

新しいパーザ PAMPSと そのコンパイラについて

上原邦昭、落谷亮、河合和久、豊田順一
大阪大学基礎工学部

1. はじめに

自然言語(特に英文)を解析するパーザとして、我々は PAMPS (PAttern Matching Parsing System)を開発してきた。文献(1)、(2) このパーザは解析制御のために pattern matching を用いているところに特徴がある。pattern matching (unification)は、LISPでの data selector (car, cdr)、data constructor (cons)等を単一の機能として実現できるため、記述形式に一貫性があり、更には side effect がないため、明晰性が保てるという利点がある。またパーズングの制御に必要な機能の大部分は統語/意味情報の転送と判定であるが、これらも pattern matching 機能に統一できるという利点がある。

またパーズングを効率よく行なうためには、Marcus も指摘しているように、入力データの利用 (bottom-up)、予測の利用 (top-down)、先読みデータの利用 (look-ahead)が必要である。文献(3) PAMPS では、Earley のアルゴリズムを上述の目的に合うように変形したものを採用している。Earley のアルゴリズムでの Predictor に相当する部分は、LINGOLでの oracleと同様に文法規則を前処理して、ビットテーブルの形で保存している。bottom-up 過程で成長した部分構文木は、このテーブルにより filter にかかけられ、将来成長する可能性のあるもののみが残される。Earleyのアルゴリズムで用いられている先読み機能は、不用な部分構文木の生成を抑制することしかできない。我々のアルゴリズムでは更に“不用な予測”をも先読み機能により削除することにより、パーズングの探索空間/時間を縮小している。ここで“不用な予測”とは、予め文法規則で予測されていた非終端記号のうちで、先読み語を最左端とする終端記号列として導出できないものをいう。言い換えると、予測非終端記号を設定するとき、その予測非終端記号が次の先読み語の親となりうるかどうか調べ、もし親となりえないならば、その予測非終端記号は構文木を組み立てるうえで不用になるということである。この先読み機能についても予め先読み語を計算しておくことができる。結果は oracleと同様に、ビットテーブルの形で保存している。

2. PAMPS のアルゴリズム

まず、Earleyのアルゴリズムでの項 (item) に対応してゴールとフレーズを定義する。項 $[A \leftarrow \alpha \cdot B\beta, i, j]$ をゴール $(B\beta, i, A, j)$ と呼ぶ。項 $[A \leftarrow \alpha \cdot, i, j]$ をフレーズ (A, i, j) と呼ぶ。 $A \xrightarrow{*} w \dots$ という導出が存在するとき、即ち A から導出される任意の文形式のうち w が最左端に現われるとき、関係 $R^*(A, w)$ が成り立つ。 $S \xrightarrow{*} \dots A w \dots$ という導出が存在するとき、即ち w が A の直接右隣りに並ぶとき、関係 $Follow(A, w)$ が成り立つという。

パーズングは左から右へ1語ずつ進行し、部分構文木は bottom-up に作られる。フレーズ (B, i, j) が生成された段階で、以下の4種類の動作が実行される。

- 1) try C: ゴール($BC\alpha, i, A, R$) ($\alpha = \epsilon$ でもよい) がある場合、先読み語 w_{j+1} について $R^*(C, w_{j+1})$ ならば、ゴール($C\alpha, j, A, R$) を設定する。
- 2) try R: ゴール(B, i, A, R) がある場合、先読み語 w_{j+1} について $Follow(A, w_{j+1})$ ならば、フレーズ(A, R, j) を生成する。
- 3) try D: ゴール($C\alpha, i, E, R$) かつ $R^*(C, A)$ となる規則 $A \leftarrow B$ がある場合、先読み語 w_{j+1} について $Follow(A, w_{j+1})$ ならば、フレーズ(A, i, j) を生成する。
- 4) try L: ゴール($C\alpha, i, E, R$) かつ $R^*(C, A)$ となる規則 $A \leftarrow BD\alpha$ がある場合、先読み語 w_{j+1} について $R^*(D, w_{j+1})$ ならば、ゴール($D\alpha, j, A, i$) を設定する。

以上の4種類の動作を繰り返し、もはや実行すべきことがなくなったならば、更に1語読み込み(w_{j+1})動作を続ける。新たに先読み語は w_{j+2} となる。
 説明の都合上、実際のアルゴリズムと若干異なるので、詳しくは文献(2)を参照のこと。以上のアルゴリズムをEarleyの文献[4]で示された例題に適用した場合以下のような効果があることがわかる。

表1. 先読みの効果

	文法	AE	UBDA	BK
	入力文	$a+axa$	x^4	x^4
Earley 先読みなし	フレーズ生成数	14	18	24
	ゴール設定数	8	10	4
Earley (1語先読み)	フレーズ生成数	11	15	21
	ゴール設定数	7	10	4
PAMPS (1語先読み)	フレーズ生成数	11	15	21
	ゴール設定数	4	6	3

3. 文法と辞書の記述

文法規則は、各非終端記号に引数を与えた拡張CFG rule である。引数として定数、変数、更に変数と定数からなる複合項(データ構造)が許される。

```
SENTENCE(ques) <-
  auxil(*atype,**,**,*nums),np(*nums),vp(*atype,**,**).
```

辞書項目は、終端記号、非終端記号、特徴素リストからなる。特徴素リストは特徴素タイプ、特徴素のペアからなる。

```
(they n(npos,ns(x,o,x)),them)
((ntype pro) (pers third)))
```

特徴素リストは、システム組み込み関数 $F-get$ (終端記号、非終端記号、特徴素タイプ、特徴素が代入される変数)によりアクセスされる。

4. 解析制御

自然言語を解析する際、CFGでは能力が弱いので、何らかの補強が必要となる。このため、大別して1) 部分構文木の上位レベルノードに、あるいは同レベルの

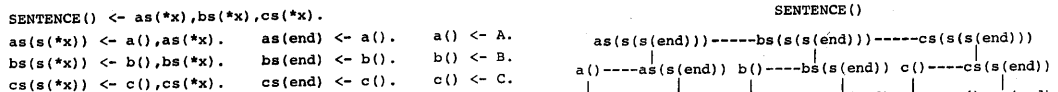
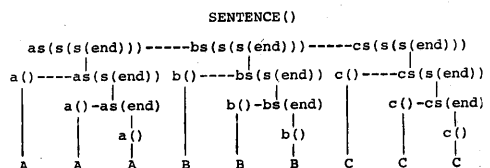


図1. $A^n B^m C^n (n \geq 1)$ の導出木



他ノードに統語/意味情報を送ること、2) 送られた情報を判定してその規則の適否を決めることが必要である。これらの機能は *pattern matching* として単一化することが可能である。図7に言語 $L(G) = \{a^n b^n c^n | n \geq 1\}$ を導出するための規則と、その導出例を示す。各非終端記号に付け加えられた引数により、情報の転送、判定がスマートに表現できる。実際の言語解析で *pattern matching* をどのように用いるかについては、文献(2)でいくつかの例が示されている。

5. Extra Condition⁽⁵⁾

文法記述においては、文法規則数、カテゴリ数の増加を防ぐため、規則内にその適用条件を埋め込み、条件の判定により規則の適否を決定することが必要である。PAMPSでは、この規則適用条件記述を *procedure call* として規則内に含むことを許している。

```
head(*ntype,*nums,third,*word) <-
  n(pos,*nums,*word),F-get(*word,ntype,*ntype), in(*ntype,[comm,prop] .
  in(*x,[*x *1]).
  in(*x,[** *1]) -> in(*x,*1).
```

上の文法規則では、中カッコでくくられている部分が規則適用条件となっている。中カッコの中はホーン節集合である。この記法を文献(6)の定義に従い、*local* ホーン構成と呼ぶ。中カッコの中が *local* に証明されれば、適用条件を満足したとして、更にパーシングは進行する。尚中カッコの中では *top-down*、*depth-first* で証明が行なわれる。

6. PAMPSコンパイラ

PAMPSのような述語論理的文法記述の場合、宣言的意味論と手続き的意味論の両方を持つことが従来より指適されている。文献(5) 文法規則が宣言的意味論を持つと見なした場合、文法規則と共にそれを解釈実行するためのメカニズムが必要となる。前章までの考察は主としてこのメカニズムをどのように構成するかについての議論であった。本章以降では、文法規則を手続き的意味論に基づき、文法規則自身を低レベル機械語に翻訳し、それを直接実行するようなメカニズム、即ちPAMPSコンパイラと仮想機械PAMPSマシンについて言及する。また、2章のアルゴリズムによって文法規則を解釈するメカニズムをPAMPSインタプリタと呼ぶ。

6.1. PAMPSコンパイラの場合 文法規則をインタプリタにより解釈実行すれば、(1) 適用可能な文法規則を規則集合の中から逐一検索しなければならず、(2) *pattern matching* も統一的なアルゴリズムで実行しなければならない。

まず(1)の問題について考える。次の文法は、入力記号列によって駆動されるオートマトンとしてとらえることができる。即ち、ある状態からスタートして、入力記号に応じた遷移によって別の状態に移る様子をラベル付き有向グラフとして表現できる。ここで、各ノードは状態に、各アーク間のラベルは状態遷移の際に必要なアクション(先読みのチェック、トップダウン過程の予測チェック、*pattern*

matchingの処理)を示している。

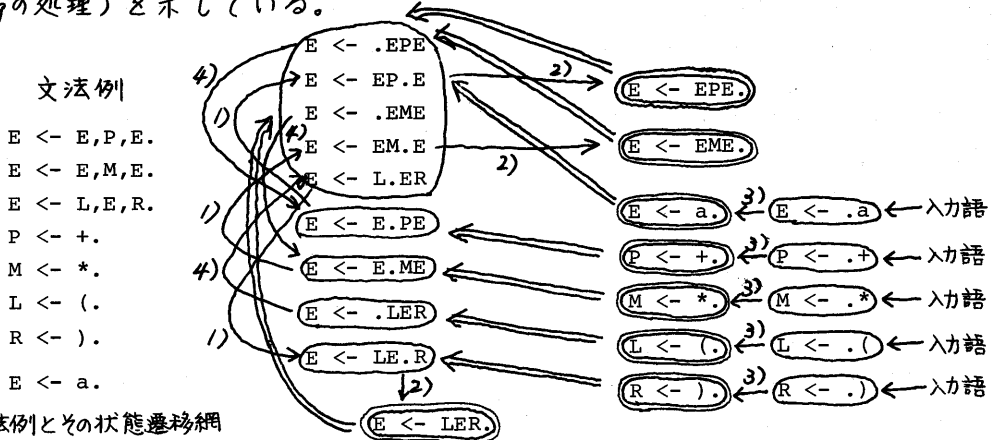


図2. 文法例とその状態遷移網

ラベルtryC等は、2章の4分類のアクションに対応している。状態2重丸ではフレーズが生成され、そのフレーズはqueueに入れられる。queueの先頭要素から順次、状態遷移(手続き呼び出し)が行なわれる。queueの要素が無くなれば、入力を1語読み込む。

上述のように文法入力時に、あらかじめ道筋に関する情報を取り出しておき、実行時には、この道筋に従ってパーズングすることで、無駄な規則検索の手間を省くことができる。

パーズング制御に使用されるpattern matchingは単純な代入と同値判定に還元できる。

```

Sentence() ← Noun-phrase(*x), Vreb-phrase(*x).
Noun-phrase(Single)

```

この場合、ゴールNoun-phraseとフレーズNoun-phraseとのpattern matchingは、 $x := \text{Single}$ という代入に帰着する。

```

Sentence() ← Noun-phrase(Single), Verb-phrase(Single).
Verb-phrase(Plaural)

```

この場合に引数どおしでの必要な操作は、Single \rightleftharpoons pluralを判定することである。このように非終端記号の各引数を、それぞれの出現状況に応じていくつかの場合に分け、それぞれに必要な操作を低レベル語に翻訳し、処理を高速化する。また、PAMPSはアルゴリズム上bottom upに動作するので、手続き処理用コード(pattern matching用コード)にコンパイルする部分は、文法規則の右辺である。左辺は、手続き呼び出し用コードにコンパイルする。以上がPAMPSコンパイルの原理である。

PAMPSマシンはFortranにより実現された仮想機械である。

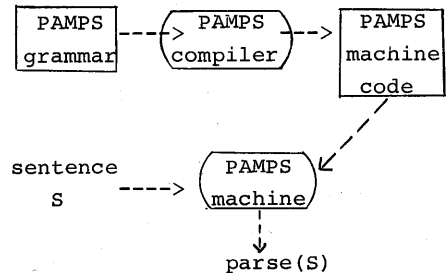


図3. PAMPSマシン概念図

6.2. PAMPSマシンコード 文法規則 $A \leftarrow E, P, T$ は図4のように翻訳される。pattern matching用コードは以下の4グループに分けられる。

- 1) 変数が、その文法規則内で初めて出現する場合は、相手引数の値を代入する。
 - 2) 変数が、文法規則内で一度しか出現しない場合は、いかなる処理も必要としない。
 - 3) 定数あるいは複合項の場合は、相手引数を調べ、同値の場合のみ成功、もし相手変数が具体化していない場合は、相手に代入する。
 - 4) 変数が既に文法規則内で出現している場合は、自分自身に代入された値を調べ、相手側の引数と pattern matching する。
- 更に、パーズングコードとして、前述の try C 等の他、手続き呼び出し用コード、手続き呼び出し制御用コードなどがある。PAMPSマシン用コードは、これら約50種類のコードから構成されている。

pattern matching (E)	} フレーズEとの pattern matching 用コード
neck	
gset	} 以下で必要な変数セルを確保
pattern matching (P)	
gset	} ゴールをメモし、手続き呼び出し符となる
pattern matching (T)	
pset	} フレーズをメモする。
call (A)	

図4. 文法規則 A ← EPT の模式コード列.

7. PAMPS マシンのデータ構造

PAMPS マシンのデータ構造は、Boyer-Moore の sharing structure 技法に基づいている。しかし、PAMPS は本質的に bottom-up、parallel で動作するため、Warren 等の Prolog で用いたデータ構造のように、1つの変数に1つの変数セルを与えるだけではうまく処理できない。たとえば、1つのフレーズが複数個のゴールと、あるいは1つのゴールが複数個のゴールと同時に pattern matching する可能性がある。この状態を競合 (conflict) と呼ぶ。我々は、各時点で行なわれた pattern matching に個々の番号を与え、更に各 pattern matching 間の祖先—子孫関係をテーブルとして表現することにより、競合問題を解決した。

図5について考える。フレーズ B(1) と規則 $A(x) \leftarrow B(y), C(x), D(x)$ が pattern matching した場合、ゴール C(x) が設定される。この場合、変数 x のセルには何も書かれていない。--- ①

ゴール C(x) とフレーズ C(2) の pattern matching が行なわれると、variable stack 上の変数 x のセルに、値 int(2) とこの時点での pattern matching 番号 2 が書き込まれる。競合はまだ起きていないので、conflict pointer は nil である。--- ②

フレーズ C(3) が生成されると、ゴール C(x) と pattern matching する。変数 x のセルには値 int(2) が書かれているが、並列性のため、この値をそのまま変数 x の値と

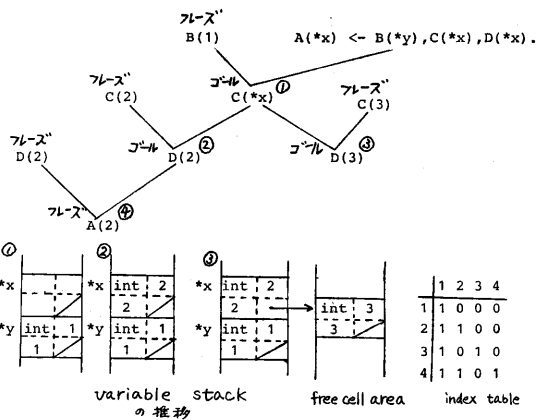


図5. データ構造例

みなすことはできない。そこで、この pattern matching での変数 α のセルとして、新しいセルが Free cell area から 1つ取られ、新しいセルに値 int (3) と pattern matching 番号 3 が書かれる。同時に variable stack 上にある変数 α のセルの conflict pointer は、Free cell area の新しいセルを指す。... ③

次にフレーズ D(2) が生成され、ゴール D(2) と pattern matching する。まず、ゴール側の引数 α の dereference が行なわれる。変数 α のセルには値 int (2) が書かれているが、これをそのまま変数 α の値と判断することはできない。(競合の可能性) その値が現在の pattern matching での真の値かどうかの判定のために、index table を用いる。index table はビットマトリックスとして構成され、i 行 j 列のビットは、pattern matching 番号 j の節点が、pattern matching 番号 i の祖先であるかどうかを示している。変数 α のセル上の値 int (2) の pattern matching 番号は 2 で、現在のゴールの pattern matching 番号は 2 であるから index table 2 行 2 列を調べる。2 行 2 列にはビットが立っているので、このセルは現在の pattern matching のゴールの祖先であることがわかる。これより α の値は int (2) となり、pattern matching は成功する。... ④

ボトムアップ、パラレル型パーサでは途中結果がすべて残されるが、自然言語入力文は高々数十語であるので、我々のデータ構造で十分対処可能である。

9. 計算機実現

PAMPS コンパイラと PAMPS マシンの有効性を確認するために、計算機実験を行なった。PAMPS インタプリタを LISP 1.9 でインプリメントし、それをコンパイルしたものと、PAMPS マシンを Fortran でインプリメントしたものを、同一文法と同一入力文を用いて速度比較した。その結果 PAMPS マシンの方が約 1.6 倍スピードアップできることがわかった。PAMPS マシンはプログラムの最適化をしていないので、まだいくらか高速化できると思われる。現在は、PAMPS コンパイラと PAMPS マシンの最適化と PAMPS マシン用の自然言語処理用文法を平行して作製している。文法規則は総数で約 150 個ほどある。

10. 参考文献

1. K. Uehara et al. "A Pattern Matching Directed Parser: PAMPS" ACL/LSA sessions on Computer Modeling of Linguistic Theory, NY (1981).
2. 河合 世. "パターンマッチングによるパーサ PAMPS について" 信学技報 AL81-94 (1982).
3. M. P. Marous "A Theory of Syntactic Recognition for Natural Language" M.I.T. Press (1980).
4. J. Earley "An Efficient Context-Free Parsing Algorithm" CACM Vol. 13, No. 2 (1970).
5. D. H. D. Warren "Implementing Prolog" D.A.I. Research Report No. 29, 40 (1977).
6. 渕 判 "述語論理的プログラミング—EPILOG の提案—" 情報処理記号処理研究会 1-2 (1977).