

## LC構文解析法を拡張した構文解析法について

椎名 広光<sup>†</sup>, 増山 繁<sup>‡</sup>

<sup>†</sup> 岡山理科大学 〒700-0005 岡山県岡山市理大町 1-1

<sup>‡</sup> 豊橋技術科学大学 〒441-8580 愛知県豊橋市天伯町雲雀ヶ丘 1-1

### あらまし

計算可能な一般の句構造言語を構文解析する場合、文法の構造上文法記号間の交差がいくつも生じるため、入力列から作られる構文解析木は複雑になる。更に、構文解析の処理中においては、構文解析木の候補者や複数の構文解析木の生成も考慮しながら計算しなければならない。そのため、構文解析の処理中に予測した構文解析木が誤っている場合は、訂正のためにバックトラック処理が必要になり、解析時間が増大する。そこで本研究では、コンパイラなどで用いられている構文解析の中でも上昇型と下降型を組み合わせた混合型に分類されるLC構文解析法の先読み文字列の定義と生成規則の出力順を変更することで拡張を行なう。それによって、一般の句構造言語を構文解析する処理よりもバックトラック処理の回数を減少させ、また、LC文法で生成される言語よりも広い言語に対して構文解析が可能とことを示す。

**キーワード** LC 構文解析, LC 文法, 句構造文法, 文脈自由文法

## The proposal of the extended *LC* parser

Hiromitsu Shiina<sup>†</sup>, Sigeru Masuyama<sup>‡</sup>

<sup>†</sup> Okayama University of Science

1-1, Ridai-cho, Okayama city, Okayama, 700-0005, Japan

<sup>‡</sup> Toyohashi University of Technology

1-1, tempaku-cho, Toyohashi city, Aichi, 441-8580, Japan

### abstract

When we parse recursive languages, the parse tree becomes complicated in structure. Furthermore, we must take into consideration the candidate of the parse tree and the generation of multiple parse trees in parsing processing. If the parser generates some wrong parse trees, then its backtracking process for correction is required. Therefore, in this paper, we propose an extension the *LC* parser as a mixed type of top-down and bottom-up parsing. We show that by using this parser the number of backtrack can be reduced and languages which is in the class larger than that of context-free languages can be proposed by changing the definition of lookahead strings.

**keywords** *LC* parser, *LC* grammar, phrase structure grammar, context-free grammar

## 1 はじめに

計算可能な一般の句構造言語を構文解析する場合、文法の構造上、入力列から作られる構文解析木には、文法記号の親が複数あって枝に交差を持つことになり、複雑になる。

更に、構文解析の処理中においては、構文解析木の候補者や複数の構文解析木生成も考慮しながら計算しなければならない。そのため、予測した構文解析木が誤っている場合は、訂正のためにバックトラック処理が必要になり、解析時間が増大する。

一方、自然言語処理などにおいては、文脈自由言語だけではなく、それより広いクラスに属する言語を処理する必要がある。しかしながら、自然言語の文法は、部分的に文脈自由文法や句構造文法の生成規則が含まれてはいるものの、ほとんどの生成規則が文脈自由文法のもので、構文解析時に最悪のケースばかりに対応していると処理時間が遅くなってしまふ。

そこで本研究では、コンパイラなどで利用されている構文解析の中でも上昇型と下降型を組み合わせた混合型に分類されるLC構文解析法 [8] を拡張する。それによって、予測構文解析木を作成しないで処理する方法とLC文法 [8] で生成される言語よりも広いクラスの言語を生成する文法を提案する。

なお、文脈自由文法  $G = (N, T, P, S)$ ,  $N$ : 非終端記号の集合,  $T$ : 終端記号の集合,  $P$ : 生成規則の集合,  $S$ : 開始記号と入力文字列  $a_1 a_2 \dots a_n \$$ ,  $a_1, a_2, \dots, a_n \in T$  を与える。また構文解析実行の便宜上、入力文字列の最後に終りを示す特殊記号  $\$$  ( $\notin N \cup T$ ) を用意する。

## 2 準備

本章では、準備として、LC文法の定義およびLC構文解析法について、簡単に述べる。なお、本論文におけるLC文法の定義は、文献 [8] に準拠する。

### 2.1 左隅出現でない非終端記号

最左導出の過程で得られる文法記号列である左文形式  $wA\delta$ ,  $w \in T^*$ ,  $A \in N$ ,  $\delta \in (T \cup N)^*$ , の最左位置の  $A$  が、生成規則の左隅の記号として導入されたものでない時、左隅でない非終端記号  $A$  という。左隅出現でない左文形式  $wA\delta$  の非終端記号  $A$  を  $S \xrightarrow{lc} wA\delta$  と表す。

### 2.2 LC文法 [8]

次の条件を満たす文脈自由文法  $G$  を  $LC(k)$  文法という。左隅出現でない任意の非終端記号  $A$ ,  $S \xrightarrow{lc} wA\delta$  と長さ  $k$  の先読み  $u$  に対して、いつも以下のような生成規則  $B \rightarrow \alpha$  が高々 1 個存在し、かつ  $A \xrightarrow{*} B\gamma$  で次のいずれかが成り立つ。

- (i)  $\alpha = C\beta$  なら、 $u$  は  $FIRST_k(\beta\gamma\delta)$  の要素である。  
(ii) さらに、 $C = A$  なら、 $u$  は  $FIRST_k(\delta)$  の要素でない。
- $\alpha$  が非終端記号で始まらないなら、 $u$  は  $FIRST_k(\alpha\gamma\delta)$  の要素である。

ここで、 $w \xrightarrow{*} w_1 w_2 \dots w_n$ ,  $w_1, w_2, \dots, w_n \in T$  とするとき、 $FIRST_k(w)$  は次のように定義される。

$$FIRST_k(w) = \begin{cases} w_1 w_2 \dots w_k, & \text{if } n \geq k, \\ w_1 w_2 \dots w_n, & \text{if } n < k. \end{cases}$$

### 2.3 LC構文解析法 [8]

LC構文解析は、上昇型と下降型を組み合わせた混合型に分類される構文解析法である。その解析の実行においては、木の第1部分木を下から上へ、第2部分木から最右部分木までそれぞれ上から下への走査を再帰的に行ないながら構文解析木を決定する。以下に、LC構文解析の動作を示す。

1.  $T$ が葉のみの場合,
  - 訪問 (根 ( $T$ ))
2.  $T$ が葉のみでない場合,
  - LC走査 (第1部分木 ( $I$ ));
  - 訪問 (根 ( $T$ ));
  - 出力 (根 ( $T$ ));
  - LC走査 (第2部分木 ( $I$ ));
  - ...
  - LC走査 (最右部分木 ( $I$ ));

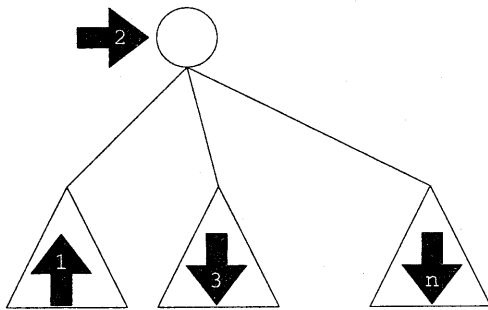


図 1: LC走査

例えば、LC(1)文法の例  $G = (N, T, P, S)$ ,  $N = \{S, A, B\}$ ,  $T = \{a, b, c\}$ ,  $P = \{S \rightarrow AC, C \rightarrow CB, C \rightarrow B, A \rightarrow a, B \rightarrow b, B \rightarrow c\}$  に対して入力文字列を

$I = abbb$  とする。この時、LC構文解析は、(1) $A \rightarrow a$ , (2) $S \rightarrow AC$ , (3) $B \rightarrow b$ , (4) $C \rightarrow B$ , (5) $B \rightarrow b$ , (6) $C \rightarrow CB$ , (7) $C \rightarrow CB$ , (8) $B \rightarrow b$  の順で生成規則を当てはめて構文解析木を作成する (図 2)。特に、(2)の  $S \rightarrow AC$  の後に (3)の  $B \rightarrow b$  を当てはめるのが LC構文解析の特徴である。

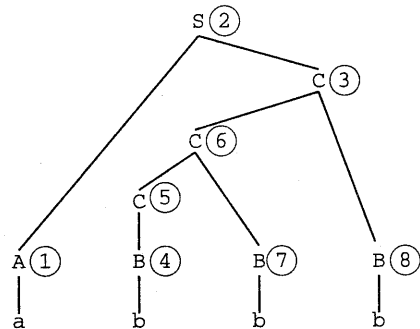


図 2: LC構文解析の例

### 3 LC構文解析の拡張 (ELC 構文解析)

文脈自由言語以下のクラスの言語の場合、生成規則の矢印の左側が複数の文字列である条件  $\alpha \rightarrow \beta \in P$ ,  $|\alpha| > 1$ ,  $|\beta| > 0$  を満たす生成規則の適用を含まないため、構文解析木の兄弟は同じ入力文字を子供とすることはない。そのため、それぞれの文法記号から生成される部分文字列は文法の構造から決定できる。

一方、文脈自由言語より大きいクラスの言語に対する構文解析木には矢印の左側が複数文字列である生成規則の適用がされ、解析木中に図 3 に示すような交差をもつ (図 3,  $AB \rightarrow CDE$ )。そのため、文脈自由言語より大きいクラスの言

語では、構文解析木の兄弟が同じ入力文字の先祖となる場合が存在する。

例えば、図3のCから生成される部分文字列とDから生成される部分文字列とは重複する部分があり、CとDは入力文字aの先祖である。

よって、Cを決定した時点でDを決定するためにDの先頭に引き戻す処理方法やEを決定した後DとCを決定する処理方法[10]をとらなければならない。

これらのどちらの処理方法もバックトラックの回数が多い。そこでバックトラックを避けるため、本研究では先読み文字列の定義と部分文字列の有効範囲を定義することによって、構文解析の左隅を決定するだけで生成規則が決定できる拡張LC構文解析(以下ELC構文解析)およびELC(k)文法を提案する。

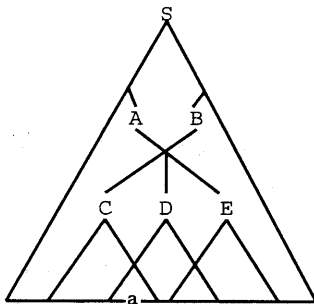


図3: 構文解析木の交差

### 3.1 文法定義の準備

先に、ELC(k)文法を定義するのに必要な用語を第3.1.1節と第3.1.2節で定義する。

#### 3.1.1 最左導出

文脈自由文法における最左導出は、文形式中の最も左側に位置する非終端記号を書き換える導出であったが、文脈自由言語より広いクラスの言語を生成する文法では、非終端記号の代わりに書き換えられる文法記号列に変更する。そして、この最左導出の過程で得られる文法記号列である左文形式  $wA\delta$  の左隅出現でない左文形式  $wA\delta$  の非終端記号  $A$  を  $S \xrightarrow{*} wA\delta$  と表す。

#### 3.1.2 生成される文字列の範囲

非終端記号  $A$  から生成される文字列を定義する。文脈自由文法では生成規則の矢印の左側は1つの文法記号であるので、 $A$  から生成される文字列は文形式の  $A$  の両側の文字に影響されない。しかし、文脈自由文法より広い場合は生成規則の左側に複数の文字列がくる可能性があるため、文形式の両側の文字に影響される。そこで、始めに非終端記号に  $\cdot$  をつけて、元々始めに  $\cdot$  のついていた非終端記号から展開された文字列を明らかにする。例えば、 $\alpha_1\beta_1\dot{A}\beta_2\alpha_2$  の  $A$  に関する展開を調べたい時は、記号  $\Rightarrow$  を用いて、次のように表す。

$$\begin{aligned} \alpha_1\beta_1\dot{A}\beta_2\alpha_2 &\Rightarrow \alpha_3\beta_3\dot{B}\dot{C}\dot{D}\beta_4\alpha_4, \\ \beta_1A\beta_2 &\rightarrow \beta_3BCD\beta_4 \in P, \\ \alpha_1, \alpha_2, \alpha_3, \alpha_4, \beta_1, \beta_2, \beta_3, \beta_4 &\in (T \cup N)^*. \end{aligned}$$

また、0回以上の展開の繰り返しを  $\dot{A} \xrightarrow{*} \dot{\gamma}, \gamma \in (T \cup N)^*$  と表す。

すべての展開に対して、文形式中の文法記号  $A$  から最左導出によって展開し、 $A$  から生成された文字列の直後の  $k$  個の文字列を  $FIRST_k$  で定義する。

$$FIRST'_k(A, \alpha_2) =$$

$$\left\{ \text{head}_k(\gamma_2) \left| \begin{array}{l} S \xrightarrow{*} \alpha_1 A \alpha_2 \xrightarrow{*} \\ \gamma_1 \beta \gamma_2 \\ \alpha_1, \alpha_2, \beta \in (T \cup N)^* \\ \gamma_1, \gamma_2 \in T^* \end{array} \right. \right\}$$

### 3.2 ELC(k) 文法の定義

ELC(k) 文法は生成規則の矢印の右側 (木の左隅) の始めの文字が決定されると、その生成規則が適用される文法である。

任意の非終端記号  $A$ ,  $S \xrightarrow{*} wA\delta$  と長さ  $k$  の先読み  $u \in T^*$  に対して、以下のような生成規則  $\alpha_1 \rightarrow \alpha_2$ , ( $\alpha_1 \alpha_2 \in (T \cup N)^*$ ), が常に高々 1 個存在し、かつ  $A \xrightarrow{*} \alpha_1 \gamma$ , ( $\gamma \in (T \cup N)^*$ ), で次のいずれかが成り立つ。

1. (i)  $\alpha_2 = C\beta$  なら,  $u$  は  $FIRST'_k(C, \beta\gamma\delta)$  の要素である。  
(ii) さらに,  $C = A$  なら,  $u$  は  $FIRST'_k(A, \delta)$  の要素でない。
2.  $\alpha_2$  が非終端記号で始まらないなら,  $u$  は  $FIRST_k(\alpha\gamma\delta)$  の要素である。
3. 終端記号を生成規則中に持つのは,  $A \rightarrow a$  の形式だけである。

ただし, 3. は本質的な制約ではない。構文解析の処理において, 終端記号の取り扱い手順を簡単にするための制約である。

### 3.3 解析表 LC

構文解析を進めていく上で, 次に適用する生成規則を決定するのに, 先頭記号, 現在処理している文法記号, 先読み文字列の引数が必要となる。それを表の

形でまとめたのが解析表 LC である。解析表 LC は次の手順で作成する。

[解析表 LC の作成]

- (i)  $\alpha_2 = C\beta$  の場合,  $u$  は  $FIRST'_k(C, \beta\gamma\delta)$  の要素で  $S \xrightarrow{*} wA\delta$ ,  $A \xrightarrow{*} \alpha_1 \gamma$  なら  
 $LC(A, C, u) := “\alpha_1 \rightarrow \alpha_2”$  とする。
- (ii)  $\alpha_2$  の先頭が終端記号で, その終端記号を  $a$  とする場合,  $u$  は  $FIRST'_k(A, \beta\gamma\delta)$  の要素で,  $S \xrightarrow{*} wA\delta$ ,  $A \xrightarrow{*} \alpha_1 \gamma$  なら,  $LC(A, a, u) := “\alpha_1 \rightarrow \alpha_2”$  とする。
- (iii)  $\alpha_2$  の先頭が終端記号で, その終端記号を  $a$  とする場合,  $u$  は  $FIRST_k(\delta)$  の要素で,  $S \xrightarrow{*} wA\delta$ ,  $A \xrightarrow{*} \alpha_1 \gamma$  なら,  $LC(A, a, u)$  を未定義とする。

例えば, 文法  $G = \{N, T, P, S\}$ ,  $N = \{S, A, B, C, D\}$ ,  $T = \{a, b, c\}$ ,  $P = \{S \rightarrow ASCD, S \rightarrow ABD, DC \rightarrow CD, BC \rightarrow BB, A \rightarrow a, B \rightarrow b, D \rightarrow c\}$  は ELC(1) の文法である。この文法に対する解析表 LC を表 1 に示す。

表 1: 解析表 LC の例

$LC(S, A, a) = “S \rightarrow ASCD”$
$LC(S, A, b) = “S \rightarrow ABD”$
$LC(S, a, b) = “A \rightarrow a”$
$LC(S, a, a) = “A \rightarrow a”$
$LC(S, b, b) = “B \rightarrow b”$
$LC(B, b, b) = “B \rightarrow b”$
$LC(B, b, c) = “B \rightarrow b”$
$LC(B, B, b) = “BC \rightarrow BB”$
$LC(D, C, c) = “DC \rightarrow CD”$
$LC(D, c, c) = “D \rightarrow c”$
$LC(D, c, \$) = “D \rightarrow c”$

## 4 $ELC(k)$ 構文解析の手順

構文解析の処理では `input_stack`, `ancestor_stack`, `closure_stack` の3つのスタックを利用する。以下にそれぞれのデータ構造を定義する。

1. `input_stack` は、入力文字列及び生成規則の左側の文字列を積むスタック。  

$$\text{input\_stack} = (T \cup N) \times (T \cup N) \times \dots \times (T \cup N).$$
2. `ancestor_stack` は、先祖と予測される文法記号を積むスタック。  

$$\text{ancestor\_stack} = N \times N \times \dots \times N.$$
3. `closure_stack` は、解析が終了した位置を表すために生成規則に  $\cdot$  をつけた項を積むスタック。  

$$\text{closure\_stack} = ([\alpha_1 \rightarrow \alpha_2 \cdot \alpha_3] \times [\alpha_4 \rightarrow \alpha_5 \cdot \alpha_6] \times \dots \times [\alpha_7 \rightarrow \alpha_8 \cdot \alpha_9]), \alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5, \alpha_6, \alpha_7 \in (T \cup N)^*.$$

次に構文解析の手順を示す。

**Step 1:** `ancestor_stack` に開始記号  $S$  を積む。

**Step 2:** `input_stack` の先頭が  $a_i$ , 先読みが  $a_{i+1} \dots a_{i+k+1}$ , `ancestor_stack` の先頭を  $H$  とする。

**Step 2.1:** もし  $LC(H, a_i, a_{i+1} \dots a_{i+k+1}) = "A \rightarrow a_i"$  なら

**Step 2.1.1:** `input_stack` から  $a_i$  を降ろす。

**Step 2.1.2:** " $A \rightarrow a_i$ " を出力する。

**Step 2.1.3:**  $A$  を `input_stack` に積む。

**Step 3:** `input_stack` の先頭を非終端記号  $A$ , 先読みを  $a_{i+1} \dots a_{i+k+1}$ , `ancestor_stack` の先頭を  $H$  とする。

**Step 3.1:** もし `closure_stack` の先頭が  $[\alpha_1 \rightarrow \alpha_2 \cdot A B \alpha_3]$  なら

**Step 3.1.1:** `closure_stack` の先頭  $[\alpha_1 \rightarrow \alpha_2 \cdot A B \alpha_3]$  を  $[\alpha_1 \rightarrow \alpha_2 A \cdot B \alpha_3]$  に入れ換える。

**Step 3.1.2:** `ancestor_stack` に  $B$  を積む。

**Step 3.1.3:** `input_stack` から  $A$  を降ろす。

**Step 3.2:** もし `closure_stack` の先頭が  $[\alpha_1 \rightarrow \alpha_2 \cdot A]$  なら

**Step 3.2.1:**  $\alpha_1 \rightarrow \alpha_2 A$  を出力する。

**Step 3.2.2:** `closure_stack` の先頭  $[\alpha_1 \rightarrow \alpha_2 \cdot A]$  を降ろす。

**Step 3.2.3:** `input_stack` から  $A$  を降ろし,  $\alpha_1$  の逆順で積む。

**Step 3.2.4:** `ancestor_stack` から1つ降ろす。

**Step 3.3:** もし  $LC(H, A, a_{i+1} a_{i+2} \dots a_{i+k+1}) = "a_1 \rightarrow A a_2"$  なら

**Step 3.3.1:** `closure_stack` に  $[\alpha_1 \rightarrow A \cdot a_2]$  を積む。

**Step 3.3.2:** `ancestor_stack` に  $a_2$  の先頭の文字を積む。

**Step 3.4:** それ以外はエラーを出力して終了する。

**Step 4:** Step 2 から Step 3 を繰り返す。`input_stack` の内容が  $S\$,$  `closure_stack` が空であるなら, "`accept`" を出力して終了する。

## 5 構文解析の動作例

入力を  $I = aabbcc$  とした時の構文解析の動作を表2に、また作成される構文解析木を図5に示す。

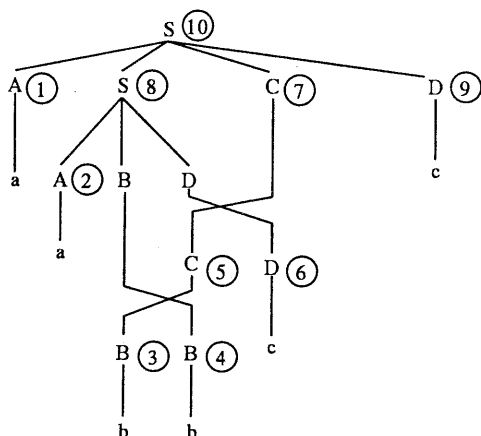


図4: 入力  $I = aabbcc$  に対する構文解析木

## 6 むすび

本研究では、構文解析木の文法記号が支配する文字列の記号の範囲を定義することにより、先読み文字列の定義を拡張した。それにより、LC構文解析法を文脈自由言語のクラスより大きいクラスに属する言語も構文解析ができる  $ELC(k)$  構文解析法に拡張し、また  $ELC(k)$  文法を提案した。今後は、 $ELC(k)$  文法の性質を調査する予定である。

表2:  $ELC$ 構文解析の動作例

input	ancestor	closure	action
1 $aabbcc\$$	S		$LC(S, a, a) = "A \rightarrow a"$ $A \rightarrow a$ を出力 ①
2 $Aabbcc\$$	S		$LC(S, A, a) = "S \rightarrow ASCD"$
3 $abbcc\$$	SS	$[S \rightarrow A \cdot SCD]$	$LC(S, a, b) = "A \rightarrow a"$ $A \rightarrow a$ を出力 ②
4 $Abbcc\$$	SS	$[S \rightarrow A \cdot SCD]$	$LC(S, A, b) = "S \rightarrow ABD"$
5 $bbcc\$$	BSS	$[S \rightarrow A \cdot BD] [S \rightarrow A \cdot SCD]$	$LC(B, b, b) = "B \rightarrow b"$ $B \rightarrow b$ を出力 ③
6 $Bbcc\$$	BSS	$[S \rightarrow A \cdot BD] [S \rightarrow A \cdot SCD]$	$LC(B, B, b) = "BC \rightarrow BB"$
7 $bcc\$$	BBSS	$[S \rightarrow A \cdot BD] [S \rightarrow A \cdot SCD]$	$LC(B, b, c) = "B \rightarrow b"$ $B \rightarrow b$ を出力 ④
8 $Bcc\$$	BSS	$[BC \rightarrow B \cdot B] [S \rightarrow A \cdot BD] [S \rightarrow A \cdot SCD]$	$BC \rightarrow BB$ を出力 ⑤
9 $BCcc\$$	SS	$[S \rightarrow A \cdot BD] [S \rightarrow A \cdot SCD]$	
10 $Ccc\$$	DSS	$[S \rightarrow AB \cdot D] [S \rightarrow A \cdot SCD]$	$LC(D, C, c) = "DC \rightarrow CD"$
11 $cc\$$	DDSS	$[S \rightarrow AB \cdot D] [S \rightarrow A \cdot SCD]$	$LC(D, c, c) = "D \rightarrow c"$ $D \rightarrow c$ を出力 ⑥
12 $Dc\$$	DSS	$[DC \rightarrow C \cdot D] [S \rightarrow AB \cdot D] [S \rightarrow A \cdot SCD]$	$DC \rightarrow CD$ を出力 ⑦
13 $DCc\$$	SS	$[S \rightarrow AB \cdot D] [S \rightarrow A \cdot SCD]$	$S \rightarrow ABD$ を出力 ⑧
14 $SCc\$$	S	$[S \rightarrow A \cdot SCD]$	
15 $Cc\$$	S	$[S \rightarrow AS \cdot CD]$	
16 $c\$$	DS	$[S \rightarrow ASC \cdot D]$	$LC(D, c, \$) = "D \rightarrow c"$ $D \rightarrow c$ を出力 ⑨
17 $D\$$	DS	$[S \rightarrow ASC \cdot D]$	$S \rightarrow ASCD$ を出力 ⑩
18 $S\$$			accept

## 参考文献

- [1] S. Y. Kuroda, *Classes of languages and linear-bounded automata*. Inform. Control, 7, pp207-223, 1967.
- [2] J. Earley, *An efficient context-free parsing algorithm*, Commun. ACM, 13, 2, pp94-102, 1970.
- [3] J. Loeckx, *The parsing for general phrase-structure grammars*, Inform. Control, 16, pp443-464, 1970.
- [4] G. Sh. Vold'man, *A parsing algorithm for context-sensitive grammars*, Program. Comput. Software, 7, pp302-307, 1981.
- [5] L.A. Harris, *SLR(1) and LALR(1) parsing for unrestricted grammar*, Acta Inform., 24, pp191-209, 1987.
- [6] M. Tomita, *An Efficient Augmented-Context-Free Parsing Algorithm*, Computational Linguistics, 13,1-2, pp31-46, 1987.
- [7] M. Tomita ed, *Generalized LR parsing*, Kluwer Academic Publisher, 1991.
- [8] 徳田 雄洋, 言語の構文解析, 共立出版, 1995.
- [9] K. Sikkel, *Parsing Schemata*, Springer, 1997.
- [10] H. Shiina and S. Masuyama, *Proposal of the unrestricted LR(k) grammar and its parser*, Mathematica Japonica, 46, 1, pp129-142, 1997.