$見出し語 "技術")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 50983))  
1)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "の") (\$読み "の"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 31)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (表層格名称 "の") (修飾関係 "連体"))  
(レコード番号 441))  
5)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "だ") (\$読み "だ"))  
(\$語尾次数 1)  
(\$漢字部 2)  
(\$読み "うそくだ")  
(\$形態素情報 (\$形態品詞助)  
(\$前後接情報 2)  
(\$後接情報 43))  
(\$原文-意味情報 (\$格助詞 (形動詞)  
(\$漢字部分類 "性質・状態")  
(\$格パターン "A 1")  
(\$辞書属性 "尺度")  
(\$名詞性程度 "一"))  
(\$格支配情報 (格パターン "A 1")  
(\$パターン "A 1")  
(\$表現格 "が" "深層格" "S U B" "名詞意味コード" "必須性" "1")  
(\$表現格 "は" "深層格" "C O R" "名詞意味コード" "必須性" "0"))  
(レコード番号 103))  
19)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "成長")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 52183))  
1)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "により") (\$語基 1 1 2) (\$読み "により"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 27)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (表層格名称 "により") (修飾関係 "連用")  
(\$漢字部 "原因-理由" "手段-道具" "条件")  
(レコード番号 437))  
5)  
(\$区切り符 ".")  
(\$未知語 "I E E")  
(\$未知語 "および")  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "電気工学教育")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 54376))  
1)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "の") (\$読み "の"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 31)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (表層格名称 "の") (修飾関係 "連体"))  
(レコード番号 441))  
5)  
(\$未知語 "今")  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "パターン")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 60954))  
1)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "で") (\$読み "で"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)

(\$前後接情報 32)  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 31)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (表層格名称 "で") (修飾関係 "連用")  
(\$漢字部 "手段" "場所" "原因-理由" "手段-道具" "材料" "構成要素"  
"方法" "条件" "観点" "その他"))  
(レコード番号 369))  
25)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "の") (\$読み "の"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 31)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (表層格名称 "の") (修飾関係 "連体"))  
(レコード番号 441))  
5)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "重点")  
(\$形態素情報 (\$形態品詞名)  
(レコード番号 52633))  
1)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "が") (\$読み "が"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞細分類格)  
(\$前後接情報 27)  
(\$後接情報 5))  
(\$原文-意味情報 (\$格助詞 (表層格名称 "が") (修飾関係 "連用") (深層格 "主体" "対象"))  
(レコード番号 351))  
5)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "移動する")  
(\$漢字部 1)  
(\$読み "いどうする")  
(\$形態素情報 (\$形態品詞動)  
(\$動詞活用形変)  
(\$前後接情報 2)  
(\$後接情報 13))  
(\$原文-意味情報 (\$分野コード "電気")  
(\$漢字部分類 "動詞")  
(\$パターン "V 1")  
(\$アスペクト "瞬時")  
(\$意味 "後")  
(\$格支配情報 (格パターン "V 1")  
(\$表現格 "が" "深層格" "S U B" "名詞意味コード" "必須性" "1")  
(\$表現格 "に" "へ" "深層格" "S T O" "名詞意味コード" "必須性" "0")  
(\$表現格 "から" "深層格" "S P R" "名詞意味コード" "必須性" "0")  
(\$表現格 "に" "へ" "深層格" "S T O" "名詞意味コード" "必須性" "0")  
(\$意味 "変化" "可能")  
(\$漢字部 "有")  
(\$格支配情報 (格パターン "V 2")  
(\$表現格 "が" "深層格" "S U B" "名詞意味コード" "必須性" "1")  
(\$表現格 "に" "へ" "深層格" "O B J" "名詞意味コード" "必須性" "1")  
(\$表現格 "に" "へ" "深層格" "S T O" "名詞意味コード" "必須性" "0")  
(\$表現格 "から" "深層格" "S P R" "名詞意味コード" "必須性" "0")  
(\$表現格 "に" "へ" "深層格" "S T O" "名詞意味コード" "必須性" "0")  
(レコード番号 512))  
2)  
(\$辞書情報  
(\$見出し情報 (\$見出し語 "た") (\$語尾次数 1) (\$読み "た"))  
(\$形態素情報 (\$形態品詞助)  
(\$助詞活用形変)  
(\$後接情報 不定))  
(\$特殊活用形 (不定形 (形 "たろ") (後接情報 30))  
(\$終止形 (形 "た") (後接情報 16))  
(\$語尾 (形 "た") (後接情報 16))  
(\$決定形 (形 "たら") (後接情報 5))  
(\$原文-意味情報 (\$漢字部分類 "動詞")  
(\$語尾 "過去" "完了")  
(レコード番号 100))  
16)  
(\$辞書情報  
(\$区切り符 ".") (\$読み ""))

### 図3 形態素解析の出力例