

151
電子通信学会(1974年1月23日)
資料刊転載

資料番号AL73-71

自然言語分析のためのプログラミング言語—PLATON— とそれによる名詞句の分析

A PROGRAMMING LANGUAGE FOR NATURAL LANGUAGE ANALYSIS AND ITS
APPLICATION FOR JAPANESE NOUN PHRASE ANALYSIS

長尾 真 辻井 潤一

MAKOTO NAGAO JUN-ICHI TSUJII

京都大学工学部
KYOTO UNIVERSITY

A Programming Language for Tree Operation -PLATON- is constructed. This accepts strings, trees and lists and transforms them in arbitrary way. Based on the argued transition network of W. Woods, this has the additional capability of matching strings, trees and lists, variable assigning, transformational rewriting. For an efficient rule application, the pre-checking of the applicability of a rule is introduced, and also the exit destination can be specified when a return is made from a lower state to an upper state. This avoids a tedious time-consuming back-tracking process. An example of using PLATON is given, which analyzes Japanese noun phrases containing post-position 'NO'

1. はじめに

W.A. Woodsらの augmented transition networkは、push-down automatonに種々の拡張機能をつけることにより、Turing Machineと同等の認識能力を持たせたモデルであるが、自然言語の認識機構としても、次のような幾つかのすぐれた特徴を持っている。

- (1) 任意の transformational grammar を記述できる。
- (2) ルールの追加, 変更, 削除が容易なため, 実験システムに向いている。
- (3) ルールを表として与えるより効率がよい。
- (4) 任意の条件 check を LISP 関数として組み込むことが可能であるから, 種々の意味論モデルとの結合が可能である。
- (5) 文法の明晰さが保存されている。

しかしながら, Woodsらのモデルでは入力シンボルを順次読み込むことにより, 状態遷移をおこなうため, 語順が比較的自由的な日本語に適

用するには困難があると考えられる。日本語においても, 意味のまとまりの単位である句内の語順はある程度一定している。しかしながら, 句間の関係は, 主として英語では語順によって表現されるのに対して, 日本語では格助詞によって示されるため, 語順自体にはかなりの自由度がある。

したがって, 英語の解析の際に, しばしば用いられてきた, 入力文章の単語を先頭から順々に scan して解析をすすめるゆく方法は日本語に対しては不適切であると考えられる。

日本語の解析においては, pattern-matching 的な手法で処理単位を切り出し, その処理単位内での構造を解析し, その結果の構造を, さらに大きな構造に組み上げてゆくといった方法が有効であると考えられる。

我々はこの様な考え方を基本として, 自然言語(日本語)処理用のプログラミング言語 PLATON (Programming Language for Tree Operation) を開発した。

一般に言語解析には Top-Down 的手法と、Bottom-Up 的手法があるが、network モデルは本質的に Top-Down 的手法になっている。Top-Down 的手法では、ある構造の存在を予測して、対応するルールを適用してゆくわけであるが、そのルールが不適切であった場合、それをどのように回復するかが問題となる。失敗した原因を知ることによって、適切なルール適用に切替るメカニズムをモデルは備えていなければならない。Woods のモデルでは状態内のルールの順位によって、このメカニズムを implicit に表現しているが、前のルール適用失敗の原因によって処理の流れを柔軟に変化させるには不十分であると考えられる。PLATON では失敗原因を failure-message として返すことによって、より柔軟な処理が可能になっている。

この資料では、PLATON 言語の各種機能と、その仕様を解説し、PLATON を使用して書かれた具体的なプログラム（名詞を中心とした句の処理）について説明する。なお、PLATON は記号処理用言語 LISP を用いて coding されている。

2. Tree 構造と List 構造

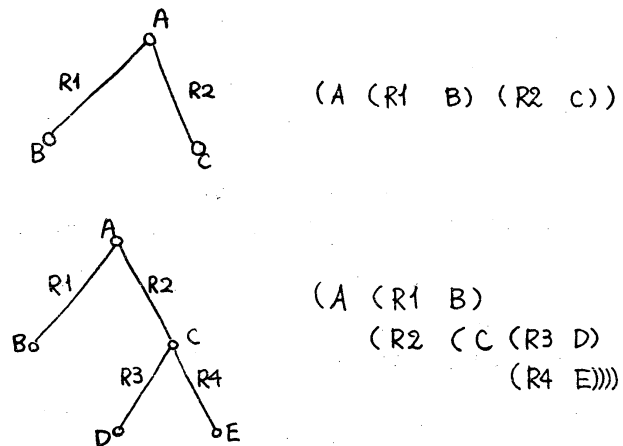
この節では PLATON が使用する基本的なデータ構造について述べる。

自然言語の処理においては、自然言語文の解析結果を tree 構造として表現することが普通おこなわれる。そこで PLATON が tree 構造をどのように表現するかを次に述べる。

[PLATON における Tree 構造の定義]

- ① <tree> ::= <node> | (<node> <branches>)
- ② <branches> ::= <branch> | <branch> <branches>
- ③ <branch> ::= (<relation> <tree>)
- ④ <node> ::= 任意の LISP 原子記号
- ⑤ <relation> ::= 任意の LISP 原子記号

以上の形式的定義を、具体的な tree 構造とその表現形式で示すと(図-1)になる。



<図-1> Tree 構造と表現形式

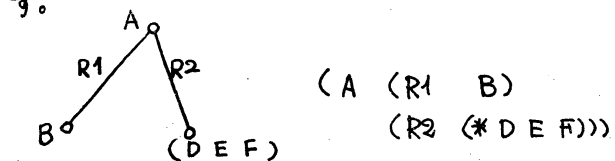
自然言語の解析は一次的な単語配列(list)から、それに対応する tree 構造(parsing tree)を作り出す過程であると考えられる。その途中段階では、すでに解析が済んで tree 構造に整理された部分と、まだ list のままの部分とが混在することになる。そこで tree 構造と list 構造の混在を許した表現方法が必要になる。この tree 構造と list 構造の混在を許した data 構造は我々のシステムの基本となるものである。以下このようなデータ構造を単に構造という。

次にこの構造の形式的定義を与えるが、tree 構造中の node 自体が list 構造になることをゆるすために、先の定義の④が変更される。

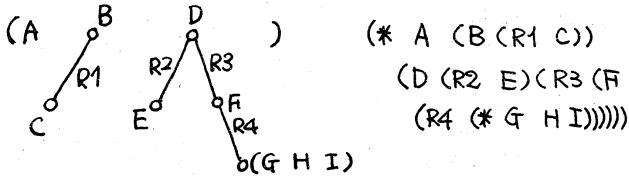
[PLATON における構造の定義]

- ⑥ <structure> ::= <tree> | <list>
- ⑦ <list> ::= (* <structures>)
- ⑧ <structures> ::= | <structure> | <structure> <structures>
- ④' <node> ::= 任意の LISP 原子記号 | <list>

具体的な構造と、その表現形式を<図-2>に示す。



(A B C) (* A B C)



(図-2) 構造と表現形式

3. パターン・マッチングの機能と方法

自然言語の文法の各ルールには、そのルールを適用するべき‘構造の pattern’がある。すなわち、書換法則の左辺にあたる部分がある。PLATONでは、この部分を変数を含む構造によって表現する。そして、この変数を含む構造と処理対象である具体的な構造との matching をとることによって、そのルールを適用すべき pattern があるかどうかを判断する。さらに処理対象の構造内での変数に対応する部分を切り出して、その部分を他のルールで処理するように指示することも出来る。

次のような変数が、2で定義した <structure> の中で使用出来る。

[PLATONにおける構造変数]

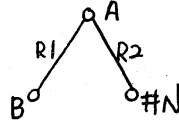
- #N 数字 ; 任意の mode と matching 可能。
 - #P 数字 ; 任意の structure と matching 可能。
 - #B 数字 ; 任意の branch と matching 可能。
 - #K 数字
 - #I 数字
 - #J 数字
- }; 任意個の list 要素との matching 可能。

(*) 数字は任意の自然数である。

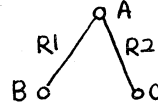
PLATON のシステム関数 MATCH は上記のような変数を含んだ structure と、任意の structure との matching をとり、matching の成否を値としてかえす。その際 matching が成功した時には各変数が具体的な構造の中で、どの部分と対応がとれたかを、ドット対としてかえす。具体的な例がこのことを示す。(<図-3>)

例1。

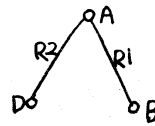
<変数を含む structure>



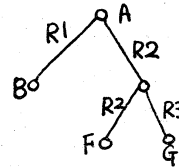
<具体的な structure と matching の結果>



成功 ; ((#N. C))



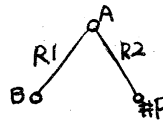
成功 ; ((#N. D))



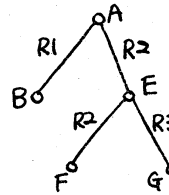
失敗

例2。

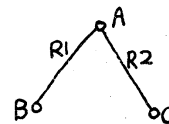
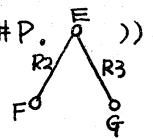
<変数を含む structure>



<具体的な structure と matching の結果>



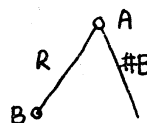
成功 ; ((#P. E))



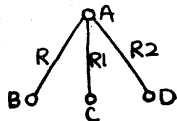
成功 ; ((#P. C))

例3。

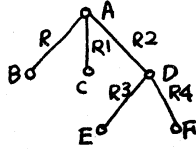
<変数を含む structure>



<具体的な structure と matching の結果>



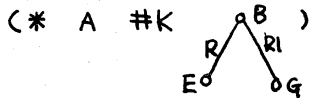
成功; ((#B.(/R1 /R2)))



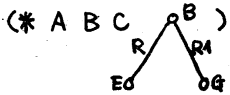
成功; ((#B.(/R1 /R2)))

例4.

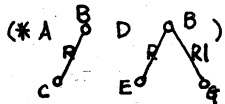
<変数を含む structure>



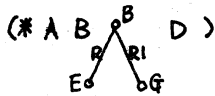
<具体的な structure と matching の結果>



成功; ((#K.(* B C)))



成功; ((#K.(* B D)))



失敗

<図-3> matching の例1~4

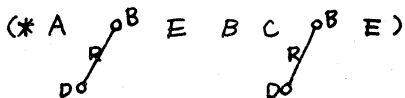
もちろん、上記変数の複数個が1つの構造の中で使用されてもかまわない。また同じ変数を2度以上使用することによって同一の構造が離れて存在することを表現することも可能である。<図-4>にこの例を示す。

例5.

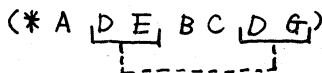
<変数を含む structure>

(* A #K B C #K)

<具体的な structure と matching の結果>



成功; ((#K.(* B E)))



<図-4> matching の例5 失敗

tree構造内のbranchにはすべてラベルがついているために、branchの順序が違っても matching は成功する。

ただし、list変数(#K, #I, #J)の使用には、連続して使用出来ないという制限があり、

(* #K #I A)

は不可である。

pattern-matchingのための関数MATCHは、指定された構造パターンがどんなに深いレベルにあっても、見つけ出すことが出来るようになっている。

4. 状態遷移とルールの表現法

プログラムで文法を記述するのと比較して、一様な syntax を持つルールで文法を記述することの利点としては、

- (1) ルールの変更、追加が容易である。
- (2) 文法の明快さ、明晰さが保持される。

がある。計算機による言語処理の立場から見れば、現在ある自然言語のすべての文法は不完全であると考えられる。したがって、その都度、変更、追加をすることによって、よりすぐれた計算機のための文法を作つてゆかねばならないと思われる。このことから、種々の文法を計算機に implement し、test し、簡単に改良できるような system が必要になる。(1)、(2)の特徴はこのような system が是非とも備えていなければならぬものだと考えられる。

逆に一様なルールの集合を基本とする解析システムの欠点としては、

- (1) 柔軟な処理、特に semantics との結びつきを表現することがあつかしい。
- (2) 複雑な高度な処理をしようとする、例外的な文章を扱うためのルールが増加して、効率が悪くなる。

特に(1)の欠点のために、従来の解析システムでは syntax analysis の部分を一様なルールで処理し、その結果を独立した semantic analyzer

に引渡すという手法をとっていた。

しかしながら、実際には *syntax analysis* と、*semantic analysis* とは交互に、相互依存的に働く必要がある。この欠点を一様なルールを使用するシステムでは解消する必要がある。

(2)の点に関しては、Woods 流の *transition network model* では、例外的なルールは例外的な状態(*state*)に置かれているために、ルール数の増大が必ずしも効率の著しい低下に結びつかない利点がある。

我々のシステムでは Woods 流の *transition network model* を採用し、一様な *grammar* の持つ利点を保持しながら、*user* が自由にプログラムを記述する部分を *rule formal* の中に設けることにより、柔軟な処理も表現出来るようになっている。

< 4.1 > PLATON のルール記述の定義

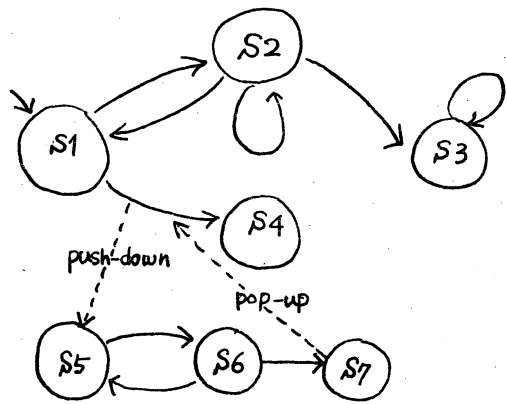
PLATON の基本的なモデルをまず説明する。

モデルの基本的な部分は、複数個の状態(*state*)とそれらを結ぶ方向性を持った枝(*branch*)からなる *network* である。各 *branch* は 1 個の書換ルールに対応しており、そのルールを使用し得るための条件、そのルールが使用され得るための構造パターン(書換ルールの左辺)、及びそのルール適用後の構造(書換ルールの右辺)が示めされている。そして、ある *branch* に対応するルールを適用した後は、その *branch* が入ってゆく *state* に状態を移し、その *state* に付いたルール(*branch*)は、再び同じように適用してゆく。(< 図-5 > 参照)

各 *state* から複数個の *branch* が出ており、各 *branch* は優先度にしたがって一列に並べられている。ある *state* に *control* が移ると、システムはその *state* から出てゆく *branch* を優先度にしたがって順々に *check* してゆく。まず *branch* についての条件を *check* し、処理対象となる構造がその *branch* についての構造パターンを満足するかどうかを *check* し、満足しておればそのルールを適用して構造を交換し、次の *state* へ *control* を移す。

PLATON では、上述の基本的な機能の他に、

push-down automaton 的な機能が付加されている。即ち、*branch* についての構造パターンと処理対象の構造との *pattern matching* によって切り出された特定の部分構造だけを、他の *state* へ移って処理をする。(*push-down*) そして、その部分の処理が終わると、*push-down* の生じた *branch* に戻って処理を続けてゆく。(*pop-up*) この能力によって、特定の部分を切り出して、その部分の処理をし、その処理結果をより大きな構造の一部として埋め込むことが出来る。



< 図-5 > 状態遷移と *push-down*, *pop-up*

branch には通常の *state-transition* を示す *branch* の他に、処理の終了を示す *branch*、処理が失敗したことを示す *branch* がある。これを整理して示すと下記のようになる。

- (i) NEXT-branch : 指示された *state* に状態を移して処理を進める。
(*state-transition* のための *branch*)
- (ii) POP-branch : その *level* の処理が終了して、指定された構造を持って *push-down* の生じた枝に戻る。もし、現在が最高 *level* の処理であれば、全体の処理が終了したことになる。
- (iii) FM-branch : その *level* の処理が失敗したことを示す *branch*。指定された *error-message* を持って *push-down* のおこった枝に戻る。

ここで(iii)の branch は Woods のモデルにはないものである。transition network のモデルは、状態を移すことにより使用する規則を制限するという意味で top-down 的な手法である。また、特に PLATON では、pattern matching により特定の部分構造を切り出し、その部分に特定の構造があるべきだと仮定して、その処理を行う。したがって、その仮定が誤った場合の処理が問題となる。

そこで、PLATON では仮定が何故不合理であったかを(iii)の branch が error-message として上位 level に返し、上位の処理では error-message を解釈して、適当な state へ処理の control を移すことが出来るようになってくる。

下に PLATON プログラムの形式的な定義を与える。

[PLATON プログラムの形式的定義]

```

< PLATON program > ::= ( < states > )
< states > ::= < state > | < state > < states >
< state > ::= ( < state name > < rules > )
< state name > ::= 任意の LISP 原子記号
< rules > ::= < rule > | < rule > < rules >
< rule > ::= ( < pcon > < strx > < con >
              < trans > ) < act > < end > )
< trans > ::= < transit > | < transit > < trans >
< transit > ::= ( ( < state name > < structure >
                  < vari > ) < errors > )
< vari > ::= < local variable > | < register name >
< local variable > ::= #N | #N1 | #N2 ·····
                   #P1 | #P11 | #P2 ·····
                   #B1 | #B11 | #B2 ·····
                   #K1 | #K11 | #K2 ·····
                   #I1 | #I11 | #I2 ·····
                   #J1 | #J11 | #J2 ·····
< register name > ::= / 任意の LISP atom
< errors > ::= | < error > < errors >
< error > ::= ( < error-message > < act >
               < pro > )
< error-message > ::= 任意の LISP 原子記号
< pro > ::= ( EXEC ( < trans > ) ) | ( TRANS
                < state name > )

```

```

< end > ::= ( NEXT < state name > < strx > )
           | ( POP < strx > )
           | ( FIM < error-message > )
< act > ::= ( SETR < register name > < structure > )
           | ( SENDU < register name > < structure > )
           | ( SENDD < register name > < structure > )
           | ( GETR < register name > )
           | ( GETU < register name > )
           | ( GETD < register name > )
< strx > ::= < structure 2 > | /
< strx > ::= < structure 1 > | /
< pcon > } ::= LISP 論理関数 & ユーザー定義
< con > |   の LISP 関数

```

(注) < structure 1 > は 3 が定義した変数名を、2 が定義した < structure > の中で使用を許した構造。 < structure 2 > は変数名以外に、register 名を使用した構造をいう。

上記の定義から解るように、< rule > は、< pcon > < strx >、< con >、< act >、< trans >、< end > からなる 6 字組によって表現される。以下の < 4-2 >、< 4-3 >、< 4-4 > においては、これら各 part の意味と使用法について説明する。

< 4-2 > < pcon >、< con >、< act > について
 前述したように、柔軟な処理(例えば syntax と semantics の相互参照的処理)を行うためには、ルールの中に user が自由に記述できるプログラムの部分が必要である。これを保障するのが < pcon >、< con >、< act > の各 part である。

ルール適用の結果をすべて構造の中に埋め込んで処理を進めてゆくのは、処理やルールをいわずに複雑にするだけであろう。そこで処理の中間結果をレジスタに保持しておいて、他のルールからそのレジスタを参照することが考えられる。即ち、ルール間、state 間の情報交換をすべて処理対象の tree-list 構造を遍いておこなわず、レジスタによっておこなうわけがある。レジスタのセットは主として < act > が行われる。

現在用意されているレジスタ用の < act > は、

<4-1>の定義で述べた *SETR*, *SENDU*, *SENDL*, *GETR*, *GETU*, *GETD* である。

処理途中で *push-down* がおこると、その時点での処理の流れから飛び出し、全く別の *state* から部分構造の処理が始まる。この部分構造の処理は全体の *control* の流れからは *level* が1段下の処理であると考えられる。このように *push-down* がおこると、*control* が1段下の処理に移り、*pop-up* がおこると、1段上位の *level* に *control* が戻ってくる。 *SETR* は現在処理中の *level* のレジスタをセットし、*SENDU* は1段上の、*SENDL* は1段下のレジスタに値をセットする関数である。同様に、*GETR*, *GETU*, *GETD* の各関数は、それぞれ現在の処理 *level* と同じ *level* のレジスタ、一段上位の *level* のレジスタ、一段下の *level* のレジスタの値を引出すための関数である。

この様に、レジスタには処理の *level* が付いており、*<register name>* が同じであっても処理 *level* の異なるレジスタは別の物であるとみなされる。このレジスタを使用することにより、*state* 間の情報交換、上下 *level* 間の情報交換が円滑におこなわれる。

<pcon>, *<con>* としては、PLATONのものほなにも関数を用意していない。特に *<con>* は *user* が意味モデルとして、どのようなモデルを採用するかによって、それに適した関数を、LISP 関数として定義する必要がある。5. において、PLATON を使用して書かれた解析プログラムの具体的な例として、'名詞句の解析'を行うプログラムを説明する。ここでは、*<con>* に使用される幾つかの LISP 関数が説明されることになる。これらの関数は、'名詞句の解析'に際して、採用した意味モデルを扱うものであって、PLATON 固有の関数ではない。

<4-3> *<strx>*, *<stry>* について

<strx> は *branch* についてのルール、書換ルールに対応していえば、左辺にあたる構造パターンを規定する部分である。*<strx>* には *<structure 1>* か / が書かれる。*<structure 1>* は <4-1> で述べたように、変数を含んだ構

造パターンである。PLATON システムは、この *<strx>* に書かれた構造パターンと、現在処理中の構造との *pattern-matching* をおこない、構造パターン *<strx>* 中の変数が、処理中の構造のどの部分に対応するかを知って、変数の結合を行う。/ は、この *pattern-matching* をおこなめないことを示す記号である。

<stry> は書換ルールの右辺にあたる構造パターンを規定する部分である。*<stry>* には *<structure 2>* か / が書かれる。*<structure 2>* は、*<structure>* の適当な部分を変数名あるいはレジスタ名にした構造パターンである。例えば、(** A #K /REG*) と書くと、*#K* の部分には、*<strx>* との *pattern-matching* の際に変数 *#K* と *binding* された部分構造が代入される。*/REG* の部分には、現在の処理 *level* のレジスタ名が *REG* であるレジスタの内容が、代入される。*<end>* part が (*NEXT <state-name> <strx>*) であれば、こうして代入された構造を持って、指定された *state* に *control* が移る。*(pop <stry>)* であれば、*control* は1段上の *push-down* の生じた *branch* に戻るが、その際 *<stry>* にまつて指定された構造を持って戻る。このように *<structure 2>* の中に、(*/<register-name>*) の表現を許すことによって、中間の処理結果をレジスタに保存し、構造の任意の場所にそれを埋込むことが可能になる。

<end> part が (*FM <error-message>*) の場合は、指定の *<error-message>* を持って *push-down* の生じた枝へ戻る。

<4-4> *<trans>* について

<trans> 部分は切出した部分構造を全く別の *state* から処理を始めるように指示する部分である。例えば、*<trans>* 部分が

((*(S1 #K #K)*)(*(S2 (* #I #J) /REG)*))

であるとする。ここで *#K*, *#I*, *#J* は *<strx>* 中に使用された変数である。*control* がこの部

分にくると、*push-down* がおこり #K に対応する部分を *state S1* から始まる処理に、(* #I #J) に対応する構造を *state S2* から始まる処理にそれぞれ手渡すことになる。そして、それぞれの処理が *pop-up* によって終わると、その処理結果をそれぞれ #K, レジスタ REG に入れる。

このように *push-down*, *pop-up* をおこなうことにより部分構造の処理をするのであるが、この部分構造の処理が失敗した時には *pop-up* ではなく、*FM-branch* によって *error-message* が上位の *branch* に戻ってくる。そこで <trans> の部分に、この *error-message* を解釈し、適当な処理に *control* の流れを変えるように指示することが出来る。例えば <trans> 部分が、

```
( ((S1 #K #K) (ERR1 (EXEC ((S5 #K #K)
                          ((S6 (* #I #J) /REQ))))
  (ERR2 (TRANS S8))
  (ERR3 (EXEC ((S7 #K #K))))
  ((S2 (* #I #J) /REQ)) )
```

であるとする。#K は *push-down* が生じて、*state S1* が処理される。この処理が戻ってくる戻り方は3種類ある。

- (1) *Pop-Up* に戻る。(*Pop-branch* に戻る)
- (2) *error-message* に戻る。(*FM-branch* に戻る。)
- (3) *NIL* に戻る。(適用すべきルールがなくなり、処理が行き詰って戻る。)

push-down を生じた上位の *branch* では、この(1)(2)(3)に依り次のような処理をする。

- (1) の場合； #K の処理は成功し、したがって次の処理 (*S2* で (* #I #J) の処理) に進む。
- (2) の場合； 帰ってきた *error-message* の解釈をする。すなわち、*error-message* が *ERR1* があれば、*S5* で #K, *S6* で (* #I #J) の処理をする。又 *error-message* が *ERR2* があればこのルール適用全体を放棄して、全体の処理を *state S8* からや

り直す。また、対応する *error-message* がない (*ERR1*, *ERR2*, *ERR3* 以外の *error-message* が戻ってきた) 時は、このルールの適用は失敗したことになる。

(3) の場合； このルールの適用は失敗したことになる。

なお左欄の例では、*S2* で (* #I #J) の処理をおこなうが、この処理には *error-message* 解釈部分がないので、(2) の戻り方をしても対応する処理がない。したがってこの場合は *NIL* が戻ってきたのと同じことになる。

control がある *state* の処理に入ってくると、その *state* から出てくる *branch* のルールが順々に適用可能かどうか *check* される。適用可能なルールがあれば、そのルールを適用し、そのルールの指定する *state* に *control* が移る。*NEXT-branch* であれば、指定された *state* に *control* が移り処理の *level* は変わらない。*Pop-branch* であれば、*push-down* を生じた *branch* へ *control* が戻り処理 *level* は一段上がる。*FM-branch* の場合は指定された *error-message* を持つ *push-down* の生じた *branch* へ *control* が戻り、この場合も処理 *level* は一段上がる。

その *state* から出てゆくルールが、すべて現在の処理対象に対して適用不可能な場合には、この処理 *level* での処理は失敗したことになる。NIL を持つ *push-down* の生じた *branch* へ戻る。現在の処理が最高 *level* での処理の場合には処理全体が失敗したことになる。

この <trans> の部分は、例えば左欄の例では、#K と (* #I #J) の処理をそれぞれ *S1*, *S2* から始まる処理に手渡すというように、*user* にとっては *parallel* に処理をしていると考えてプログラム出来る。もちろん実際には #K の処理を終了してから (* #I #J) の処理に進むわけだが、このことを考慮に入れてプログラムすることも可能である。

5. 具体的な解析プログラム

現在、我々は中学校理科の教科書を対象として文章分析のプログラムを作成中であるが、この章では名詞句の処理方式とそれが PLATON においてどのようにルール化されるかを説明する。

PLATON が書かれる言語処理プログラムには基本的には2つの part がある。

- (1) <strx>, <stry> が表現される構造変換をおこなう part ; これは $X \rightarrow Y$ という通常の書換ルールを表現する part が主として文章の syntax 処理に関係する。
- (2) <con>, <act> etc が表現される意味 check の part ; ここには user が採用する意味論モデルに対応する check が行われる。

従来の言語分析プログラムにおいては、文章の持つ意味 (semantics) と文章の構造 (syntax) との相互依存性が必ずしも十分にとり入れられてはいなかった。これには意味が計算機にかかるほど精密に定義されにくいという理由の他に、文章構造の解析途中の適切な個所で意味処理をおこなう (syntax 処理と semantic 処理の) *interactivet* な処理プログラムを書くことが非常に困難であるという理由が考えられる。このためプログラムは全体の syntax 処理を終了した後、semantics 処理に移るといったリアルな処理がおこなわれた。しかしながら PLATON においては、

- (1) 上記の (1), (2) の part が分離しているため意味分析のプログラムは独立してプログラム出来る。

- (2) PLATON の1つのルールの適用は上記 (1), (2) の処理結果によって行なわれること

から、syntax 処理と semantics 処理を独立してプログラム出来、かつルール適用に際してはそれらの処理が同時に行なえるという利点を持っている。

<5-1> 意味的処理

現在までの文章分析においては「名詞十の名詞型の句の処理は困難な問題であり、ほとんど手がつけられていない」といってよい。しかしながら分野を限って、その分野に関する知識を計算機に与えておいてやれば、これらの処理もある程度可能になると思われる。辞書の中には、動詞がどのような格 (Case) 関係を支配するかについて記述するだけでなく、名詞についてもかなり詳しく記述しておく必要がある。

「形容詞十名詞」型の構文も「(名詞 (MODIFY 形容詞))」の形式に整理するだけでは十分な意味解析とはいえない。また、「形容詞十名詞十の名詞」型構文において形容詞がどちらの名詞を形容しているのかを解析するには、名詞と形容詞の関係を辞書内に規定しておく必要がある。

我々はこういった名詞を中心とした構文を解析するために、名詞-名詞間、形容詞-名詞間の関係を記述した辞書を作成し、その辞書を検索する基本的な意味処理関数を用意した。そしてこの基本的な関数を PLATON のルールの <con> 部分に使用することにより、辞書で示された意味モデルと syntax 処理とを結合した解析プログラムを作成した。

(A) 辞書の構成

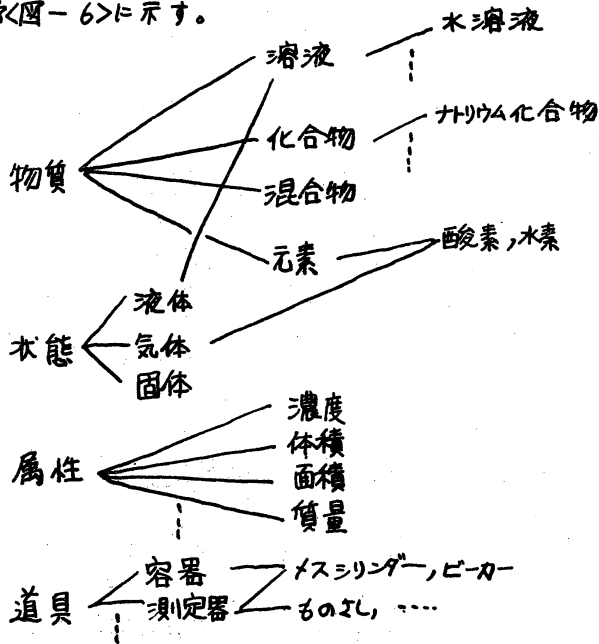
我々は日本語文の処理は「動詞中心の処理」が有効であると考え、教科書中の動詞について各動詞の格関係 (Case-Frame) を記述した辞書を構成し、それを利用した文章分析プログラムを PLATON によって作成中であるが、ここでは名詞句の処理を中心にして、名詞の辞書構成について述べる。

名詞辞書の entry には次の名辞素が記述される。

- (1) 上位概念
- (2) 下位概念
- (3) 属性 (4) 部分
- (5) 直接かかり得る形容詞
- (6) 他の名詞との結びつき。[(1)(2)(3)(4)以外の]

(1), (2) は従来の名詞の category 分類と同一役割

を果すものであるが、この分類は階層的な構造を持っている。教科書を対象として分類の一部を図-6に示す。



<図-6> 名詞の階層構造

‘物質の質量’、‘物質の体積’、‘物質の温度’
 或中の質量、体積、温度の名語は‘物質’につ
 き得る属性だと考えられる。したがってこれ
 らの名詞は‘物質’の辞書entry中の(3)の項
 目に入れられる。又‘人間の手’、‘人間の足’
 或中の手、足は‘人間’の部分を示す名詞であ
 り(4)の項目に入れられる。この様に1個の名
 詞は意味分節という見方から上位概念、下位概
 念の名語と関連し、その名詞の持つ属性、部分
 という見方から(3)(4)の各語と関係する。

例えば、

- 重い物質 → (質量が重い) 物質
- あたたかい物質 → (温度があたたかい) 物質
- 大きな物質 → (体積が大きな) 物質

のように、‘重い’、‘あたたかい’、‘大きな’の各形容
 詞は‘物質’を直接形容しているのではなく、‘物
 質’の‘質量’、‘温度’、‘体積’という各属性をそれぞれ
 形容していると考えられる。したがってこれらの
 形容詞は、‘物質’の辞書エントリの(5)ではなく、
 それぞれ‘質量’、‘温度’、‘体積’の辞書エントリ

の(5)に入っている。

次に、具体的な辞書entryの一例を示す。

<例>

液体(M)

superset	物質, 状態
subset	溶液, 水, アルカリ
属性	(体積 MODT) (質量 MODH) (温度 MODO)

溶液(M)

superset	液体, 混合物
subset	水溶液, ...
属性	(濃度 MODN)
部分	(溶質 物質) (溶媒 液体)

体積(M)

superset	属性
subset	
直接か 形容句	MOD T

MODT(K)

superset	MOD
subset	大きい, 小さい 数+cc, 数+cm ³

濃度(M)

superset	属性
subset	
直接か 形容句	MODN

MODN(K)

superset	MOD
subset	濃い, うすい 数+%

(*) Mは名詞, Kは形容詞を示す。この記法は以下でも使用する。

(B)意味処理関数

上記の辞書を検索して、名詞間の意味関係をcheckする関数として次の様な関数が用意されている。

(1)ATR; M1 ∈ ATR[M2]の関係をcheckする。
 即ち、名詞M1が名詞M2の属性語かどうかをcheckする。

例えば、M1:=濃度, M2:=溶液の場合には、この関係が成立する。

(2)PW; M1 ∈ PW[M2]のcheck。名詞M1が名詞M2の部分を示す名詞かどうかをcheckする。

(3)ARV; M1 ∈ ARV[M2]のcheck。例えば、M2:=溶液の時、溶液の属性としては、濃度、体積、温度などが考えられる。そして、属性‘濃度’の値としては‘濃い’、‘うすい’、‘80%’、属性‘体積’の値としては‘大きい’、‘小さい’ etc がある。このように名詞M1が名詞M2の属性の値に

になつていようかどうかを check する関数が ATRV である。

- (4) PWV: $M1 \in PWV[M2]$ を check する。木は、溶液の「部分」である溶質、溶媒のうち、溶媒の「値」となり得るから
 $木 \in PWV[溶液]$ が成立する。

以上4つの関数を用いて名詞句処理のプログラムが書かれる。

<5-2> 処理の考え方

名詞を中心とした構文を処理する方式として、「M+の+M」型の処理と、「(M+) + K + M」型の処理の2つに分けて説明し、次にこれらの各処理を統合して全体の処理がどのように流れるかを説明する。

(A) 「M+の+M」型の処理

この型には構文の名詞の種類によって基本的に次の4種類の処理がある。

(1) 機能語の処理

(M1+の+M2)型の句において、M2にある特定の名詞がくることによつて、1つの宛ま、た意味の変換、付加をM2がM1に対しておこなう。

- (例) 水の粒の間、物質の中、物質の一部、化学変化の前後

これらの機能語の特長は、この語の前後が意味的なまとまりがあるということである。即ち、語の前の部分を独立に切り出して処理し、その処理結果に対してこれらの語が作用すると考え処理する。これは PLATON のパターンマッチング機能を用いれば簡単に実現できる。

(2) 名詞辞書を利用し、名詞間の関係が処理する場合

M1とM2の関係として次の様な場合が考えられる。

(i) $M2 \in ATRV[M1]$

- (例) 溶液の体積、物質の比熱

この場合は、M1とM2は結合されて(M2 (*ATR M1))の形に整理される。

(ii) $M2 \in PWV[M1]$

- (例) 溶液の溶質、化合物の構成元素

処理は(i)の場合に準いる。

(ii) $M1 \in ATRV[M2]$

- (例) 100ccの木、赤色の溶液

この場合M1とM2は直接つながっているのではなく、M2の属性(例えば体積、色)の値としてM1がつながる場合がある。処理結果は、それぞれ(木(体積 100cc))、(溶液(色 赤色))となる。

(iv) $M1 \in PWV[M2]$

- (例) 塩化ナトリウムの木溶液、銅の化合物

(ii)の場合に準いる。処理結果は(木溶液(溶質 塩化ナトリウム)) etc となる。

(3) M2がActive要素を持つ場合

名詞の中には「する」がついて動詞化する名詞がある。このような名詞はサ変名詞と呼ばれ、この名詞が単独で使用される時も、やはり目的格、主体、道具などの動詞的要素を支配することが多い。理科教科書中でのこれらの名詞としては、「分解」「変化」「検出」「合成」「実験」 etc がある。これらの名詞の辞書エントリには、その動詞化に対する格支配関係(Case-Frame)が用意されているので、これを利用して「の」の処理を行う。

(4) 時間、場所を示す名詞がきた場合

時間、場所の名詞がM1として使用された場合には、その時間、場所における名詞M2と解釈出来ると思われる。

(B) 「(M+) + K + M2」型の処理

「K + M」型の処理は比較的簡単である。即ち、前述辞書の中で、KがMのentry中の(b)「直接かかれる形容詞、名詞」の中に入っているか、又はMの属性、部分にあたる名詞の中でKによって直接形容される名詞があるかどうかを check する。例えば、

- (a) 濃い溶液, (b) 赤い色

例(b)ではKがMを直接形容出来るから処理結果は(色(SPEC 赤い))となる。例(a)においては、溶液の属性である濃度が、Kによって形容されるから処理結果は(溶液(濃度 濃い))となる。

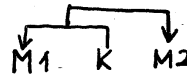
「 $M1+K+M2$ 」型の係受関係は $M1$ と $M2$, & K と $M1$, K と $M2$ の関係によって、大別して次の2種類がある。

(I) $M2 \in \text{ATR}[M1]$ 又は $M2 \in \text{PW}[M1]$ であり、かつ $K \in \text{MOD}[M2]$ である場合
係受関係は下図のようになる。



(例) 水の高い沸点, 銅板の赤い部分

(II) $M1 \in \text{ATR}[M2]$ 又は $M1 \in \text{PW}[M2]$ であり、かつ $K \in \text{MOD}[M1]$ である場合
係受関係は下図のようになる。



(例) 濃度の濃い溶液, 沸点の高い塩化ナトリウム
これらの処理結果は、それぞれ(溶液(濃度 濃い)), (塩化ナトリウム(沸点, 高い))となる。

もちろん、「 $M1+K+M2$ 」の並びであっても、 $M1$ と $M2$, K と $M2$ が独立にかかっているように3要素組($M1, K, M2$)の関係性を調わなくとも処理が出来る場合がある。

例えば、「固体の白い物質」の場合には、「固体 \in ATRV[物質]」「白い \in ATRV[物質]」となって「固体」「白い」がそれぞれ独立に物質に係ることになる。この様な並びについては次の(C)の全体の処理の流れによって処理可能になる。

(C)全体の処理のflowについて

(A), (B)において局所的な処理の進行について述べたが、次に全体として処理がどの様に流れるかを考える。実際の外々の名詞句においては、例えば、(1)「濃い濃度の食塩の溶液」($K, M1$ の $M2$ の $M3$) (2)「濃度の濃い食塩の溶液」($M1$ の $K, M2$ の $M3$) (3)「食塩の濃度の濃い溶液」($M1$ の $M2$ の $K, M3$) (4)「食塩の濃い濃度の溶液」($M1$ の $K, M2$ の $M3$) の様に、多数個の M と K からなり、その係受のパターンもその中の名詞によって変わる。

(例)

- (a) 濃度の濃い食塩の溶液

```

            graph TD
              M --- K
              M --- M1[M]
              M1 --- M2[M]
            
```
- (b) 食塩の濃い濃度の溶液

```

            graph TD
              M --- K
              M --- M1[M]
              M1 --- M2[M]
              M1 --- M3[M]
            
```

この様な多数個の名詞や形容詞を含む場合の処理は、基本的には(A)(B)が述べた処理の繰返しであるが、PLATONのpush-down 機能を用いて、柔軟に係受関係が調_合られる。例を示す。

<例1> 「濃度の濃い食塩の溶液」

(M, K, M)のpatternになっ_ているが、 $M1$ と $M2$ には(B)が述べたような関係が成立しない。したがって($M1, K$)の部分は処理を保留して、($M2$ の $M3$)「食塩の溶液」の部分をpush-downして、処理する。この部分の処理結果は(溶液(溶液 食塩))となっ_てpop-upしてくる。したがって全体の句の様子は

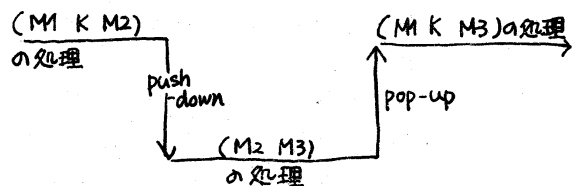
「濃度の濃い(溶液(溶液 食塩))」

となる。これに対して(M, K, M)のpatternをあてはめて、(B)が述べた関係が成立するかどうかをcheckする。すると「濃度 \in ATR[溶液]」「濃い \in MOD[濃度]」となっ_て条件が成立する。したがって全体の処理が完了して、処理結果は、

(溶液(溶液 食塩)(濃度 濃い))

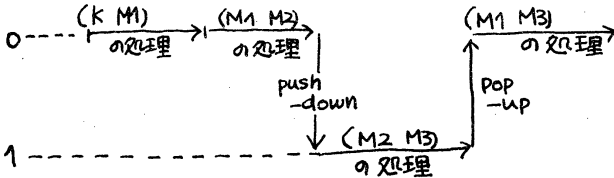
となる。この様に、一度(A)(B)が述べた処理をおこなない、その処理が失敗すると、push-downをおこなして処理を一段進め、その処理結果を埋込構造に対して、もう一度(A)(B)の処理をおこなうことを処理の基本的なflowとしている。なお、push-downの際のstateは全体の処理の出発点となるstateと同じであり、したがってpush-downされた処理も最レベルでの処理も全く同じ過程である。

この例の場合のflowを図示すると、下の様になる。



<例2> 「濃^Kい濃^{M1}度の食塩^{M2}の溶液^{M3}」

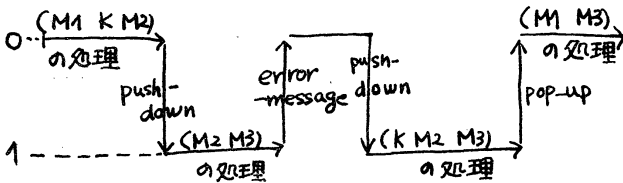
<例1>の場合と同じように図示すると下図のようになる。



<例3> 「食塩^{M1}の濃^Kい濃^{M2}度の溶液^{M3}」

この場合には、push-downした後の処理が失敗し、その失敗原因をerror-messageとして返すことにより、処理の流れがcontrolされる。すなわち「濃度の溶液」という句は解析不可能が失敗する。これは「濃度 ← ATR [溶液]」であるが、日本語においては、属性語(この場合では「濃度」)が、本体の語より前に単独で現われることがないからである。そして「濃い濃度の溶液」ならば解析可能である。

この場合の処理の流れは下図の様になる。



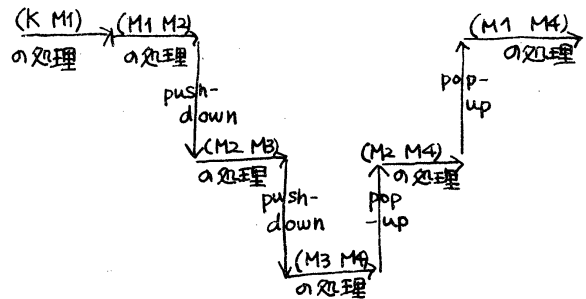
これら<例1><例2><例3>の処理結果はいずれも

(溶液 (溶質 食塩) (濃度 濃い))

となり、同じである。また以上の例ではpush-downが1段であるが、プログラムとしては、push-downしたlevel 1の処理も、level 0の処理と全く同じ処理をするので、level 1の処理途中で再びpush-downがはいてlevel 2の処理へ移ることも当然可能である。この例としては、

<例4> 「濃^Kい濃^{M1}度の食塩^{M2}の100cc^{M3}の溶液^{M4}」

この例の処理のflowは下図のようになる。



処理結果は、

(溶液 (溶質 食塩) (濃度 濃い) (体積 100cc))

となる。

自然言語処理のためのプログラミング言語を作ることの意味は次のようなところにある。(1)人間が頭で考えた処理方法がなるべく素直に表現できること、(2)試行錯誤的に行なうために書き替えが容易であること、(3)他人が見てその内容がわかりやすいこと、(4)アルゴリズムの記述能力がすぐれていること

我々の作ったPLATONはこのような条件を十分に満たしているとは言えないまでも、かなり便利なシステムになっているものと考えている。現在これを用いて日本語分析のプログラムを書いているが、システムとして特に大きな問題はない。このシステムは京都大学大型計算機センターのFACOM 会話型LISPの上で動いているが速度は会話型として許される範囲内にある。

文献

長尾真：自然言語処理のためのソフトウェア
昭和48年電気四学会論文集 No.348
W.Woods: Transition Network Grammar for Natural Language Analysis, CACM Vol.13, No.10, 1970