

TR-627

YAGLR法：Yet Another Generalized
LR Parser

田中 穂積、G. スレッシュ
(東京工業大学)

March, 1991

© 1991, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03)3456-3191~5
Telex ICOT J32964

Institute for New Generation Computer Technology

1990. 2. 6

Y A G L R 法: Yet Another Generalized LR Parser

田中穂積, K. G. スレッシュ

(東京工業大学・工学部)

概要

プログラム言語などの人工言語のコンパイラでは, Knuthの考案した L R (k) 法とよぶ統語解析アルゴリズム [Knuth 65] が使われることがある。これは先読み情報を用いて決定的に統語解析を進め、無駄なく効率的に統語解析を行うことができる。ただしそこで使われる文法は、L R 文法に限られており、自然言語の解析に用いる一般的な文脈自由文法 (CFG) を扱うことができない。

富田は L R (k) 法の持つ利点をそのまま保持し、しかも一般的な CFG が扱えるように L R (k) 法を拡張している [Tomita 86, 87]。これは一般化 L R 法 (generalized LR parsing algorithm; 拡張 L R 法) の 1 つであり、富田法とよばれている。経験的に富田法はアーリー法と比べて高速な統語解析が可能であるが [Tomita 86]。富田法に対して、アーリー法以上のオーダの解析時間を要す特殊な CFG の存在が知られている。

我々は、アーリー法の利点と富田法の利点をそのまま引き継ぐ新しい一般化 L R 法（これを Y A G L R 法とよぶ）を提案する。それによれば、先に述べた特殊な CFG に対する解析時間のオーダを n^3 に抑えることができる。また、富田法ではグラフ構造化スタック (GSS) を用いて、使用記憶空間の節約を図っているが、GSS の構造は必ずしも単純でなく、実装が容易ではない。これに対して Y A G L R 法では木構造化スタック (TRS S) を用いており、経験的に、富田法の GSS と比較して同程度の使用記憶空間に抑えることが可能であり、構造が単純であるため扱い易い。これは、Y A G L R 法では富田法以上にスタックのマージが可能になるためである。本論文では、Y A G L R 法による統語解析アルゴリズムの定義と動作例を説明する。また、富田法と比較して、Y A G L R 法はどの様な利点があるかを具体例を用いて説明する。

田中穂積, K. G. スレッシュ
東京工業大学・工学部

概要

プログラム言語などの人工言語のコンパイラでは, Knuthの考案したLR(k)法とよぶ統語解析アルゴリズム [Knuth 65] が使われることがある。これは先読み情報を用いて決定的に統語解析を進め、無駄なく効率的に統語解析を行うことができる。ただしそこで使われる文法は、LR文法に限られており、自然言語の解析に用いる一般の文脈自由文法 (CFG) を扱うことができない。

富田はLR(k)法の持つ利点をそのまま保持し、しかも一般のCFGが扱えるようにLR(k)法を拡張している [Tomita 86, 87]。これは一般化LR法 (generalized LR parsing algorithm; 拡張LR法) の1つであり、富田法とよばれている。経験的に富田法はアーリー法と比べて高速な統語解析が可能であるが [Tomita 86]、富田法に対して、アーリー法以上のオーダーの解析時間を要す特殊なCFGの存在が知られている。

我々は、アーリー法の利点と富田法の利点をそのまま引き継ぐ新しい一般化LR法（これをYAGLR法とよぶ）を提案する。それによれば、先に述べた特殊なCFGに対する解析時間のオーダーを n^3 に抑えることができる。また、富田法ではグラフ構造化スタック (GSS) を用いて、使用記憶空間の節約を図っているが、GSSの構造は必ずしも単純でなく、実装が容易ではない。これに対してYAGLR法では木構造化スタック (TRS) を用いており、経験的に、富田法のGSSと比較して同程度の使用記憶空間であり、構造が単純であるため扱い易い。これは、YAGLR法では富田法以上にスタックのマージが可能になるためである。本論文では、YAGLR法による統語解析アルゴリズムを与え、富田法との比較を具体例により説明する。

[Abstract]

We have developed a new generalized LR parsing algorithm called YAGLR in which it combines some of the advantages of both Earley's and Tomita's algorithms. We show that for some special CFGs the time consumed by YAGLR parsing algorithm is less than Tomita's parsing algorithm and is in the order of $O(n^3)$. While YAGLR uses a tree structured stack, the structure is not so complex compared to the one of Tomita's graph-structured stack. The reason is that YAGLR parsing algorithm enables us to make much more merge operations on the tree-structured stack.

1. はじめに

プログラム言語などの人工言語のコンパイラでは、Knuthの考案したLR(k)法とよぶ統語解析アルゴリズム [Knuth 65] が使われることがある。これは先読み情報を用いて決定的に統語解析を進め、無駄なく効率的に統語解析を行うことができる。ただしそこで使われる文法は、LR文法に限られており、自然言語の解析に用いる一般の文脈自由文法 (CFG) を扱うことができない。

富田はLR(k)法の持つ利点をそのまま保持し、しかも一般のCFGが扱えるようにLR(k)法を拡張している [Tomita 86, 87]。これは一般化LR法 (generalized LR parsing algorithm: 拡張LR法) の1つであり、富田法とよばれている。経験的に富田法はアーリー法と比べて高速な統語解析が可能であるが [Tomita 86]。富田法に対して、解析すべき文の長さをnとしたとき、 n^m ($m > 3$) のオーダーの解析時間を要す特殊なCFGの存在が知られている [Johnson 89] [Kipps 89]。一方、アーリー法は、いかなるCFGに対しても、解析時間のオーダーが n^3 であるから、この特殊なCFGに対して、富田法はアーリー法以上の解析時間を要すことになる。

本稿で提案する新しい一般化LR法では、アーリー法の利点と富田法の利点をそのまま引き継ぐ。この新しい一般化LR法をYAGLR法 (Yet Another Generalized LR parsing) とよんでいる。YAGLR法によれば、先に述べた特殊なCFGに対する解析時間のオーダーを n^3 に抑えることができる。また、富田法ではグラフ構造化スタック (GSS) を用いて、使用記憶空間の節約を図っているが、GSSの構造は必ずしも単純でなく、実装が容易ではない。これに対してYAGLR法では木構造化スタック (TRS) を用いるが、経験的に、富田法のGSSと比較して使用記憶空間は同程度であり、構造が単純であるため扱い易い。これは、YAGLR法では富田法以上にスタックのマージが可能になるためである。

YAGLR法では、解析中に統語解析木は作らず、統語解析木を構成する部品のみを作成する。この部品を逆ドット項と呼ぶが、これはアーリー法で作成するドット項 (アーリー項) と対称な構造をしている。2章では、逆ドット項を説明する。3章では、YAGLR法の鍵となるTRSのマージ操作 (統合化操作)、シフト操作、レデュース操作を説明し、最後にYAGLR法の統語解析アルゴリズムを説明する。レデュース操作の過程で逆ドット項を作成するが、逆ドット項を作成することの利点も3章で説明する。また富田法とYAGLR法による解析例を示し、富田法とYAGLR法の比較を行う。4章では、特殊なCFGに対しても、YAGLR法の統語解析時間のオーダーが n^3 であることを証明する。5章では、まとめと今後の研究課題について述べる。

2. 逆ドット項(Dot Reverse Item)

本章では、YAGLR法による統語解析で重要な役割を果たす逆ドット項を説明する。富田法がレデュース操作時に部分統語解析結果を圧縮統語森として作成するのに対して、YAGLR法ではアーリー項と対称な逆ドット項の集合を作成する。逆ドット項の定義を以下に与える。

- (1) 解析対象文wを $w = w_1 w_2 \dots w_n \in T^*$ であるとする。ここに T^* は、終端記号の集合Tの要素を0個以上並べたものを表す。終端記号 (語と考えてよい) w_i と w_{i+1} の間に番号iを振り、これを位置番号とよぶ (iを、語 w_i の位置番号とよぶ)。ただし w_i の左には位置番号0を、また w_n の右には位置番号nを振る。
- (2) CFG規則 $A \rightarrow \alpha$ があり、 α が $\beta \gamma$ と書けるとき、 $0 \leq j \leq n$ なるjに対して、 $[A \rightarrow \beta \cdot \gamma, j]$ は、wに対する逆ドット項である。特に、 $[A \rightarrow \alpha \cdot, j]$ または $[A \rightarrow \alpha \cdot \cdot, j]$ もwに対する逆ドット項である。
- (3) 逆ドット項の集合 I を次のように定義する。 $0 \leq i \leq j \leq n$ なるiとjに対して、

$$[A \rightarrow \alpha \cdot \beta, j] \in I \text{ iff } S \Rightarrow^* \alpha \beta \delta, \beta \Rightarrow^* w_{i+1} w_{i+2} \dots w_j, \text{ and } \delta \Rightarrow^* w_{j+1} w_{j+2} \dots w_n$$

アーリー項も逆ドット項も、ドットの位置は項の集合Iの添字iとして表す。アーリー項と逆ドット項の違いは、項中の位置番号jの解釈にある。上記した(3)の定義から明らかなように、逆ドット項では、ドット記号の右の部分が β として解釈済みであるとみなす。項内のCFG規則の右辺の直後の位置番号がjであり、そこから先頭方向にiまで逆戻りした部分 $w_{i+1} w_{i+2} \dots w_j$ が β として解釈済みであると解釈する。ちなみにアーリー項では解析対象文の α が解釈済みであるとされていた。逆ドット項を作成する利点は3、6で説明する。

3 YAGLR法

本章ではTRSの構造と節点のマージ操作を説明し、次にシフト操作、レデュース操作を説明する。最後にYAGLR法による統語解析アルゴリズムを説明する。

3. 1 識別子付き位置番号(Position Number with an Identifier)

YAGLR法で用いるTRSの各節点には、レデュース操作時に逆ドット項、 $[A \rightarrow \alpha \cdot \beta, j] \in I$ を作成するために、状態の他に、jとiを得る情報が含まれていなければならない。さらに、マージを効率よく行うために、識別子付き位置番号を導入する。識別子付き位置番号は全解析過程でユニークな番号であり、後述のマージ操作、シフト操作、レデュース操作により、スタックトップの節点が新たに作られる度に作り出され、スタックトップの節点の位置番号にこれを付加する。識別子が付加された位置番号のことを識別子付き位置番号と呼ぶ。位置番号がiで識別子がxなら、この識別子付き位置番号をxiと書くことにする。識別子はマージの歴史を保存しておくためのもので、それにより重複した統合操作を避けることができる (3, 2)。レデュース操作時の逆ドット項の生成には、識別子付き位置番号のなかの、位置番号のみが使われることに注意したい (3, 3)。

YAGSSで用いるTRSの節点の構造を以下に示す：

[識別子付き位置番号の集合、状態]

たとえば、 $(^32, ^43, ^54), (^12, ^32, ^6)$ は識別子付き位置番号の集合であり、 $((^32, ^43, ^54), 2), ((^12, ^32, ^6), 4)$ はTRSの節点である。

3. 2 節点のマージ操作

YAGLR法は富田法以上にスタックのマージを進めることができる。TRSSの構造は富田法で用いるGSSの構造ほど複雑にならない。スタックトップの節点中の識別子付き位置番号の集合には、最後にシフトした語に対する識別子付き位置番号が一つ含まれるだけである。以下では状態が同一の二つの節点のマージ操作について説明するが、これを三つ以上の節点のマージの場合に容易に拡張することができる。(M1)と(M2)により、状態が等しい二つの節点をマージして一つの節点(統合節点)にすることができる。

(M1) スタックトップの節点 $[i_1, s]$ と $[i_2, s]$ は、節点 $[i_1, s]$ にマージする。但し、 i_2 は新たに作り出されたユニークな識別子である。また統合前の二つの節点の持つ全ての子節点を、統合節点の子節点の候補としマージの対象とする。

(M2) スタックトップ以外の二つの節点 $[M, s]$ と $[N, s]$ の場合、節点 $[M \cup N, s]$ にマージする。ただし、 M が N の部分集合なら、節点 $[N, s]$ の子節点を統合節点の子節点とする。さもなければ、統合前の二つの節点の持つ全ての子節点を、統合節点の子節点とし、マージの対象とする。前者の場合、統合節点の子節点は節点 $[N, s]$ の子節点に等しく、それらの状態は全て異なるので、子節点同士のマージを更に進める必要がない。この事実は、TRSSのマージ操作を中途で中断して良いことを意味し、TRSSのマージ操作でマージの重複計算を避けることができるので重要である。なお(M2)の妥当性については、レデュース操作が関係するので、レデュース操作の説明後に与える。TRSSマージ操作アルゴリズムのPASCAL風の定義を以下に与える。

[マージ操作のアルゴリズム]

```

procedure main(TRSS: set of trees)
  var current_set: set of nodes;
  begin
    current_set := all the top nodes of trees in TRSS;
    merge(current_set, "top")
  end main;

function merge(current_set: set of nodes, flag: char): char
  var new_current_set: set of nodes;
  begin
    if all the nodes in current_set have distinct states
      then case flag of
        "top" : "fail";
        "otherwise" : "success"
      end
    else begin
      Form groups of nodes with the same state in the current_set;
      For each group of nodes
        do begin
          case flag of
            "top" : apply (M1) repeatedly to get a merged node;
            "otherwise" : apply (M2) repeatedly to get a merged node
          end;
          new_current_set := the children nodes of the merged node;
          merge(new_current_set, "otherwise")
        end
      end
    end merge;

```

[TRSSマージの例]

- (a) $\cdots - [i_2, i_3, 8] - [i_2, i_3, 75, 9] - [i_6, 1] (ト, フ)$
- (b) $\cdots - [i_2, i_3, 64, 8] - [i_2, i_3, 75, 9] - [i_6, 1] (ト, フ)$
↓ (マージ)
- (c) $\cdots - [i_2, i_3, 64, 8] - [i_2, i_3, 4, 75, 9] - [i_6, 1] (ト, フ)$

上記したマージ操作のアルゴリズムに従い、(a)と(b)のマージを行う。スタックトップ節点間の統合であるから、(M1)によるマージを行い、統合節点の位置番号6に新しい識別子10を付加する。次に子節点間のマージを行うが、両者の状態が9で等しいから(M2)によるマージを葉の方向に向けて進める。次の孫節点同士も同一の状態8を持つから(M2)によるマージを行う。ただしこの場合、両節点の識別子付き位置番号の集合間に、部分集合関係が成立するので、(M2)により(b)の節点 $[i_2, i_3, 64, 8]$ 以下をそのままマージの結果として終了する。こうしてマージのための再計算を避けることができる。富田法ではこうした深いマージは不可能である。

3.3 シフト操作とレデュース操作

[シフト操作]

(a) のスタックトップにシフト操作 "sh u" を行った結果を(b)に示す。それと共に、(c)に示す逆ドット項を作る。シフト操作により、スタックトップの節点の位置番号が1増える。

- (a) $\cdots - [M, s] - [(i_j, t)] (ト, フ) "sh u"$
- (b) $\cdots - [M, s] - [(i_j, t)] - [(i_{(j+1)}, u)] (ト, フ)$
- (c) $i_j \rightarrow [X \rightarrow \cdot w_{j+1}, j+1]$

[レデュース操作]

レデュース操作の指定するCFG規則 $A \rightarrow X_1 X_2 \dots X_n$ により、スタックは(a)から(b)に変わると共に、(c)に示す逆ドット項を作成する。

$$(a) \cdots - [N_k, s_k] - [N_{k+1}, s_{k+1}] - \cdots - [N_{k+n}, s_{k+n}] (\text{トガ' }) \\ \downarrow \\ (b) \cdots - [N_k, s_k] - [X, t] (\text{トガ' })$$

ここで状態 t は、レデュース操作後にCFG規則の左辺の非終端記号 A と状態 s_k とから GOTO 表を参照して決めた新しい状態である。また、 $N_k = \{\alpha, \beta, \dots\}$ ， $N_{k+1} = \{\gamma, \delta, \dots\}$ ， \dots ， $N_{k+n-1} = \{\epsilon, \zeta, \dots, \eta\}$ ， $N_{k+n} = \{\nu\}$ ， $N = \{\nu\}$ であるとする。スタックトップの節点の位置番号の集合 N_{k+n} は、最後にシフトした語の識別子付き位置番号 ν を含み、レデュース後のスタックトップの節点の位置番号の集合 N は、新しいユニークな識別子付き位置番号 ν を含むことに注意したい。これは、レデュース後のスタックのトップに新しい節点が作られるためである。レデュース操作の前後で位置番号 j は変わらないことにも注意したい。

(c) 逆ドット項の作成：

$$\begin{array}{ll} l_a & a [A \rightarrow X_1 X_2 \dots X_n, j] \\ l_b & a [A \rightarrow X_1 X_2 \dots X_n, j] \\ \dots & \\ l_c & a [A \rightarrow X_1 \dots X_n, j] \\ l_d & a [A \rightarrow X_1 \dots X_n, j] \\ \dots & \\ l_e & a [A \rightarrow X_1 X_2 \dots X_n, j] \\ l_f & a [A \rightarrow X_1 X_2 \dots X_n, j] \\ \dots & \\ l_g & a [A \rightarrow X_1 X_2 \dots X_n, j] \end{array}$$

レデュース操作では、操作の指定するCFG規則の右辺の非終端記号の数だけ、スタックトップから節点をポップする。ポップする各節点には、CFG規則の右辺の非終端記号が対応するので、上記した(1)の場合、 X_1, X_2, \dots, X_n のそれぞれには、スタックの節点、 $[N_{k+1}, s_{k+1}], [N_{k+2}, s_{k+2}], \dots, [N_{k+n}, s_{k+n}]$ が順に対応する。TRS に含まれる節点の位置番号は、どこまでがシフト済み（処理済み）かを示しているので、これを用いて逆ドット項のドット位置（逆ドット項の集合 I の添字）を知ることができる。逆ドット項内に書かれる位置番号は全て、スタックトップの節点の持つ位置番号 j であり、位置番号に付加された識別子は、逆ドット項の作成には寄与していないことに注意。

3.4 節点のマージ操作 (M2) の妥当性

二つの節点 $[M, s]$ と $[N, s]$ に対して (M2) が成立することは、次の(1), (2), (3)から証明できる。

- (1) YAGLR 法では、全てのレデュース操作が終了してからシフト操作を行うため、TRS の全てのスタックトップの節点の持つ位置番号は、レデュース操作中は最後にシフトした語の位置番号に等しく同一である。従って、レデュース操作で作成する逆ドット項中の位置番号はマージ前もマージ後と変わらず、最後にシフトした語の位置番号である。
- (2) 逆ドット項の集合 I の添字は、逆ドット項中のドットの位置番号を表す。この添字は、TRS 中の節点の持つ位置番号の集合から 1 つずつ取り出されるもので、マージ前の位置番号の集合から個々に取り出しても、マージ後の 1 つの和集合から取り出しても変わらない。
- (3) M が N の部分集合でない場合に、統合前の二つの節点の持つ全ての子節点を、統合節点の子節点として良いことは明らかである。一方、 M が N の部分集合である場合に、節点 $[N, s]$ の子節点をそのまま統合節点の子節点として良い理由は次の通りである。 M が N の部分集合であるということは、節点 $[M, s]$ を得るために行ったのと全く同じマージの履歴が、節点 $[N, s]$ に残されていることを意味している。従って、節点 $[M, s]$ と節点 $[N, s]$ 以降のマージの結果は、既に節点 $[N, s]$ の子節点以降に含まれているので、両節点の統合節点を $[N, s]$ とし、その子節点を、 $[N, s]$ の子節点に一致させて良い。

3.5 YAGLR 法による統語解析アルゴリズム

YAGLR 法は以下のアルゴリズムにより入力文を解析する。

- (1) TRS の初期状態を、 $(\#ト)$ + $(^0, 0)$ (トガ')
- (2) スタックトップの節点の状態と先読み語の品詞により、LR 表から実行すべき操作を求めて実行する。操作にはシフト操作、レデュース操作、受理、エラーがある。受理なら解析成功として、またエラーなら解析失敗として終了。それ以外は(3)へ。
- (3) TRS のスタックに対する操作が全てシフト操作なら、全てのシフト操作実行前と実行後に TRS にマージ操作を施し(2)へ。さもなければ(4)へ。
- (4) レデュース操作を必要とする一つのスタック S に対して全てのレデュース操作を施し、その結果えられる新しいスタックと、S 以外の既存のスタックとの間でマージを行い(3)へ。

3.6 なぜ逆ドット項か

レデュース操作を説明したので、YAGLR 法で逆ドット項を作る意義が理解できる。実は TRS を利用してレデュース操作時にアーリー項を作ることができるのであるが、ここで逆ドット項を作る理由を述べてみたい。今、以下の(r1)から(r3)に示す CFG を用いて、 $a You _ walk _ z _ in _ a _ the _ park _ s$ の解析を試みて見よう。

(r1) $S \rightarrow NP VP$ (r2) $VP \rightarrow V PP$ (r3) $VP \rightarrow V$
 解析が進み、"You walk" までをシフトした段階では、規則(r1)を用いたレデュース操作を行い、図3. 1に示すアーリー項を作ることができる。解析がさらに右に進み、"in the park" が PP (前置詞句) として受理され、規則(r2)により "walk in the park" が VP として受理されたとしよう。すると再び規則(r1)を用いたレデュース操作を行い図3. 2に示すアーリー項を作ることができる。

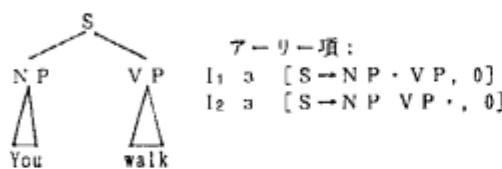


図3. 1 部分統語解析木とアーリー項(1)

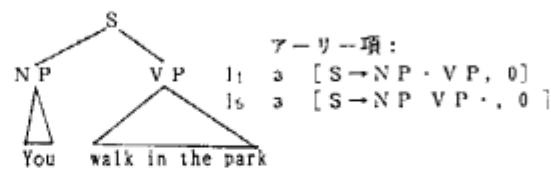


図3. 2 部分統語解析木とアーリー項(2)

ここで、項の集合 I_1 に所属させようとするアーリー項は、 I_1 の節点として既に登録済みのアーリー項と同じで重複することに注意したい。このことは、項の集合への節点の重複登録を避けるために、過去に作成したアーリー項の集合を調べる必要があることを意味している。従って過去に作成したアーリー項の全集合が、任意の時点での検索可能でなければならない。

一方逆ドット項を作成する場合には、逆ドット項中の位置番号は、レデュース操作時のスタックトップの節点のもつ位置番号である。この位置番号は、新たなシフト操作がなされぬ限り変わらないだけでなく、過去のシフト操作時の位置番号とも異なる。そのため逆ドット項の重複検査は、現在実行中の全てのレデュース操作により作られる項の集合を調べるだけでもよい。この様にして、重複する逆ドット項を検査する範囲を局所化することが可能になり、容易にしかも高速な重複処理を行うことができる。新たなシフト操作を行う度に、重複検査のために、それまでに作られた逆ドット項を参照する必要はないのである。上記例の場合、逆ドット項なら以下が生成され、いずれも互いに重複しない。

図3. 1 に対応した逆ドット項: $I_1 a [S \rightarrow NP . VP, 2]$
 $I_2 a [S \rightarrow NP VP, 2]$

図3. 2 に対応した逆ドット項: $I_1 a [S \rightarrow NP . VP, 5]$
 $I_2 a [S \rightarrow NP VP, 5]$

これまでの説明から明らかなように、YAGLR法では、 PP 付加規則に対して、アーリー項なら重複した項生成が行われるが、逆ドット項なら重複した生成が行われない。一般に、重複する項の生成は、文法規則に依存するので、逆ドット項の生成が有利かアーリー項の生成が有利かは一概には決められない。しかし、逆ドット項は重複した項の検査を局所化することができることを説明した。これは、LR法が最右導出をベースとしていることと関連している。

3. 7 富田法とYAGLR法による統語解析例

本節では図3. 3に示すCFGを用いて、富田法とYAGLR法のそれぞれに対する統語解析例を示し、両者の比較を行う。解析対象文は「文化がきたから伝わったから...」であるとし、解析中途まで双方のトレースを行う。図3. 3のCFGから図3. 4に示すLR表が得られる。なお、富田法のスタッフの構造は、[部分木番号、状態]である。

- | | |
|--------------------------|--------------------------|
| (1) $S \rightarrow PP S$ | (5) $v \rightarrow きた$ |
| (2) $S \rightarrow v$ | (6) $v \rightarrow 伝わった$ |
| (3) $PP \rightarrow n p$ | (7) $n \rightarrow きた$ |
| (4) $PP \rightarrow S p$ | (8) $n \rightarrow 文化$ |
| | (9) $p \rightarrow から$ |
| | (10) $p \rightarrow が$ |

図3. 3 簡単な日本語CFG

状態	アクション部				GOTO部	
	n	p	v	\$	PP	S
0	sh4		sh3		2	1
1		sh5		受理		
2	sh4		sh3		2	6
3		re2		re2		
4		sh7				
5	re4		re4			
6		sh5 rel		rel		
7	re3		re3			

図3. 4 図3. 3のCFGから得られたLR表

[富田法による解析]

- T1) $[_, 0]-$
- T2) $[_, 0]-[0, 4]-$
- T3) $[_, 0]-[0, 4]-[1, 7]-$
- T4) $[_, 0]-[2, 2]-$
- T5) $[_, 0]-[2, 2]-[3, 4]-$
 $\quad\quad\quad [4, 3]-$

- | | |
|----------------|--------------|
| 文化(n) sh4 : | |
| が(p) sh7 : | 0 [n 文化] |
| きた(n; v) re3 : | 1 [p が] |
| きた(n) sh4 : | 2 [PP (0 1)] |
| きた(v) sh3 : | |
| から(p) sh7 : | 3 [n きた] |
| re2 : | 4 [v きた] |

T6)	$\boxed{[., 0]} - [2, 2] \boxed{[3, 4]} -$ $\boxed{[5, 6]} -$	から(p) sh7 sh5, rel	5 [S (4)]
T7)	$\boxed{[., 0]} - [2, 2] \boxed{[3, 4]} -$ $\boxed{[5, 6]} -$ $\boxed{[6, 1]} -$	から(p) sh7 sh5	6 [S (2 5)]
T8)	$\boxed{[., 0]} - [2, 2] \boxed{[3, 4]} - [7, 1] -$ $\boxed{[5, 6]} - [7, 5] -$ $\boxed{[6, 1]} -$	伝わった(v) re3 re4	7 [p から] sh5 (at T8)によるマージ
T9)	$\boxed{[., 0]} - [2, 2] \boxed{[8, 2]} -$ $\boxed{[9, 2]} -$ $\boxed{[10, 2]} -$	伝わった(v) sh3	8 [PP (3 7)] 9 [PP (5 7)] 10 [PP (6 7)]
T10)	$\boxed{[., 0]} - [2, 2] - [11, 2] - [12, 3] -$ $\boxed{[10, 2]} -$	から(p) re2	マージ (圧縮統括森11) sh3 (at T9)によるマージ
T11)	$\boxed{[., 0]} - [2, 2] - [11, 2] - [13, 6] -$ $\boxed{[10, 2]} -$	から(p) sh5, rel	11 [PP (3 7) (5 7)] 12 [v 伝わった] 13 [S (12)]
T12)	$\boxed{[., 0]} - [14, 1] -$ $\boxed{[15, 6]} -$ $[2, 2] - [11, 2] - [13, 6] -$ $\boxed{[10, 2]} -$	から(p) sh5 sh5, rel sh5	14 [S (10 13)] 15 [S (11 13)]
T13)	$\boxed{[., 0]} - [14, 1] -$ $\boxed{[16, 1]} -$ $[2, 2] - [11, 2] - [13, 6] -$ $\boxed{[10, 2]} -$	から(p) sh5 sh5	16 [S (2 15)]
T14)	$\boxed{[., 0]} - [17, 1] -$ $[2, 2] - [11, 2] - [13, 6] -$ $\boxed{[10, 2]} -$	から(p) sh5 sh5	マージ (圧縮統括森17) 17 [S (2 15) (10 13)]

(以下省略)

[YAGL法による解析]

Y1)	$\boxed{[\{^0\}, 0]} -$	文化(n) sh4	10 [n → 文化, 1]
Y2)	$\boxed{[\{^0\}, 0]} - [\{^1\}, 4] -$	が(p) sh7	11 [p → が, 2]
Y3)	$\boxed{[\{^0\}, 0]} - [\{^1\}, 4] - [\{^2\}, 7] -$	きた(n;v) re3	10 [PP → n p, 2] 11 [PP → n p, 2]
Y4)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] -$	きた(n) sh4 きた(v) sh3	12 [n → きた, 3] 12 [v → きた, 3]
Y5)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^3\}, 4] -$ $\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^5\}, 3] -$	から(p) sh7 sh7, rel	13 [S → v, 3] 10 [S → PP S, 3] 12 [S → PP S, 3]
Y6)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^6\}, 4] -$ $\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^6\}, 6] -$	から(p) sh7 sh5, rel	13 [p → から, 4]
Y7)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^3\}, 4] -$ $\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^6\}, 6] -$ $\boxed{[\{^0\}, 0]} - [\{^7\}, 1] -$	から(p) sh7 sh5 sh5	マージ
Y8)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^4\}, 4] - [\{^6\}, 7] -$ $\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^6\}, 6] - [\{^9\}, 5] -$ $\boxed{[\{^0\}, 0]} - [\{^7\}, 1] -$	伝わった(v) re3 re4	12 [PP → n p, 4] 13 [PP → n p, 4] 12 [PP → S P, 4] 13 [PP → S P, 4] 10 [PP → S P, 4] 13 [PP → S P, 4]*
Y9)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^4\}, 2] -$ $\boxed{[\{^0\}, 0]} -$	伝わった(v) sh3	15 [v → 伝わった, 5]
Y10)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^4\}, 2] - [\{^12\}, 3] -$ $\boxed{[\{^0\}, 0]} -$	から(p) re2	14 [S → v, 5]
Y11)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^4\}, 2] - [\{^13\}, 6] -$ $\boxed{[\{^0\}, 0]} -$	から(p) sh5, rel	12 [S → PP S, 5] 14 [S → PP S, 5] 10 [S → PP S, 5]
Y12)	$\boxed{[\{^0\}, 0]} - [\{^2\}, 2] - [\{^4\}, 2] - [\{^13\}, 6] -$ $\boxed{[\{^0\}, 0]} -$ $\boxed{[\{^0\}, 0]} - [\{^2\}, 2] -$ $\boxed{[\{^0\}, 0]} - [\{^4\}, 6] -$ $\boxed{[\{^0\}, 0]} -$ $\boxed{[\{^5\}, 1] -}$	から(p) sh5 sh5, rel sh5	10 [S → PP S, 5]* 12 [S → PP S, 5]

$Y_{13}) [^{10}0, 0] - [^{12}, 2] - [^{12, 14}, 2] - [^{175}, 6] -$
 $[^{10}0, 0] - [^{10}, 0] - [^{185}, 1] -$
 (以下省略)

から(p) sh5 | マージ
 sh5 |
 注) *付きの項は、重複項を表す。

ここで富田法とYAGL R法の比較をしてみたい。T12)からT14)までとY11)からY13)までを比較してみよう。YAGL R法は、TRSの構造が単純な分だけ、実装が容易である。また、YAGL R法は富田法以上にスタックのマージを進めることができるので、GSSを用いざとも計算効率と計算時に使用する記憶空間の縮小をはかることが可能になり、解析結果に含まれる曖昧性が増大するにつれて、YAGL R法と富田法の違いが際だってくる。4章ではその一端を説明する。

4 解析対象文の長さに対する計算の複雑性

Johnsonによれば、次の一連の文法Lnに対して、解析対象文の長さをn+2としたとき、解析に要する時間は $\Omega(n^*)$ である[Johnson 89][Kipps 89]。

- (1) $S \rightarrow a$
- (2) $S \rightarrow S S$
- (3) $S \rightarrow S^{n+2}$

ただし S^{n+2} は記号Sがn+2個連結したもの省略形である。解析対象文を $a^{n+2} S$ とする(ただし $n > m$)。上記した規則に対して右のTRS表が得られる。今、 $i (m+3 < i < n+2)$ 番目のaをシフトしrelを施した後のTRSの状態を図4.1に示す。ここで $N_{p,q}$ は節点を表し、各節点の持つ識別子付き位置番号の集合間に、矢印で示す包含関係がある。節点XとYの間に $X \rightarrow Y$ があれば、Xの識別子付き位置番号の集合がYの部分集合であることを示す。

容易に示すことができるよう、節点 $N_{p,q}$ の持つ状態 $T_{p,q}$ は、

$$T_{p,q} = \begin{cases} q & \text{if } 1 \leq q \leq m+3, \\ m+3 & \text{if } m+3 \leq q, \\ 0 & \text{if } q = 1. \end{cases}$$

図4.1に対して、次のシフト操作を施すまで、レデュース操作を次々に施す必要があるが、そのための計算コストを考えてみよう。まず、規則(2)を適用してレデュース操作を進める場合の方が、規則(3)を適用する場合に比べてコストがかかることは明らかである。したがって、計算コストの

状態	アクション部		GOTO部
	a	s	
0	sh1		2
1	rel	:rel	
2	sh1	acc	3
3	sh1/re2	re2	4
4	sh1/re2	re2	5
....
m+1	sh1/re2	re2	m+2
m+2	sh1/re2	re2	m+3
m+3	sh1/re2/re3	re2/re3	m+3

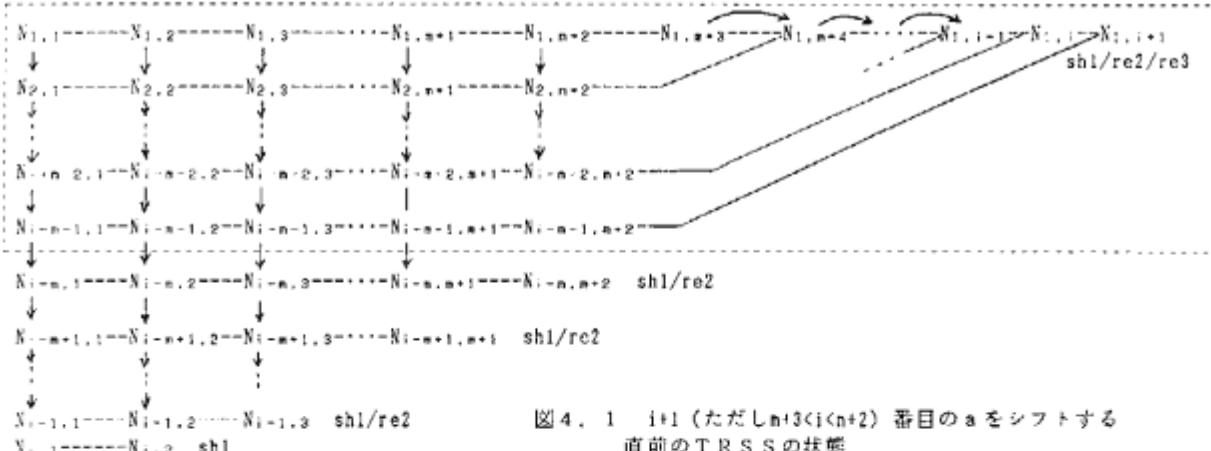
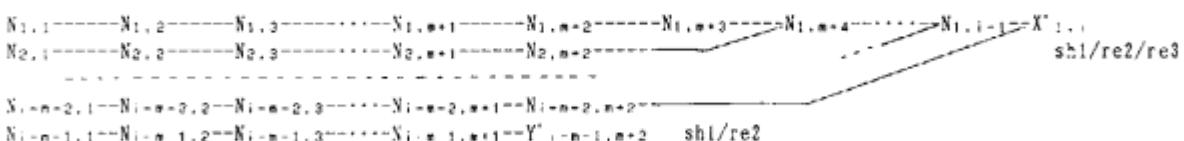


図4.1 $i+1$ (ただし $m+3 < i < n+2$) 番目のaをシフトする直前のTRSの状態

オーダーを求めるためには、規則(2)を適用するレデュースを考えれば十分である。また、最上部の点線で囲まれたTRSの部分のみが、入力文の長さに依存しているので、これに対するレデュース操作の計算コストを考えれば十分である。今、図4.1の最上部のスタックに対して、規則(2)によるレデュース操作を一度施すと、下図の二つのスタック構造が得られる。



上図のスタックは、図4.1の統合化可能なスタックとマージされる。たとえば、sh1を待つ最長スタック同士のスタックトップの状態はn+3で等しいのでマージされる。このとき、図4.1に記述した節点間の識別子付き位置番号の集合間に、包含関係が成り立つため、マージ操作は先頭から二つの目までの節点同士($< N_{1,i+1} \rightarrow X_{1,i+1} >$, $< N_{1,i} \rightarrow N_{1,i+1} >$, $< N_{1,-n-1,n-2} \rightarrow N_{1,-n-2,n-2} >$)で行えよう。先頭から3番目の節点同士($< N_{1,-n-1} \rightarrow N_{1,-n-2} >$ または $< N_{1,-n,n+1} \rightarrow N_{1,-n-1} >$)

$\langle \dots \rangle$ のマージなどは、識別子付き位置番号の集合間に部分集合関係が成立しているので、それ以上マージを莫の方向に向けて進める必要はない。

以上の考察から、re2操作後のマージには先頭から3番目のマージまでを考慮すればよいことが分かる。しかも、それぞれのマージには、識別子付き位置番号の集合の要素が高々 $i+1$ 含まれるだけであるから、マージに要す計算コストは i のオーダーである。一方、re2は繰り返し行われるが、その回数は高々 i 回である。したがって、レデュース操作とマージ操作の総和は、 $\Omega(i^2)$ である。以上より、 a^{n+2} の入力文を解析するために要する総計算コストのオーダーは、 $\sum_{i=2}^{n+2} \Omega(i^2) = \Omega(n^3)$ となり、証明された。

次に図4、1の各節点中の識別子付き位置番号の数は：節点 $N_{i,j}$ に対して、 $2 \leq j \leq m+3$ の場合1、 $m+4 \leq j \leq i-1$ の場合2、 $i \leq j \leq i+1$ の場合3、節点 $N_{k,j}$ ($2 \leq k \leq i-m-1$, $2 \leq j \leq m+2$)に対してkである。したがって、点線で囲まれた部分の節点中の識別子付き位置番号の総和は、 $1*(n+4)+2*(i-m-4)+(m+2)(2+3+\dots+(i-m-1))$ である。この総和のオーダーは i^2 である。これは、点線で囲まれた部分以外の節点に含まれる識別子付き位置番号の総和のオーダーと等しい。以上から、入力文 a^{n+2} を解析するために必要なTRSSの空間のオーダーは n^2 であり、また、作成される逆ドット項の数の総和のオーダーは、 n^3 であることが証明された。

5 おわりに

本論文ではYAGLR法と呼ぶ新しい統語解析アルゴリズムを提案した。YAGLR法は、既存の方法と比べていくつかの利点を持っている。

- (A) Johnsonらの与えた特殊な文法に対しても、解析時間のオーダーは、解析対象文の長さ n に対してアーリー法と同等の n^3 のオーダーに抑えることができる。
- (B) レデュース操作時に、アーリー項と対称な逆ドット項を作成しながら統語解析を行う。逆ドット項を生成する利点について考察した。
- (C) 先読み情報を用いているので、最終的に生成される逆ドット項の数はアーリー法より少ない。
- (D) 富田法以上にスタックのマージを行うことができるため、グラフ構造化スタックより構造が単純で操作し易い木構造化スタックを用いることができる。それによりスタック操作が単純になる。
- (E) 解析終了後に、作成した逆ドット項を組み合わせて統語解析木を作る。

今後の検討課題として、以下のものを挙げることができる。

- (1) YAGLR法の統語解析時間のオーダーが、一般的CFGに対して n^3 のオーダーであることの証明。
- (2) YAGLR法による統語解析に要する空間量の評価。
- (3) YAGLR法の並列解析アルゴリズムの研究。
- (4) 逆ドット項からの統語構造の並列抽出アルゴリズムの研究。

(1)は重要な研究課題であるが、(2)に述べたように、解析に要す空間の大きさの評価も行う必要がある。Kippsは富田法に比較的単純な修正を施すことで、一般的CFGに対する計算時間のオーダーを見 (n^3) に抑えることができるることを示した[Kipps 89]。しかし、このアルゴリズムは富田法以上の空間を要す。これに対してYAGLR法では、スタックのマージを徹底的に行うアルゴリズムであるから、計算空間を抑えることができる。一般的CFGに対するYAGLR法の使用記憶空間の大きさの詳細な検討を行う必要がある。

(3)に述べたように、富田法と同様、YAGLR法の並列解析アルゴリズムを検討することも意味があるだろう[Tanaka 89]。本稿で述べたアルゴリズムは構型探索の逐次型アルゴリズムであり、できるだけ無駄な解析を省いたアルゴリズムである。この長所を生かしつつ、シフトレデュース・コンフリクトが発生したときの待ち合わせができるだけ少なくなる並列解析アルゴリズムを研究する必要がある。それには沼崎らの研究が参考になろう[Numazaki 90]。

YAGLR法では、統語解析後に得られた逆ドット項の集合から統語解析木を作り出さねばならない。これはアーリー法と双対なアルゴリズムを用いれば良い。一つの統語解析木を計算するための時間は n^2 のオーダーであることが知られている[Aho 72]。この統語解析木の計算にも多くの並列性が観察されるので、統語解析木抽出用の並列アルゴリズムについて今後研究する必要がある。

参考文献

- [Aho 72] Aho, A. V. and Ullman, J. D.: The Theory of Parsing, Translation and Compiling. Prentice-Hall, New Jersey(1972).
- [Earley 70] Earley, J.: An Efficient Augmented-Context-Free Parsing Algorithm. Comm. of ACM, 13, 1-2, 95-102(1970).
- [Johnson 89] Johnson, M.: The Computational Complexity of Tomita's Algorithm. International Parsing Workshop '89, Carnegie-Mellon University, 203-208(1989).
- [Kipps 89] Kipps, J. W.: Analysis of Tomita's Algorithm for General Context-Free Parsing. International Parsing Workshop '89, Carnegie-Mellon University, 193-202(1989).
- [Numazaki 90] Numazaki, H. and Tanaka, H.: A new Parallel Algorithm for Generalized LR Parsing. COLING'90 305-310(1990).
- [Suresh 90] Suresh, K. G.: Yet Another Generalized LR Parser-Report4. Tanaka Lab., Tokyo Inst. of Technology (1990. 4. 16).
- [田中 89] 田中穂積：自然言語解析の基礎、産業図書(1989)。
- [Tanaka 89] Tanaka, H. and Numazaki, H.: Parallel Generalized LR Parser Based on Logic Programming. 1st Australian-Japan Joint Symposium on Natural Language processing, 201-211(1989).
- [Tomita 86] Tomita, M.: Efficient Parsing for Natural Language. Kluwer, Boston, Mass(1986).
- [Tomita 87] Tomita, M.: An Efficient Augmented-Context-Free Parsing Algorithm. Computational Linguistics, 13, 31-46(1987).