

Research Reports on Mathematical and Computing Sciences

An Efficient Parsing for Highly Ambiguous
Context-Free Grammars Based on Pruning

Shin-ichi Morimoto

Jan 2007, C-240

Department of
Mathematical and
Computing Sciences
Tokyo Institute of Technology

SERIES **C**: Computer Science

An Efficient Parsing for Highly Ambiguous Context-Free Grammars Based on Pruning

Shin-ichi Morimoto
NEC Aerospace Systems, Ltd.
s-morimoto@ah.jp.nec.com

In this report, I introduce an efficient parsing algorithm even for highly ambiguous context-free grammars, which is intended to be applied to natural language processing. It uses a graph-structured stack as those by Tomita (1991), Kipps (1991), and Nederhof (1993). Here, a graph-structured stack is a directed graph whose nodes are sets of pairs of an item and an input string, and whose edges are pointers (parent pointers), each pointing to a parent node in the parse tree. The new technique here is to prune the pointers pointing to the nodes having the same item part. By this pruning, we cannot obtain all parse trees, but we can reduce the size of the set of parent pointers to a constant, thus making the time complexity of this algorithm $O(n^2)$. In this report, I first define the algorithm and show its correctness, and then I analyze its time complexity. Then I show that by this parsing a class of typical highly ambiguous grammars given by Kipps (1991) that need $O(n^3)$ by his algorithm can be parsed in $O(n^2)$. Finally some study is made on the range of prunable grammars and application for natural language processing.

1 Introduction

Many parsing algorithms for general context-free grammars have been presented, in which there are mainly two approaches: the tabular approach [1], [4], and the stack approach [3], [5], [11]. The tabular approach uses an $n \times n$ matrix for an input $a_1 \dots a_n$ where the (i, j) element of this matrix corresponds to a set of items that represents recognition of the part of the input $a_{i+1} \dots a_j$. On the other hand, the stack approach uses stacks whose elements are states [3], [11] or items [5]. In order to parse a grammar with ambiguity, these algorithms use a set of stacks, or a graph structured stack that is a directed graph whose nodes correspond to stack elements and edges correspond to pointers to previous elements, which are called parent pointers. An example of stacks whose top elements are D and a corresponding graph structured stack is shown in the left part of Fig. 1.

As for efficiency, the tabular approaches are generally better than stack approaches, since time complexity of the former is $O(n^3)$, and the latter is $O(n^{p+1})$ in the original implementation [11], where p is the longest length of the right-hand side of productions. There are modified versions of the stack approaches as in [3] and [5], whose theoretical time complexities are $O(n^3)$. However, the implementation in [5] is presented only in its outline, and according to [3], its implementation is not practical. For practical situations the stack approaches seem to be better than the tabular approaches [8].

The algorithm presented here, named Pruned Graph-Structured Stack Parsing, is a modification of the algorithm by Nederhof [5]. In this algorithm when a new node is generated during parsing, we prune some of parent pointers of the new generated node if possible. The graph structured stack after pruning and its corresponding stack set is shown in the right part of Fig. 1. The graph structured stack in the right part is same as the one in the left part except that the

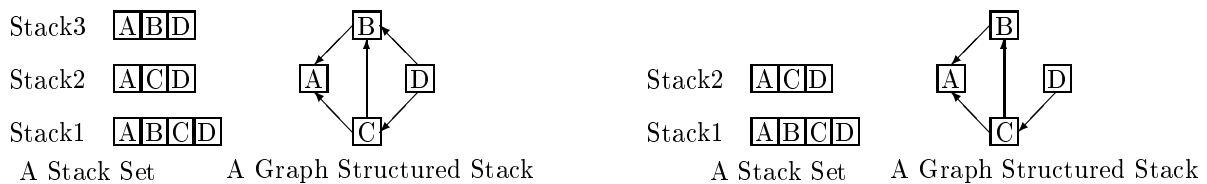


Figure 1: Stack sets and corresponding graph structured stacks

edge from node D to node B is pruned. Pruning this edge corresponds to delete stack 3 from the stack set in the left part.

In this algorithm stack elements are extended items that are pairs of an item and an input string, whereas elements in stack are items and tokens in the algorithm by Nederhof [5]. By adding more information of input strings to items, we can distinguish items generated from different input strings. Thus, we can prune parent pointers in a reduce action, so that the parent pointer set for every stack element does not have more than one element whose items are same and inputs are different. By this pruning we can reduce the time complexity to $O(n^2)$, since the size of a set of parent pointers becomes less than constant (= the size of the set of items), although we cannot obtain all possible parse trees.

In Sect. 2, I introduce definitions and notations used in this report. I state an algorithm based on a set of parse stacks in Sect. 3, and modify it to a graph-structured stack in Sect. 4. The methods in Sect. 3 and 4 seem rather known, and the chief contribution here is to formulate them in formal descriptions. In Sect. 5, by introducing the idea of pruning into the algorithm in Sect. 4, I state the improved algorithm, and its correctness. In Sect. 6 I analyze the time complexity of this algorithm. In Sect. 7 as an example of this algorithm, I show that grammars called S_m ($m \geq 3$) whose time complexity of parsing is $O(n^{m+1})$ in [11] and $O(n^3)$ in [3], can be parsed in $O(n^2)$ by this algorithm. In Sect. 8 some study is made on the range of grammars that can be applicable for this algorithm. In Sect. 9, application of this algorithm for natural language processing is stated. Finally in conclusion, I state characteristics of pruning approach compared with the current approaches to improve the complexity. In Appendix I prove that the parsing algorithms in Sect. 3 and Sect. 4 corresponds exactly to each other which are rather routine.

2 Definitions and Notations

Definition 1 (Context-free grammar). A context-free grammar (CFG) G is a 4-tuple (V_N, V_T, S, P) where V_N is a finite set of nonterminal symbols, V_T a finite set of terminal symbols, $S \in V_N$ a start symbol, and P is a set of productions. For $r \in P$, $\#r$ denotes the length of the right part of a production r , r_j the j -th element in the right part of r . Thus, a production r is expressed as $r_0 : r_1 r_2 \dots r_{\#r}$.

As usual we augment a CFG with the new start symbol S' and a production $S' : S \dashv$. Throughout this report we assume that all CFGs are ε -free.

Example 1. Let $G1$ be a context-free grammar with the following productions.

$$\begin{array}{llll}
 (r^0) S' : S \dashv & (r^1) S : X Y d & (r^2) X : a & (r^3) X : a b \\
 (r^4) Y : Z e & (r^5) Z : c & (r^6) Z : b c &
 \end{array}$$

Parse trees of $G1$ are shown in Fig. 2. For $G1$, ambiguity exists for an input "abcd".

This grammar will be used as a primary example throughout this report.

Definition 2 (item). For $r \in P$ and a nonnegative integer n ($0 \leq n \leq \#r$), an item is a pair $\langle r, n \rangle$, and the set of all items is denoted by I_G . For a production r of the form $r_0 : r_1 r_2 \dots r_{\#r}$,

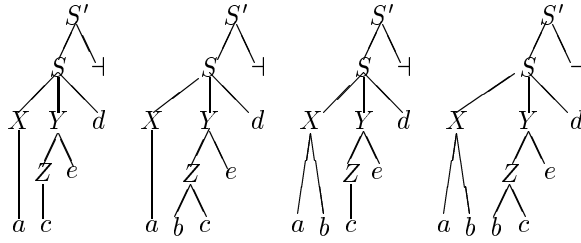


Figure 2: Parse trees for $G1$

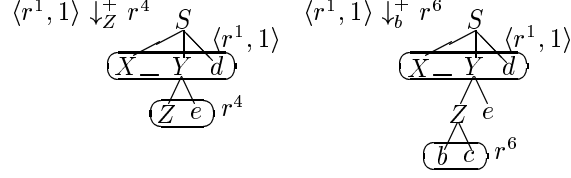


Figure 3: Examples of \downarrow_w^+

an item $\langle r, n \rangle$ can be expressed as $r_0 : r_1 r_2 \dots r_{n-1} r_{n+1} \dots r_{\#r}$. In particular, the item $\langle S' : _ S _ \rangle$ is denoted by i_0 .

Definition 3 (extended item or e-item). For $t \in I_G$ and $\alpha \in V_T^*$, an extended item (or an e-item) is a pair $\langle t, \alpha \rangle$, and the set of all e-items are denoted by J_G . For $x = \langle t, \alpha \rangle \in J_G$, we define $\text{body}(x)$ and $\text{input}(x)$ by $\text{body}(x) = t$ and $\text{input}(x) = \alpha$, respectively. In particular, the e-item $\langle i_0, \varepsilon \rangle$ is denoted by e_0 .

For a sequence η , $\#\eta$ denotes the length of η , η_j the j -th element of η , and $\eta_{\#}$ the last element of η .

Definition 4 ($\downarrow, \downarrow_w^+$). For $\langle r, n \rangle \in I_G$, $r' \in P$ and $w \in V_N \cup V_T$, we define $\langle r, n \rangle \downarrow r'$ if $r_{n+1} = r'_0$. We define $\langle r, n \rangle \downarrow_w^+ r'$ if there exist $r^{m_1}, \dots, r^{m_k} = r' \in P$ ($k \geq 1$) that satisfy the following conditions:

1. $\langle r, n \rangle \downarrow r^{m_1}$,
2. $\langle r^{m_j}, 0 \rangle \downarrow r^{m_{j+1}}$, $r^{m_{j-1}} \neq w$ ($1 \leq j < k$),
3. $r^{m_k} = w$.

Example 2. Let $G1$ be a CFG in example 1. We have $\langle r^1, 1 \rangle \downarrow r^4$, $\langle r^1, 1 \rangle \downarrow_Z^+ r^4$, and $\langle r^1, 1 \rangle \downarrow_b^+ r^6$. We recall that $r^1 = S : XYd$, $r^4 = Y : Ze$, and $r^6 = Z : bc$. The situation for $\langle r^1, 1 \rangle \downarrow_Z^+ r^4$ and $\langle r^1, 1 \rangle \downarrow_b^+ r^6$ are shown in Fig. 3.

Definition 5 ($S0\text{item}_G, S1\text{item}_G$). For a CFG G , $a \in V_T$, we define

$$S0\text{item}_G(a) = \{ \langle r, n \rangle, \alpha \mid r_{n+1} = a \text{ for some } r \in P, \alpha \in V_T^*, 0 \leq n \leq \#r \},$$

$$S1\text{item}_G(a) = \{ \langle r, n \rangle, \alpha \mid \langle r, n \rangle \downarrow_a^+ r' \text{ for some } r, r' \in P, 0 \leq n \leq \#r \}.$$

Hereafter we fix a CFG G and we often omit G in notations like $S0\text{item}$.

Definition 6 (Shift0, Shift1). For $a \in V_T$, $x_0 = \langle r^0, n_0 \rangle, \alpha_0 \in S0\text{item}(a)$, $x_1 = \langle r^1, n_1 \rangle, \alpha_1 \in S1\text{item}(a)$, we define $\text{Shift0}(x_0, a)$ and $\text{Shift1}(x_1, a) \subseteq J_G$ by

$$\text{Shift0}(x_0, a) = \{ \langle r^0, n_0 + 1 \rangle, \alpha_0 \cdot a \},$$

$$\text{Shift1}(x_1, a) = \{ \langle r^1, 1 \rangle, \alpha_1 \cdot a \mid \langle r^1, n_1 \rangle \downarrow_a^+ r' \}.$$

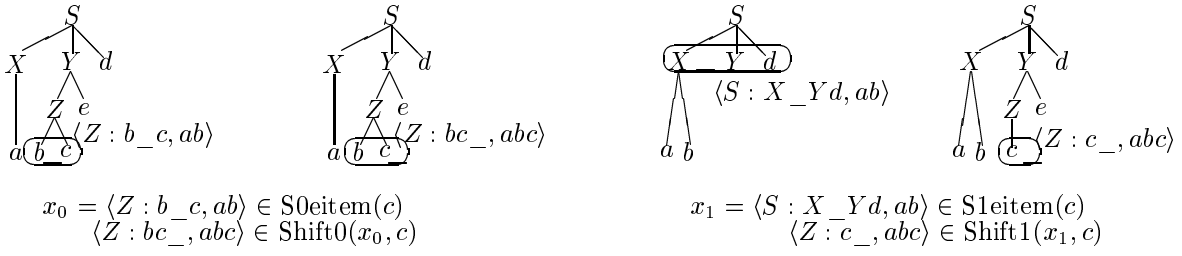


Figure 4: Examples of $\text{S0eitem}(c)$, $\text{Shift0}(x, c)$ and $\text{S1eitem}(c)$, $\text{Shift1}(x, c)$

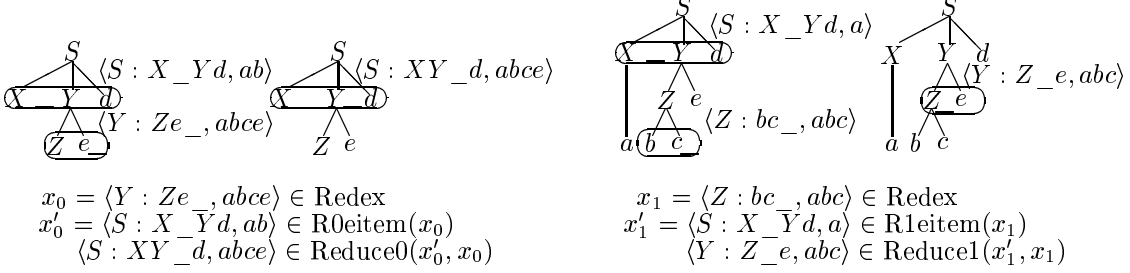


Figure 5: Example of $\text{R0eitem}(x)$, $\text{Reduce0}(x', x)$ and $\text{R1eitem}(x)$, $\text{Reduce1}(x', x)$

Example 3. Let G_1 be a CFG in example 1,

$$\begin{aligned} \langle X : a_b, a \rangle &\in \text{S0eitem}(b), & \text{Shift0}(\langle X : a_b, a \rangle, b) &= \{ \langle X : a b_ , ab \rangle \} \\ \langle S : X_Y d, a \rangle &\in \text{S1eitem}(b), & \text{Shift1}(\langle S : X_Y d, a \rangle, b) &= \{ \langle Z : b_c, ab \rangle \} \end{aligned}$$

The situation of this example is shown in Fig. 4.

Definition 7 (Redex). We define $\text{Redex}_G \subseteq J_G$ by

$$\text{Redex}_G = \{ x \in J_G \mid \text{body}(x) = \langle r, \#r \rangle \text{ for some } r \in P \}.$$

Definition 8 (R0eitem, R1eitem). For $x \in \text{Redex}$ we define $\text{R0eitem}(x)$ and $\text{R1eitem}(x) \subseteq J_G$ by

$$\begin{aligned} \text{R0eitem}(x) &= \{ \langle \langle r', n' \rangle, \alpha \rangle \mid r'_{n'+1} = r_0 \text{ for } r' \text{ and } n' \text{ where } x = \langle \langle r, \#r \rangle, \alpha \rangle \}, \\ \text{R1eitem}(x) &= \{ \langle \langle r', n' \rangle, \alpha \rangle \mid \langle r', n' \rangle \downarrow_{r_0}^+ r'' \text{ for some } r'' \in P \}. \end{aligned}$$

Definition 9 (Reduce0, Reduce1). For $x = \langle \langle r, \#r \rangle, \alpha \rangle$ and $x'_0 = \langle \langle r'^0, n'_0 \rangle, \alpha'_0 \rangle \in \text{R0eitem}(x)$, we define $\text{Reduce0}(x'_0, x) \subseteq J_G$ by

$$\text{Reduce0}(x'_0, x) = \{ \langle \langle r'^0, n'_0 + 1 \rangle, \alpha \rangle \}.$$

For $x = \langle \langle r, \#r \rangle, \alpha \rangle$ and $x'_1 = \langle \langle r'^1, n'_1 \rangle, \alpha'_1 \rangle \in \text{R1eitem}(x)$, we define $\text{Reduce1}(x'_1, x) \subseteq J_G$ by

$$\text{Reduce1}(x'_1, x) = \{ \langle \langle r'', 1 \rangle, \alpha \rangle \mid \langle r'^1, n'_1 \rangle \downarrow_{r_0}^+ r'' \}.$$

Example 4. Let G_1 be a CFG in example 1, $x = \langle Y : Z e_ , abce \rangle$, $x'_0 = \langle S : X_Y d, ab \rangle$

$$x \in \text{Redex}, x'_0 \in \text{R0eitem}(x), \quad \text{Reduce0}(x'_0, x) = \{ \langle S : X Y_d, abce \rangle \}.$$

For $x = \langle Z : b c_ , abc \rangle$, $x'_1 = \langle S : X_Y d, a \rangle$

$$x \in \text{Redex}, x'_1 \in \text{R1eitem}(x), \quad \text{Reduce1}(x'_1, x) = \{ \langle Y : Z_e, abc \rangle \}.$$

The situation of this example is shown in Fig. 5.

3 Parsing Algorithm for Parse Stacks

First we formulate a parsing method similar to [4], but using e-items instead of items. For an input string $\alpha \in V_T^*$ we define a set of sequences of e-items $L^\alpha \subseteq J_G^*$ inductively on substrings of α and an integer n , which can be seen as a parsing for α .

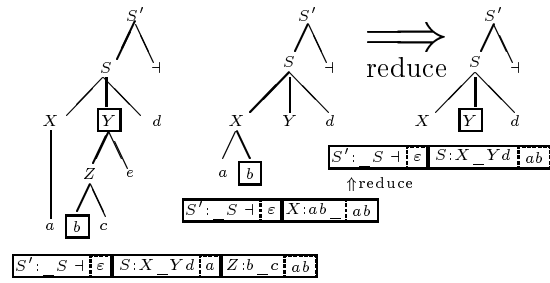


Figure 6: parse stacks in L^{ab} and corresponding parse trees

Example 5. Let $G1$ be a CFG in example 1, parse stacks in L^{ab} and their corresponding parse trees are shown in Fig. 6. In these parse trees paths from e_0 (the root) to b (a leaf) are expressed in thick lines and nodes corresponding to the elements of parse stacks are expressed in boxes.

Algorithm 1 (Parse Stack Parsing).

1. $L^\varepsilon = \{e_0\}$.

In the following case 2 and case 3 we define elements of $L^{\alpha a}$ generated by a shift action from elements of L^α .

2. For $\mu \cdot x \in L^\alpha$, $\mu \cdot y \in L^{\alpha a}$ where $x \in \text{S0item}(a)$ and $y \in \text{Shift0}(x, a)$.
3. For $\mu \cdot x \in L^\alpha$, $\mu \cdot x \cdot y \in L^{\alpha a}$ where $x \in \text{S1item}(a)$ and $y \in \text{Shift1}(x, a)$.

In the following case 4 and case 5 we define elements of $L^{\alpha a}$ generated by a reduce action from elements of $L^{\alpha a}$ by induction on n . The induction base is defined by

$$\text{Reduce}^{(0)}(L^\alpha, a) = \text{Shift0}^{(0)}(L^\alpha, a) \cup \text{Shift1}^{(0)}(L^\alpha, a)$$

where $\text{Shift0}^{(0)}(L^\alpha, a)$ is a set of sequences of e-items $(\mu \cdot y)$ defined in case 2 and $\text{Shift1}^{(0)}(L^\alpha, a)$ is a set of sequences of e-items $(\mu \cdot x \cdot y)$ defined in case 3.

4. For $\mu \cdot x' \cdot x \in \text{Reduce}^{(n)}(L^\alpha, a)$, $\mu \cdot y \in L^{\alpha a}$ where $x \in \text{Redex}$, $x' \in \text{R0item}(x)$ and $y \in \text{Reduce0}(x', x)$.
5. For $\mu \cdot x' \cdot x \in \text{Reduce}^{(n)}(L^\alpha, a)$, $\mu \cdot x' \cdot y \in L^{\alpha a}$ where $x \in \text{Redex}$, $x' \in \text{R1item}(x)$ and $y \in \text{Reduce1}(x', x)$.

Case 4 and case 5 generate $\text{Reduce}^{(n+1)}(L^\alpha, a)$ from $\text{Reduce}^{(n)}(L^\alpha, a)$ where n corresponds to the times of reduce action. $\text{Reduce}^{(n+1)}(L^\alpha, a)$ is defined by

$$\text{Reduce}^{(n+1)}(L^\alpha, a) = \text{Reduce0}^{(n+1)}(L^\alpha, a) \cup \text{Reduce1}^{(n+1)}(L^\alpha, a)$$

where $\text{Reduce0}^{(n+1)}(L^\alpha, a)$ is a set of sequences of the e-items $(\mu \cdot y)$ defined in case 4 and $\text{Reduce1}^{(n+1)}(L^\alpha, a)$ is a set of sequences of the e-items $(\mu \cdot x' \cdot y)$ defined in case 5.

Example 6. Let $G1$ be a CFG in example 1, L^α of $G1$ for an input string “ $abcd$ ” and its substrings are shown in Fig. 7. Each element of L^α is a sequence of e-items and corresponds to a parse stack in [5].

Let $\mu^2 = (e_0)$, $x^2 = \langle X : a _ b, a \rangle$, and $y^2 = \langle X : ab _ , ab \rangle$.

As $\mu^2 \cdot x^2 \in L^\alpha$, $x^2 \in \text{S0item}(b)$ and $y^2 \in \text{Shift0}_E(x^2, b)$, $\mu^2 \cdot y^2 \in L^{ab}$, by case 2 of the stack set parse algorithm.

Let $\mu^3 = (e_0)$, $x^3 = \langle S : X _ Y d, a \rangle$, $y^3 = \langle Z : b _ c, ab \rangle$.

As μ^3 , x^3 , y^3 satisfy $\mu^3 \cdot x^3 \in L^\alpha$, $x^3 \in \text{S1item}(b)$, $y^3 \in \text{Shift1}_E(x^3, b)$, $\mu^3 \cdot x^3 \cdot y^3 \in L^{ab}$ by case 3 of the stack set parse algorithm.

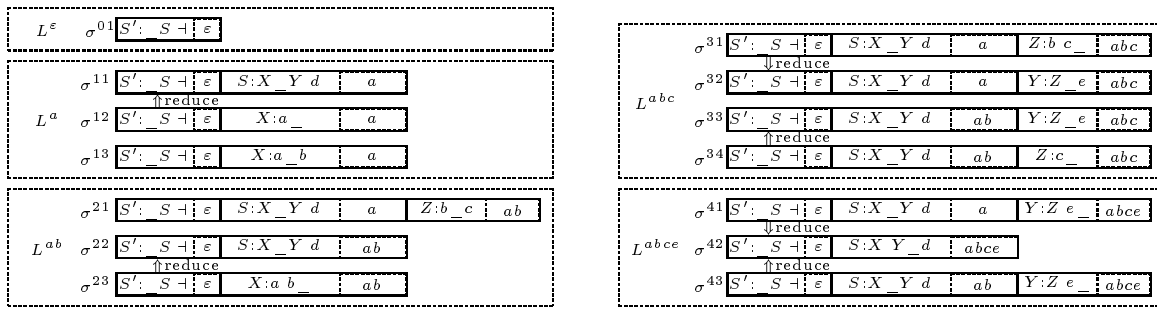


Figure 7: examples of L^α for an input “abce”

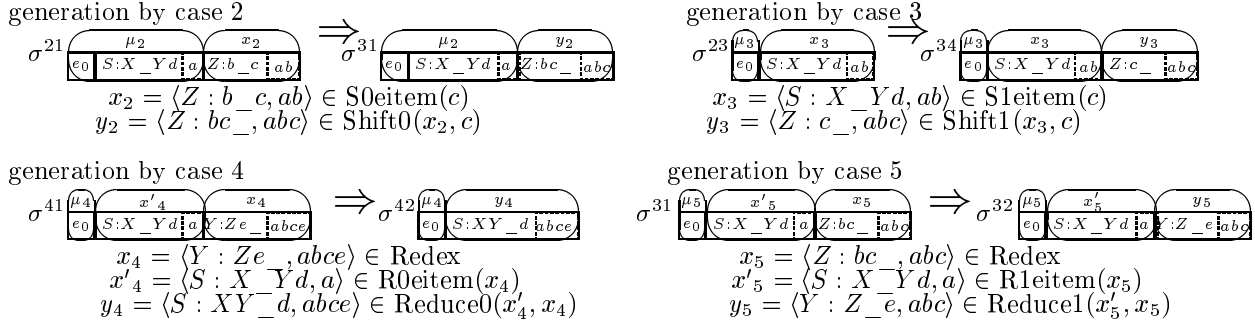


Figure 8: Example of parse stack generation

Let $\mu^4 = (e_0)$, $x'^4 = \langle S : X_Y d, a \rangle$, $x^4 = \langle Y : Z e_ , abce \rangle$, $y^4 = \langle S : X Y_d, abce \rangle$.
As μ^4, x^4, x'^4, y^4 satisfy $\mu^4 \cdot x'^4 \cdot x^4 \in \text{Reduce}^{(0)}_{E^*}(L^{abc}, e)$, $x^4 \in \text{Redex}$, $x'^4 \in \text{R0item}(x^4)$, $y^4 \in \text{Reduce0}_E(x'^4, x^4)$,
 $\mu^4 \cdot y^4 \in L^{abce}$ by case 4 of the stack set parse algorithm.

Let $\mu^5 = (e_0)$, $x'^5 = \langle S : X_Y d, a \rangle$, $x^5 = \langle Z : b c_ , abc \rangle$, $y^5 = \langle Y : Z_e, abc \rangle$.
As μ^5, x^5, x'^5, y^5 satisfy $\mu^5 \cdot x'^5 \cdot x^5 \in \text{Reduce}^{(0)}_{E^*}(L^{ab}, c)$, $x^5 \in \text{Redex}$, $x'^5 \in \text{R1item}(x^5)$, $y^5 \in \text{Reduce1}_E(x'^5, x^5)$,
 $\mu^5 \cdot x'^5 \cdot y^5 \in L^{abc}$ by case 5 of the stack set parse algorithm.

Examples of parse stack generation are shown in Fig. 8.

As shown in the example above, we can get all parse stacks for an input string by starting from L^ϵ as the initial state and construct $L^{\alpha \cdot a}$ from L^α if the current input is α and the next symbol is a .

4 Graph-Structured Stack Parsing

In this section I modify Algorithm 1 utilizing the idea of graph-structured stacks [3], [5], [11]. Each element of a stack has parents, i. e. , those pointed by the element, which are its previous elements in the parse stack. In this algorithm, the parents of a stack element, an e-item x , are a set of e-items, denoted by $\text{Parent}(x)$.

Algorithm 2 (Graph-structured stack parsing).

For $\alpha \in V_T^*$ we define $E_\alpha \subseteq J_G$, and $\text{Parent}(x)$ for $x \in E_\alpha$, inductively on elements of a string α and a positive integer n , as follows:

1. $E_\epsilon = \{e_0\}$, $\text{Parent}(e_0) = \emptyset$.

In the following case 2 and case 3 we define elements of $E_{\alpha \cdot a}$ generated by a shift action from elements of E_α .

2. For $x \in E_\alpha$, $y \in E_{\alpha,a}$ where $x \in \text{S0eitem}(a)$, $y \in \text{Shift0}(x, a)$.

For this y , $\text{Parent}(y)$ is defined by

$$\text{Parent}(y) = \{w \mid w \in \text{Parent}(z) \text{ for } z \text{ where } z \in E_\alpha, z \in \text{S0eitem}(a), \\ y \in \text{Shift0}(z, a)\}$$

3. For $x \in E_\alpha$, $y \in E_{\alpha,a}$ where $x \in \text{S1eitem}(a)$, $y \in \text{Shift1}(x, a)$.

For this y , $\text{Parent}(y)$ is defined by

$$\text{Parent}(y) = \{z \mid z \in E_\alpha, z \in \text{S1eitem}(a), y \in \text{Shift1}(z, a)\}.$$

In the following case 4 and case 5 we define elements of $E_{\alpha,a}$ generated by a reduce action from elements of $E_{\alpha,a}$ by induction on n . The induction base is defined by

$$\text{Reduce}_{(0)}(E_\alpha, a) = \text{Shift0}_{(0)}(E_\alpha, a) \cup \text{Shift1}_{(0)}(E_\alpha, a)$$

where $\text{Shift0}_{(0)}(E_\alpha, a)$ is a set of e -items (y) defined in case 2 and $\text{Shift1}_{(0)}(E_\alpha, a)$ is a set of e -items (y) defined in case 3.

4. For $x \in \text{Reduce}_{(n)}(E_\alpha, a)$, $y \in E_{\alpha,a}$ where $x' \in \text{Parent}(x)$, $x \in \text{Redex}$, $x' \in \text{R0eitem}(x)$, $y \in \text{Reduce0}(x', x)$.

For this y , $\text{Parent}(y)$ is defined by

$$\text{Parent}(y) = \{w \mid w \in \text{Parent}(z') \text{ for } z, z' \text{ where } z \in \text{Reduce}_{(n)}(E_\alpha, a), \\ z' \in \text{Parent}(z), z \in \text{Redex}, z' \in \text{R0eitem}(z), y \in \text{Reduce0}(z', z)\}$$

5. For $x \in \text{Reduce}_{(n)}(E_\alpha, a)$, $y \in E_{\alpha,a}$ where $x' \in \text{Parent}(x)$, $x \in \text{Redex}$, $x' \in \text{R1eitem}(x)$, $y \in \text{Reduce1}(x', x)$.

For this y , $\text{Parent}(y)$ is defined by

$$\text{Parent}(y) = \{z' \mid \text{for } z' \text{ where } z \in \text{Reduce}_{(n)}(E_\alpha, a), \\ z' \in \text{Parent}(z), z \in \text{Redex}, z' \in \text{R1eitem}(z), y \in \text{Reduce1}(z', z)\}.$$

Case 4 and case 5 generate $\text{Reduce}_{(n+1)}(E_\alpha, a)$ from $\text{Reduce}_{(n)}(E_\alpha, a)$ where n corresponds to the times of reduce action. $\text{Reduce}_{(n+1)}(E_\alpha, a)$ is defined by

$$\text{Reduce}_{(n+1)}(E_\alpha, a) = \text{Reduce0}_{(n+1)}(E_\alpha, a) \cup \text{Reduce1}_{(n+1)}(E_\alpha, a)$$

where $\text{Reduce0}_{(n+1)}(E_\alpha, a)$ is a set of e -items (y) defined in case 4 and $\text{Reduce1}_{(n+1)}(E_\alpha, a)$ is a set of e -items (y) defined in case 5.

Example 7. Let $G1$ be a CFG in example 1. The graph structured stack of $G1$ for an input string “ $abcd$ ” is shown in Fig. 9, where elements of $\text{Parent}(x)$ are denoted by an arrow from x to each elements when $\text{Parent}(x) \neq \{e_0\}$. For example, there is an arrow from $x = \langle Y : Z_e, abc \rangle$ to $\langle S : X_Yd, ab \rangle$ since $\text{Parent}(x) = \{\langle S : X_Yd, a \rangle, \langle S : X_Yd, ab \rangle\}$ and there is no arrow from $x' = \langle S : X_Yd, ab \rangle$ to e_0 since $\text{Parent}(x') = \{e_0\}$.

Let $x^2 = \langle X : a_b, a \rangle$, $y^2 = \langle X : a b_ , ab \rangle$.

As x^2 , y^2 satisfy $x^2 \in E_a$, $x^2 \in \text{S0eitem}(b)$, $y^2 \in \text{Shift0}_E(x^2, b)$, $y^2 \in E_{ab}$ and $\text{Parent}(y^2) = \text{Parent}(x^2) = \{e_0\}$ by case 2 of the graph structured stack parse algorithm.

Let $x^3 = \langle S : X_Y d, a \rangle$, $y^3 = \langle Z : b_c, ab \rangle$.

As x^3 , y^3 satisfy $x^3 \in E_a$, $x^3 \in \text{S1eitem}(b)$, $y^3 \in \text{Shift1}_E(x^3, b)$,

$y^3 \in E_{ab}$ and $\text{Parent}(y^3) = \{x_3\}$ by case 3 of the graph structured stack parse algorithm.

Let $x^4 = \langle S : X_Y d, a \rangle$, $x^4 = \langle Y : Z e_ , abce \rangle$, $y^4 = \langle S : X Y_d, abce \rangle$.

As x^4, x'^4, y^4 satisfy $x^4 \in \text{Reduce}^{(0)}_E(E_{abc}, e)$, $x^4 \in \text{Redex}$, $x'^4 \in \text{R0eitem}(x^4)$, $y^4 \in \text{Reduce0}_E(x'^4, x^4)$, $y^4 \in E_{abce}$ and $\text{Parent}(y^4) = \text{Parent}(x'^4) = \{e_0\}$ by case 4 of the graph structured stack parse algorithm.

Let $x^5 = \langle S : X_Y d, a \rangle$, $x^5 = \langle Z : b c_ , abc \rangle$, $y^5 = \langle Y : Z_e, abc \rangle$.

As x^5, x'^5, y^5 satisfy $x^5 \in \text{Reduce}^{(0)}_E(E_{ab}, c)$, $x^5 \in \text{Redex}$, $x'^5 \in \text{R1eitem}(x^5)$, $y^5 \in \text{Reduce1}_E(x'^5, x^5)$,

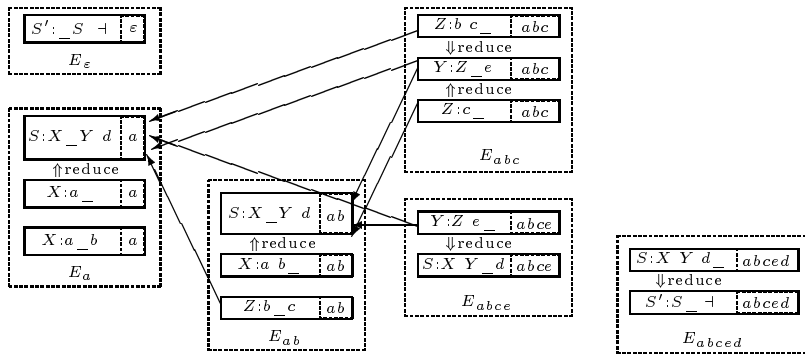
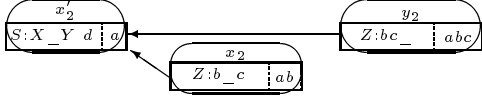


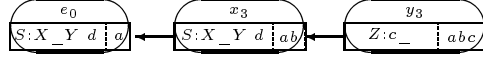
Figure 9: Graph structured stack for the input “abcd”

Generation by case 2



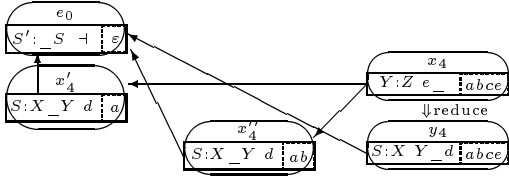
$x_2 \in \text{S0item}(c)$, $\text{Parent}(x_2) = \{x'_2\}$
 $\text{Shift0}(x_2, c) = \{y_2\}$
 $\text{Parent}(y_2) = \text{Parent}(x_2) = \{x'_2\}$

Generation by case 3



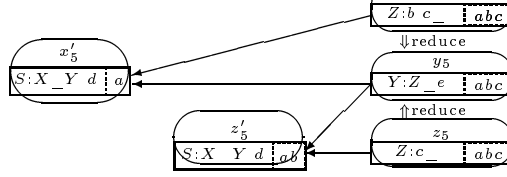
$x_3 \in \text{S1item}(c)$, $\text{Parent}(x_3) = \{e_0\}$
 $\text{Shift1}(x_3, c) = \{y_3\}$, $\text{Parent}(y_3) = \{x_3\}$

Generation by case 4



$x'_4, x''_4 \in \text{R0item}(x_4)$, $\text{Parent}(x_4) = \{x'_4, x''_4\}$
 $\text{Reduce0}(x'_4, x_4) = \text{Reduce0}(x''_4, x_4) = \{y_4\}$
 $\text{Parent}(x'_4) = \text{Parent}(x''_4) = \{e_0\}$
 $\text{Parent}(y_4) = \text{Parent}(x'_4) \cup \text{Parent}(x''_4) = \{e_0\}$

Generation by case 5



$x'_5 \in \text{R1item}(x_5)$, $\text{Parent}(x_5) = \{x'_5\}$
 $z'_5 \in \text{R1item}(z_5)$, $\text{Parent}(z_5) = \{z'_5\}$
 $\text{Reduce1}(x'_5, x_5) = \text{Reduce1}(z'_5, z_5) = \{y_5\}$
 $\text{Parent}(y_5) = \text{Parent}(x_5) \cup \text{Parent}(z_5) = \{x'_5, z'_5\}$

Figure 10: Example of Graph Structured Stack Generation

$y^5 \in E_{abc}$ and $\text{Parent}(y^5) = \{\langle S : X_Y d, a \rangle, \langle S : X_Y d, ab \rangle\}$ by case 5 of the graph structured stack parse algorithm.

Examples of graph structured stack generation are shown in Fig. 10.

The following Lemma 1 holds from definitions of $\text{Shift0}(x, a)$ and $\text{Reduce0}(x', x)$.

Lemma 1. For $a \in V_T$ and $x, y \in J_G$, if $x \in \text{S0item}(a)$ and $y \in \text{Shift0}(x, a)$, then $\text{Parent}(y) = \text{Parent}(x)$.

Proof. If y is expressed as $\langle \langle r, n \rangle, \alpha \cdot a \rangle$, x is expressed as $\langle \langle r, n - 1 \rangle, \alpha \rangle$ and is unique. Therefore, x is the only e-item that satisfies the conditions for z in the definition of $\text{Parent}(y)$ in Case 2 of Algorithm 2. \square

Definition 10 (E^α). For $\alpha \in V_T^*$, $E^\alpha \subseteq J_G^*$ is defined by

$$E^\alpha = \{\sigma \mid \text{Parent}(\sigma_1) = \emptyset, \sigma_j \in \text{Parent}(\sigma_{j+1}) (1 \leq j < \#\sigma), \sigma_{\#\sigma} \in E_\alpha\}.$$

Example 8. For $G1$ of example 1, $E^{ab} = \{\sigma^1, \sigma^2, \sigma^3\}$ where

$$\sigma^1 = (\sigma^1_1, \sigma^1_2) = (e_0, \langle X : a b _ , ab \rangle),$$

$$\sigma^2 = (\sigma^2_1, \sigma^2_2) = (e_0, \langle S : X_Y d, ab \rangle),$$

$$\sigma^3 = (\sigma^3_1, \sigma^3_2, \sigma^3_3) = (e_0, \langle S : X_Y d, a \rangle, \langle Z : b _ c, ab \rangle).$$

Relations between elements of $\sigma^1, \sigma^2, \sigma^3$ are shown in Fig. 11, where every element of $\text{Parent}(x)$ is denoted by an arrow from x to each element. For example, for $\sigma^1_2 = \langle X : ab _ , ab \rangle$, $\text{Parent}(\sigma^1_2) = \{e_0\}$ and there is an arrow from x to e_0 .

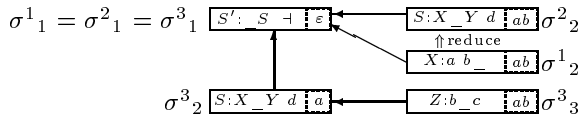


Figure 11: Elements of $\sigma^1, \sigma^2, \sigma^3$

We have the basic theorems about L^α and E^α . They are proved inductively based on definitions of case (1) to case 5 in Algorithm1 and Algorithm2.

Theorem 1. For $\alpha \in V_T^*$, if $\lambda \in L^\alpha$ then $\lambda \in E^\alpha$.

Theorem 2. For $\alpha \in V_T^*$, if $\lambda \in E^\alpha$ then $\lambda \in L^\alpha$.

5 Pruned Graph Structured Stack Parsing

5.1 Outline of Pruned Graph Structured Stack Parsing

We denote the size of a finite set H by $|H|$.

Definition 11 (\mapsto). For e -items $x, y \in J_G$, we define $x \mapsto y$ if one of the following is satisfied:

1. $x \in \text{S0item}(a)$ and $y \in \text{Shift0}(x, a)$,
2. $x \in \text{S1item}(a)$ and $y \in \text{Shift1}(x, a)$,
3. $x \in \text{Redex}$ and $y \in \text{Reduce0}(x', x)$ for some $x' \in \text{R0item}(x)$,
4. $x \in \text{Redex}$ and $y \in \text{Reduce1}(x', x)$ for some $x' \in \text{R1item}(x)$.

The transitive closure of \mapsto is denoted by \mapsto^* .

Example 9. For $G1$ in example 1,

$\langle Z : b_c, ab \rangle \mapsto \langle Z : bc_ , abc \rangle$ holds because $\langle Z : bc_ , abc \rangle$ is generated from $\langle Z : b_c, ab \rangle$ from case 2 of Algorithm 2, as stated in example 6.

Similarly, $\langle S : X_Y d, a \rangle \mapsto^* \langle S : XY_d, abce \rangle$ holds because

$\langle S' : _S \neg, \epsilon \rangle \mapsto \langle X : a_ , a \rangle \mapsto \langle S : X_Y d, a \rangle \mapsto \langle Z : b_c, ab \rangle \mapsto \langle Z : bc_ , abc \rangle \mapsto \langle Y : Z_e, abc \rangle \mapsto \langle Y : Ze_ , abce \rangle \mapsto \langle S : XY_d, abce \rangle \mapsto \langle S' : S_ \neg, abced \rangle$ hold

By definitions of \mapsto^* and $\text{Parent}(x)$, we easily have:

Proposition 1. If $y \in \text{Parent}(x)$ then $y \mapsto^* x$.

Definition 12 (separable). For a set of e -items $D \subseteq J_G$, D is called separable if $\text{body}(x) = \text{body}(y)$ implies $x = y$ for $x, y \in D$.

For a separable set D we can naturally embed D into I_G by ignoring the input parts of the elements. By this embedding we obtain the following easy proposition.

Proposition 2. For a set of e -items $D \subseteq J_G$, if D is separable then $|D| \leq |I_G|$. In particular, a separable set is a finite set.

Following lemma holds for separable parent sets in case 4of Algorithm 2.

Lemma 2. For x, x' and y that satisfy $x \in \text{Redex}$, $x' \in \text{R0item}(x) \cap \text{Parent}(x)$, $y \in \text{Reduce0}(x', x)$, if $\text{Parent}(x)$ is separable then $\text{Parent}(y) = \text{Parent}(x')$.

Proof. If there exists $z' \in \text{R0item}(x) \cap \text{Parent}(x)$ satisfying $y \in \text{Reduce0}(x', x)$, $\text{body}(x') = \text{body}(z')$ holds from $y \in \text{Reduce0}(x', x)$ and $y \in \text{Reduce0}(z', x)$. Therefore $x' = z'$ holds if $\text{Parent}(x)$ is separable. It means that x' that satisfy $y \in \text{Reduce0}(x', x)$ is unique and $\text{Parent}(y) = \text{Parent}(x')$ holds from definition of $\text{Parent}(y)$ in case 4of Algorithm 2. \square

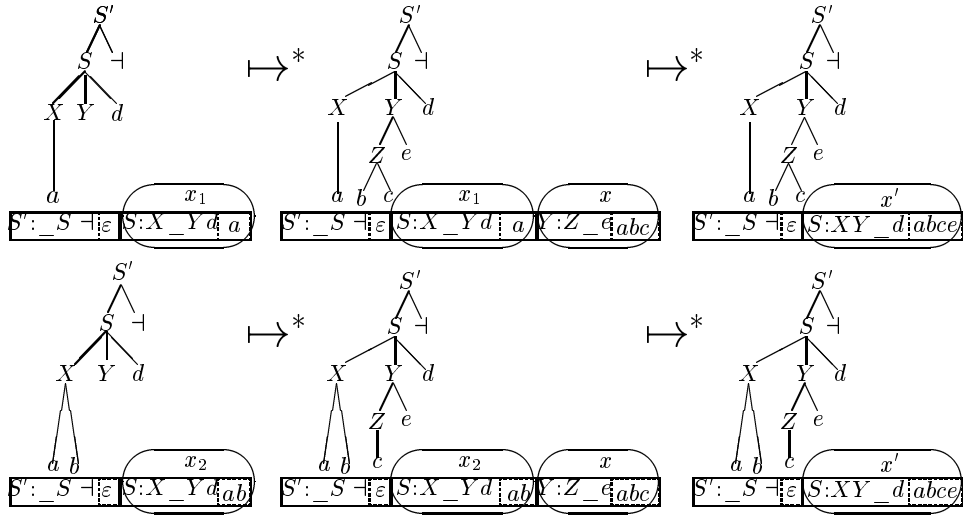


Figure 12: example of Theorem 3: Parse process for “abce”

Unfortunately, all parent sets are not separable.

Example 10. As is seen from Fig. 2, $\text{Parent}(\langle Y : Z_e, abc \rangle)$ is not separable, because it contains $\langle S : X_Yd, a \rangle$ and $\langle S : X_Yd, ab \rangle$.

For parent sets which are not separable, we have the following important theorem.

Theorem 3. *Let $x_1, x_2 \in \text{Parent}(x)$ for an e -item x , and assume that $\text{body}(x_1) = \text{body}(x_2) = \langle r, n \rangle$ and $\text{input}(x_1) \neq \text{input}(x_2)$. If $x \mapsto^* \langle \langle r, n + 1 \rangle, \alpha \rangle$ for some $\alpha \in V_T^*$, then $x_1 \mapsto^* \langle \langle r, n + 1 \rangle, \alpha \rangle$ and $x_2 \mapsto^* \langle \langle r, n + 1 \rangle, \alpha \rangle$.*

Proof. $x_1 \mapsto^* x$ and $x_2 \mapsto^* x$ hold because $x_1, x_2 \in \text{Parent}(x)$. Therefore, $x_1 \mapsto^* x \mapsto^* \langle \langle r, n + 1 \rangle, \alpha \rangle$ and $x_2 \mapsto^* x \mapsto^* \langle \langle r, n + 1 \rangle, \alpha \rangle$ are satisfied. \square

Example 11. Let x_1, x_2, x and x' be $x_1 = \langle S : X_Yd, a \rangle, x_2 = \langle S : X_Yd, ab \rangle, x = \langle Y : Z_e, abc \rangle$ and $x' = \langle S : XY_d, ab \rangle$. In this case $x_1 \mapsto^* x'$ and $x_2 \mapsto^* x'$ holds since $x_1, x_2 \in \text{Parent}(x)$ and $x \mapsto^* x'$.

This situation is shown in Fig. 12, in which the upper part indicates a parse process corresponding to $x_1 \mapsto^* x \mapsto^* x'$ and the lower part indicates a parse process corresponding to $x_2 \mapsto^* x \mapsto^* x'$. As $\text{Parent}(x)$ is not separable, two different parse trees exist for an input “abc” that are indicated in the center of Fig. 12. The right part of Fig. 12 shows that after the reduction of “ $Y : Ze$ ”, there exists 2 different parse trees but their parse stack is same because the contents of parse stack do not depend on the substructures of X and Y after the reduction of productions for X and Y .

As top elements of stacks in the middle part are same (x), moves of both stacks are same until their top elements are reduced. After top elements are reduced, elements under top elements (x_1, x_2) are shifted to a same element (x') and both stacks become same. Therefore moves of both stacks are same after their top elements is reduced that means moves of both stacks in the middle part are same for any input string. It means that we can prune one of stacks in the middle part.

5.2 Definition of Pruned Graph Structured Stack Parsing

Definition 13 (E). *We define E as $E = \cup_{\alpha \in V_T^*} E_\alpha$.*

Definition 14 (I^α, E^x, I^x). For $\alpha \in V_T^*, x \in E_\alpha$, we define

$$I^\alpha = \{\rho \mid \exists \sigma \in E^\alpha \text{ where } \# \rho = \# \sigma, \rho_j = \text{body}(\sigma_j) \ (1 \leq j \leq \# \rho)\}.$$

$$E^x = \{\sigma \mid \sigma \in E^\alpha, \sigma_{\#} = x\}.$$

$$I^x = \{\rho \mid \rho \in I^\alpha, \rho_{\#} = x\}.$$

Elements of I^α are body parts of elements of E^α and E^x is a subset of E^α whose element σ satisfies $\sigma_{\#} = x$.

Example 12. For $G1$ in example 1, $I^{abc} = \{\rho^1, \rho^2, \rho^4\}$ where

$$\rho^1 = (\rho^1_1, \rho^1_2, \rho^1_3) = (i_0, S : X_Yd, Z : bc_).$$

$$\rho^2 = (\rho^2_1, \rho^2_2, \rho^2_3) = (i_0, S : X_Yd, Y : Z_e),$$

$$\rho^4 = (\rho^4_1, \rho^4_2, \rho^4_3) = (i_0, S : X_Yd, Z : c_).$$

For $x = \langle Y : Z_e, abc \rangle$, $I^x = \{\rho^2\}$.

From the definition of I^x , we have

Lemma 3. For $x \in E$,

$$I^x = \bigcup_{y \in \text{Parent}(x)} \{\rho \cdot \text{body}(x) \mid \rho \in I^y\}.$$

From proposition 11 in Sec 5.3, we can prune x if $I^x \subseteq I^y$ holds for $x, y \in \text{Parent}(z)$. We define a binary relation that shows $I^x \subseteq I^y$ over $x, y \in \text{Parent}(z)$ that satisfy $\text{body}(x) = \text{body}(y)$.

Definition 15 (\prec). For $x, y \in E$ that satisfy $\text{body}(x) = \text{body}(y)$,

we define $x \prec y$ if the following condition is satisfied,

$$\forall x' \in \text{Parent}(x), \exists y' \in \text{Parent}(y) . \quad x' = y' \text{ or } x' \prec y'$$

Proposition 3. If $x \prec y$ then $I^x \subseteq I^y$.

Proof. As context free grammars in this report are ε -free, we have (A),(B) from Algorithm 2.

(A) For $x, y \in E$, if $x \in \text{Parent}(y)$ then $\# \text{input}(x) < \# \text{input}(y)$.

(B) For $x \in E$ if $\# \text{input}(x)=0$ then $x = e_0$. Let $L(x, y)$ be $L(x, y)=\max(\# \text{input}(x), \# \text{input}(y))$ and we prove this proposition by induction on $L(x, y)$.

1. $L(x, y) = 1$ As $\# \text{input}(x)=\# \text{input}(y)=1$ holds in this case, $\text{Parent}(x) = \text{Parent}(y) = \{e_0\}$ holds from (A), (B) and this proposition is satisfied.
2. We assume that this proposition is satisfied for $L(x, y) \leq k$ and prove this proposition for $L(x, y) = k + 1$. From $x \prec y$, for any $x' \in \text{Parent}(x)$ there exists $y' \in \text{Parent}(y)$ that satisfies $x' = y'$ or $x' \prec y'$.

For case of $x' = y'$, $I^{x'} \subseteq I^{y'}$ holds.

For case of $x' \prec y'$, $\# \text{input}(x') \leq k, \# \text{input}(y') \leq k$ holds from $x' \in \text{Parent}(x), y' \in \text{Parent}(y)$ and (A). Therefore $L(x', y') \leq k$ holds and $I^{x'} \subseteq I^{y'}$ holds from the induction hypothesis.

In both cases $I^{x'} \subseteq I^{y'}$ holds. From lemma3, for any $\rho \in I^x$, there exist $x' \in \text{Parent}(x)$, $\rho' \in I^{x'}$ that satisfy $\rho = \rho' \cdot \text{body}(x)$. From $I^{x'} \subseteq I^{y'}$, $\rho' \in I^{y'}$ holds. This proposition holds from $\rho = \rho' \cdot \text{body}(x) = \rho' \cdot \text{body}(y) \in I^y$.

□

Definition 16 (representative, representable). For $D \subseteq E$, $x \in D$ and $D_x = \{y \mid y \in D, \text{body}(y) = \text{body}(x)\}$, $p_x \in D_x$ is called the representative of x if and only if $z \prec p_x$ holds for any $z \in D_x$. D is called representable if and only if the representative of x exists for any $x \in D$.

D_x is a set of elements whose body parts are same as the body part of x and the representative of x is the unpruned element of D_x . If D is representable, we can prune elements of D to their representative.

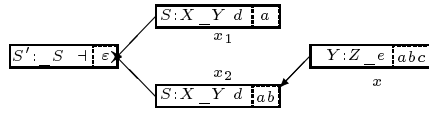


Figure 13: Prune(Parent((Y : Z_e, abc))) in G1

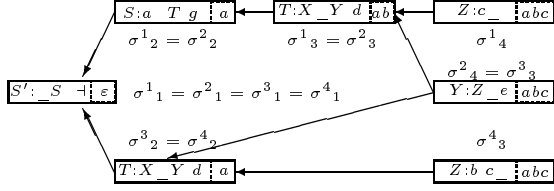


Figure 14: Unprunable case in G2

Example 13. For x, x_1, x_2 in example 11, x_1 and x_2 are both representative of x because $\text{Parent}(x) = \{x_1, x_2\}$, $x_1 \prec x_2, x_2 \prec x_1$ holds. The case for $\text{Prune}(\text{Parent}(x)) = \{x_2\}$ is shown in Fig. 13.

From proposition 3 and Definition 16, we have

Proposition 4. *If p_x and p'_x are representatives of x then $I^{p_x} = I^{p'_x}$.*

Definition 17 (Prune). *For a representable set $D \subseteq E$, $\text{Prune}(D)$ is a set whose elements are one of representative for each element of D .*

As shown in example 13, there may be more than one representatives for one element of D . However we can choose any one of them as the element of $\text{Prune}(D)$ from proposition 4.

From definition 17, we have

Proposition 5. *For $D \subseteq E$, if D representable then $\text{Prune}(D)$ is separable.*

We define context free grammars for this parsing algorithm.

Definition 18 (prunable). *A context free grammar G is called prunable if and only if $\text{Parent}(x)$ is representable for any $x \in E$.*

Example 14. $G1$ is prunable because $\text{Parent}(x)$ is representable from example 13.

Example 15. Let $G2$ be a context free grammar with following productions.

$$\begin{aligned} (r^0) S' : S \rightarrow | \epsilon & \quad (r^1) S : T f & \quad (r^2) S : a T g & \quad (r^3) T : X Y d \\ (r^4) X : a & \quad (r^5) X : b & \quad (r^6) Y : Z e & \quad (r^7) Z : c & \quad (r^8) Z : b c \end{aligned}$$

Relations of elements of $E^{abc} = \{\sigma^1, \sigma^2, \sigma^3, \sigma^4\}$ are shown in Fig. 14. From Fig. 14, neither $\langle T : X_Y d, a \rangle \prec \langle T : X_Y d, ab \rangle$ nor $\langle T : X_Y d, ab \rangle \prec \langle T : X_Y d, a \rangle$ hold. Therefore $G2$ is not prunable because $\text{Parent}(\langle Y : Z_e, abc \rangle)$ is not representable.

Algorithm 3 (Pruned Graph-Structured Stack Parsing). *An algorithm in which "Parent(y) = $P^l(y)$ " in case 5 of Algorithm 2 is modified to "Parent(y) = Prune($P^l(y)$)" is called a Pruned Graph-Structured Stack Parsing algorithm.*

Proposition 6. *For x, y in Algorithm 2, if $\text{Parent}(x)$ is separable then $\text{Parent}(y)$ defined in case 2, 3, 4 is separable.*

Proof. In Case 2, $\text{Parent}(x) = \text{Parent}(y)$ from lemma 1 and $\text{Parent}(y)$ is separable because $\text{Parent}(x)$ is separable. In Case 3, $\text{Parent}(y) \subseteq E_\alpha$ and $\text{Parent}(y)$ is separable because E_α is separable. In Case 4, $\text{Parent}(x)$ is separable and $\text{Parent}(y)$ is separable from lemma. \square

From this proposition, $\text{Parent}(y)$ defined in case 2, 3, 4, 5 of Pruned Graph-Structured Stack Parsing are separable.

From proposition 3, definition 16, definition 17, we have

Theorem 4. For $x \in E$ satisfying D is representable,

$$\bigcup_{y \in \text{Parent}(x)} I^y = \bigcup_{y \in \text{Prune}(\text{Parent}(x))} I^y.$$

From this theorem, $I^x = \bigcup_{y \in \text{Parent}(x)} I^y$ is conserved if $\text{Parent}(x)$ is modified to $\text{Prune}(\text{Parent}(x))$. It means that the result of Pruned Graph-Structured Stack Parsing in which $\text{Parent}(x)$ of Graph-Structured Stack Parsing is modified to $\text{Prune}(\text{Parent}(x))$ is same as the result of Graph-Structured Stack Parsing.

5.3 Correctness of the Pruned Graph-Structured Stack Parsing

In this section, I show the correctness of Pruned Graph-Structured Stack Parsing by showing that the result (whether the input is a word of the grammar) of Graph-Structured Stack Parsing is not changed if $\text{Parent}(x)$ is modified to $\text{Prune}(\text{Parent}(x))$.

From definition 13, we have

Proposition 7. For $\alpha' \in V_T^*$

α is a word of G

if and only if $\langle i_0, \varepsilon \rangle = \langle S' : _ S \dashv, \varepsilon \rangle \mapsto^* \langle S' : S _ \dashv, \alpha' \rangle = \langle \text{Succ}(i_0), \alpha' \rangle$ holds

Proposition 8. For $\alpha \in V_T^*$,

there exists $\beta \in V_T^*$ that satisfies $\alpha \cdot \beta$ is a word of G

if and only if $\sigma \in E^\alpha$ that satisfy $\langle i_0, \varepsilon \rangle \mapsto^* \sigma_{\#} \mapsto^* \langle \text{Succ}(i_0), \alpha \cdot \beta \rangle$ exists.

Example 16. For $\alpha = ab, \beta = ced$,

$\alpha \cdot \beta = abced$ is a word of $G1$ and $\sigma^2 \in E^{ab} = E^\alpha$ in example 8 satisfies

$\langle i_0, \varepsilon \rangle = \langle S' : _ S \dashv, \varepsilon \rangle \mapsto^* \sigma_{\#}^2 = \langle S : X _ Yd, ab \rangle \mapsto^* \langle S' : S _ \dashv, abced \rangle = \langle \text{Succ}(i_0), \alpha \cdot \beta \rangle$ from example 9.

For σ that satisfies condition of proposition 8, every elements(σ_j) of σ will be popped while reading and parsing β . After σ_j will be popped, $\text{body}(\sigma_{j-1})$, which is the body part of an element of the parent set of σ_j will shift to $\text{Succ}(\text{body}(\sigma_{j-1}))$. After this shifted element will be popped, the body part of σ_{j-2} will shift to $\text{Succ}(\text{body}(\sigma_{j-2}))$ and so on. From this, we have

Proposition 9. For $\alpha \in V_T^*$,

there exists $\beta \in V_T^*$ that satisfy $\alpha \cdot \beta$ is a word of G if and only if there exist $\sigma \in E^\alpha, \beta_j \in V_T^*, \sigma^j \in E^{\alpha \cdot \beta_j}$ ($1 \leq j \leq \# \sigma - 1$) that satisfy following conditions.

1. $\# \sigma^j = \# \sigma - j$
2. $\sigma^j_k = \sigma_k$ ($1 \leq k \leq \# \sigma^j - 1 = \# \sigma - j - 1$)
3. $\sigma^j_{\#} = \langle \text{Succ}(\text{body}(\sigma_{\# \sigma - j})), \alpha \cdot \beta_j \rangle$
4. $\beta_{\# \sigma - 1} = \beta$

Example 17. For $\alpha = abc \in V_T^*$,

there exists $\beta = ed \in V_T^*$ that satisfy $\alpha \cdot \beta$ is a word of $G1$.

For σ^{21}, σ^{42} in Fig. 7

$\sigma = \sigma^{21}, \beta_1 = e, \beta_2 = \beta = ed$,

$\sigma^1 = \sigma^{42} \in E^{abce} = E^{\alpha \cdot \beta_1}, \sigma^2 = (\langle S' : S _ \dashv, abced \rangle) \in E^{abced} = E^{\alpha \cdot \beta_2}$

satisfy conditions of proposition 9. Relations between elements of $\sigma, \sigma^1, \sigma^2$ are shown in Fig. 15.

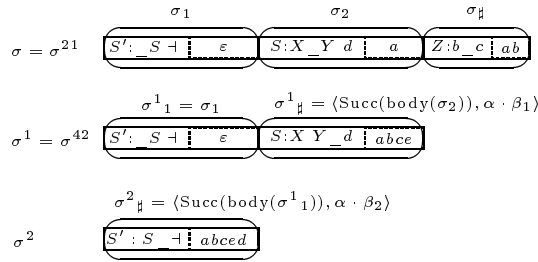


Figure 15: example of proposition.9 in $G1$

Between top elements of σ and σ^j ,
 $\sigma_{\#} \mapsto^* \sigma^1_{\#} = \langle \text{Succ}(\text{body}(\sigma_{\#\sigma-1})), \alpha \cdot \beta_1 \rangle$,
 $\langle \text{Succ}(\text{body}(\sigma_{\#\sigma-j})), \alpha \cdot \beta_j \rangle = \sigma^j_{\#} \mapsto^* \sigma^{j+1}_{\#} = \langle \text{Succ}(\text{body}(\sigma_{\#\sigma-j-1})), \alpha \cdot \beta_{j+1} \rangle$ holds for $1 \leq j \leq \#\sigma - 2$ and we have

Proposition 10. For $\alpha \in V_T^*$,

there exists $\beta \in V_T^*$ that satisfy $\alpha \cdot \beta$ is a word of G

if and only if $\sigma \in E^\alpha, \beta_j \in V_T^*$ ($1 \leq j < \#\sigma$) that satisfy following conditions

1. $\sigma_{\#} \mapsto^* \langle \text{Succ}(\text{body}(\sigma_{\#\sigma-1})), \alpha \cdot \beta_1 \rangle$.
2. $\langle \text{Succ}(\text{body}(\sigma_{\sigma-j})), \alpha \cdot \beta_j \rangle \mapsto^* \langle \text{Succ}(\text{body}(\sigma_{\sigma-j-1})), \alpha \cdot \beta_{j+1} \rangle$ ($1 \leq j \leq \#\sigma - 2$).

We have a following proposition because conditions of above proposition depend on the body part of σ .

Proposition 11. For $\alpha, \beta \in V_T^*, \sigma \in E^\alpha$ that satisfy the following condition

$\langle i_0, \varepsilon \rangle \mapsto^* \sigma_{\#} \mapsto^* \langle \text{Succ}(i_0), \alpha \cdot \beta \rangle$

and for $x, y \in \text{Parent}(\sigma_{\#})$ that satisfy $I^x = I^y$,

if there exist $\gamma \in V_T^*, \sigma_x \in E^x$ that satisfy

$\sigma_{\#} \mapsto^* \langle \text{Succ}(\text{body}((\sigma_x)_{\#})), \alpha \cdot \gamma \rangle \mapsto^* \langle \text{Succ}(i_0), \alpha \cdot \beta \rangle$,

there exist $\sigma_y \in E^y$ that satisfy

$\sigma_{\#} \mapsto^* \langle \text{Succ}(\text{body}((\sigma_y)_{\#})), \alpha \cdot \gamma \rangle \mapsto^* \langle \text{Succ}(i_0), \alpha \cdot \beta \rangle$.

From proposition 11, let x, y be elements in $\text{Parent}(z)$ that satisfy $I^x = I^y$, for any $\sigma_x \in E^x$ there exists $\sigma_y \in E^y$ whose parse process is same as σ_x for all input.

Example 18. For x, x_1, x_2 in example 11,

let $\sigma_{x_1}, \sigma_{x_2}$ be $\sigma_{x_1} = (e_0, x_1), \sigma_{x_2} = (e_0, x_2)$ then $E^{x_1} = \{\sigma_{x_1}\}, E^{x_2} = \{\sigma_{x_2}\}$ hold. As $I^{x_1} = I^{x_2}$ holds,

$x = \sigma_{\#} \mapsto^* \langle \text{Succ}(\text{body}((\sigma_{x_1})_{\#})), abc \cdot e \rangle \mapsto^* \langle \text{Succ}(i_0), abc \cdot ed \rangle$ and

$x = \sigma_{\#} \mapsto^* \langle \text{Succ}(\text{body}((\sigma_{x_2})_{\#})), abc \cdot e \rangle \mapsto^* \langle \text{Succ}(i_0), abc \cdot ed \rangle$ hold for $\sigma \in E^x$.

6 Complexity of Pruned Graph-Structured Stack Parsing

The modified algorithm Pruned Parsing is shown in Algorithm 4 and Algorithm 5. In this pseudo-code the scopes of nesting structures such as loops are specified by indentation. In this algorithm Parent sets are expressed by an array whose index is an e-item and its component is a set of e-items, and Reduce is an array whose n -th component ($\text{Reduce}[n]$) corresponds to $\text{Reduce}^{(n)}(E_\alpha, a)$ in Algorithm 2. In Is_prec of Algorithm 5, duplicate calculation is suppressed by a table PrecTab . The initial value of elements of PrecTab is 0 and for $x, y \in E$ satisfying $\text{body}(x) = \text{body}(y) = i$, the value of $\text{PrecTab}[i, \#\text{input}(x), \#\text{input}(y)]$ is 1 if $x \prec y$ holds, -1 if $x \prec y$ does not hold and 0 otherwise.

I analyze the time complexity of Pruned Parsing.

```

Algorithm 4 (Pruned Graph-Structured Stack Parsing a)
Parse( $\alpha$ )/*returns true if alpha is a word, false otherwise*/
(1)  $E_{top} := \{\epsilon_0\}$ ; Parent[ $\epsilon_0$ ] :=  $\emptyset$ ;
(2) for  $i \in I_G$ 
(3)   for  $j$  from 0 to  $\# \alpha$  for  $k$  from 0 to  $\# \alpha$ 
(4)     PrecTab[ $i, j, k$ ] := 0;
(5) for  $j$  from 1 to  $\# \alpha$ 
(6)    $E_{top} := \text{Goto}(E_{top}, \alpha_j)$ ;
(7)    $E_{top} := \text{Goto}(E_{top}, \uparrow)$ ;
(8) return ( $E_{top} \neq \emptyset$ );

Goto( $E_\alpha, a$ ) /* Generate  $E_{\alpha.a}$  from  $E_\alpha$  and  $a$  */
(9) Reduce_base :=  $\emptyset$ ; PrecTab[ $i, \# \alpha + 1, j$ ] := FALSE;
(10) for  $x \in E_\alpha$ 
(11)   if  $x \in S0\text{item}(a)$  then /* For Case 2 */
(12)     E_shift0 := Shift0( $x, a$ );
           Reduce_base := Reduce_base  $\cup$  E_shift0;
(13)   for  $y \in \text{E\_shift0}$  /* Parent[y] in Case 2 */
(14)     Parent[y] := Parent[x]; /* From lemma 1 */
(15)   if  $x \in S1\text{item}(a)$  then /* For case 3 */
(16)     E_shift1 := Shift1( $x, a$ );
           Reduce_base := Reduce_base  $\cup$  E_shift1;
(17)   for  $y \in \text{E\_shift1}$  /* Parent[y] in Case 3 */
(18)     Parent[y] :=  $\emptyset$ ;
(19)   for  $z \in E_\alpha$ 
(20)     if  $z \in S1\text{item}(a) \wedge y \in \text{Shift1}(z, a)$  then
(21)       Parent[y] := Parent[y]  $\cup$  { $z$ };
(22) return Rclosure(Reduce_base);

Rclosure(D) /* case 4, 5 of Algorithm 3 */
(23)  $n := 0$ ; Reduce[0] := D; Reduce[1] :=  $\emptyset$ ; E_return := Reduce[0];
(24) loop /* The exit of this loop ((24) - (42)) is (42) */
(25)   for  $x \in \text{Reduce}[n] \cap \text{Redex}$ 
(26)     for  $x' \in \text{Parent}[x]$ 
(27)       if  $x' \in R0\text{item}(x)$  then /* For case 4 */
(28)         E_reduce0 := Reduce0( $x', x$ );
           Reduce[ $n + 1$ ] := Reduce[ $n + 1$ ]  $\cup$  E_reduce0;
(29)       for  $y \in \text{E\_reduce0}$  /* Parent[y] in case 4 */
(30)         Parent[y] := Parent[x']; /* From lemma 2 */
(31)       if  $x' \in R1\text{item}(x)$  then /* For case 5 */
(32)         E_reduce1 := Reduce1( $x', x$ );
           Reduce[ $n + 1$ ] := Reduce[ $n + 1$ ]  $\cup$  E_reduce1;
(33)       for  $y \in \text{E\_reduce1}$  /* Parent[y] in case 5 */
(34)         Parent[y] :=  $\emptyset$ ;
(35)       for  $z \in \text{Reduce}[n] \cap \text{Redex}$ 
(36)         for  $z' \in \text{Parent}[z] \cap R1\text{item}(z)$ 
(37)           if  $y \in \text{Reduce1}(z', z)$  then
(38)             Parent[y] := Parent[y]  $\cup$  { $z'$ };
(39)           Parent[y] := Prune(Parent[y]);
(40)         if Reduce[ $n + 1$ ]  $\neq \emptyset$ 
(41)           then E_return := E_return  $\cup$  Reduce[ $n + 1$ ];
                 $n := n + 1$ ; Reduce[ $n + 1$ ] :=  $\emptyset$ ;
(42)         else return E_return; /* exit for loop ((24)-(42)) */

Algorithm 5 (Pruned Graph-Structured Stack Parsing b)
Prune(D)
(43) D_result :=  $\emptyset$ ;
(44) for  $i \in I_G$  D_x[ $i$ ] :=  $\emptyset$ ;
(45) for  $x \in D$ 
(46)   D_x[body( $x$ )] := D_x[body( $x$ )]  $\cup$  { $x$ };
(47) for  $i \in I_G$ 
(48)   if |D_x[ $i$ ] | = 1 then
(49)     D_result := D_result  $\cup$  D_x[ $i$ ];
(50)   else if |D_x[ $i$ ] | > 1 then
(51)      $p_x :=$  any element of D_x[ $i$ ];
(52)     for  $x \in \text{D\_x}[i]$  /* choose representatives */
(53)       if Is_prec( $p_x, x$ ) then
(54)          $p_x := x$ ;
(55)     for  $x \in \text{D\_x}[i]$  /* error if  $p_x$  is NOT representative */
(56)       if not Is_prec( $x, p_x$ )
(57)         then Error /* This grammar is NOT prunable */
(58)     D_result := D_result  $\cup$  { $p_x$ };
(59) return D_result;

Is_prec( $x, y$ ) /* returns true if  $x < y$  holds, false otherwise */
(60)  $i := \text{body}(x)$ ;  $l_x := \# \text{input}(x)$ ;  $l_y := \# \text{input}(y)$ ;
(61) if PrecTab[ $i, l_x, l_y$ ] = 1 then /*  $x < y$  holds */
(62)   return True;
(63) if PrecTab[ $i, l_x, l_y$ ] = -1 then /*  $x < y$  does NOT hold */
(64)   return False;
(65) flag := True;
(66) for  $x' \in \text{Parent}[x]$  (67)   x_flag := False;
(68)   for  $y' \in \text{Parent}[y]$ 
(69)     x_flag := x_flag  $\vee$  ( $x' = y'$ )
            $\vee$  ((body( $x'$ ) = body( $y'$ ))  $\wedge$  Is_prec( $x', y'$ ))
(70)   flag := flag  $\wedge$  x_flag;
(71) if flag then PrecTab[ $i, l_x, l_y$ ] := 1;
(72)   else PrecTab[ $i, l_x, l_y$ ] := -1;
(73) return flag;

```

Figure 16: Pruned Graph-Structured Stack Parsing

Definition 19 (degenerated). For $D \subseteq J_G$, D is called *degenerated* if $\text{body}(x) = \text{body}(y)$ for $\forall x, y \in D$.

From definitions 12 and 14, we have

Lemma 4. Let $D, D_1, D_2 \subseteq J_G$. If D is *degenerated* then D is *separable*.

Lemmas 5 and 6 state that the time complexity of basic operations such as computation of $\text{Shift0}(x, a)$ and assignment in this algorithm are constant in n .

Lemma 5. For $x, x' \in J_G$ and $a \in V_T$, the sets $\text{Shift0}(x, a)$, $\text{Shift1}(x, a)$, $\text{Reduce0}(x', x)$ and $\text{Reduce1}(x', x)$ are *degenerated* and computed in constant time in n (exactly, $O(|I_G|)$).

Proof. Let y be an element of these 4 sets. From Definitions 6 and 9, $\text{body}(y)$ depends only on $\text{body}(x')$ and $\text{body}(x)$. As a body part of an e-item is an element of the finite set I_G , we can compute the body part of these 4 sets for each $\text{body}(x)$ and $\text{body}(x')$ before parsing, and we can get them in constant time in parsing. From Definitions 6 and 9, for each of these 4 sets, all elements of the set have the same input part. Thus all these sets are *degenerated*.

Since each of the 4 sets is *degenerated*, the number of its elements is at most $O(|I_G|)$. The body part of each element is obtained before parsing, its input part is identical, so that the set can be decided in $O(|I_G|)$. \square

Lemma 6. For $x \in J_G$ and the sets of e-items $D, D' \subseteq J_G$, if D is *separable* then both the time complexities of membership $x \in D$ and assignment $D' := D$ are constant in n .

Lemmas 7, 8 and 9 show that some of the important sets appearing in Pruned Parsing are *separable*. From Lemma 6 we obtain the following lemma.

Lemma 7. If D is *degenerated* then the sets $\text{Rclosure}(D)$ and $\text{Goto}(D, a)$ in Pruned Parsing are *degenerated*.

As the set $\{e_0\}$ is *degenerated* we obtain the following lemma from Lemma 7.

Lemma 8. E_top and E_alpha in Pruned Parsing is always *separable*.

Lemma 9. For $y \in E_top$, $\text{Parent}(y)$ is *separable*.

Proof. For $y \in E_top$, one of the following is satisfied:

- a) $y \in \text{Shift0}(x, a)$ for some $x \in E_alpha$,
- b) $y \in \text{Shift1}(x, a)$ for some $x \in E_alpha$,
- c) $y \in \text{Reduce0}(x', x)$ for some $x \in E_alpha, x' \in \text{Parent}(x)$,
- d) $y \in \text{Reduce1}(x', x)$ for some $x \in E_alpha, x' \in \text{Parent}(x)$.

For a), $\text{Parent}(y)$ is *separable* from induction hypothesis, because $\text{Parent}(y) := \text{Parent}(x)$ at (13). For b), this lemma holds because $\text{Parent}(y) = \{x' \mid x' \in E_alpha \cap \text{S1eitem}(a), y \in \text{Shift1}(x', a)\}$ from (16) – (20), therefore $\text{Parent}(y) \subseteq E_alpha$. For c), this lemma holds from the induction hypothesis, because $\text{Parent}(y) := \text{Parent}(x')$ at (30). For d), the element of $\text{Parent}(y)$ is inserted by PruneInsert at (39) that returns a *separable* set. \square

Hereafter we denote the time complexity of a function f as $|f|$.

Proposition 12. The time complexity of Algorithm 4 is $O(n^2)$ except for the time complexity of Prune .

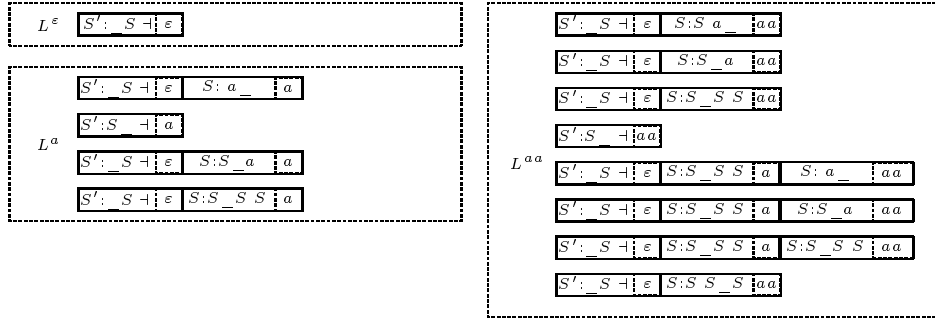


Figure 17: example of L^α in $S3$

Proof. The time complexity of the pruned parsing is $|\mathbf{Parse}|$ that is equal to $n \times |\mathbf{Goto}|$ from the loop of (2) where n is the length of α . As operations in loops of (10), (13), (17) and (19) are executed less than constant ($|I_G|$) time from lemmas 5 and 9, and the complexity of these operations are constant from lemmas 4 and 7. Therefore, $|\mathbf{Goto}|$ is equal to $|\mathbf{Rclosure}|$. The complexity of operations in the loop (24) is also constant, and the number of execution of these operations in the loop is equal to the number of reductions for each input that is at most n . Therefore, $|\mathbf{Rclosure}|$ is equal to $O(n)$. From these considerations the time complexity of this algorithm is $O(n^2)$. \square

Proposition 13. *The time complexity of Prune is $O(n^2)$.*

Proof. As $\mathbf{Reduce}[n]$ in (35) and $\mathbf{Parent}[z]$ in (36) are separable, the number of $\mathbf{Parent}[y]$ which is a parameter of \mathbf{Prune} at (39) is less than $|I_G|$. Loops in (52) and (55) does not depend on n because $|D_x[i]| \leq |I_G|^2$ from $D_x[i] \subseteq \mathbf{Parent}[y](i \in I_G)$. Therefore $|\mathbf{Prune}| = |\mathbf{Is_prec}|$ holds.

As $\mathbf{Parent}[x]$ in (66) and $\mathbf{Parent}[y]$ in (68) are separable, the time complexity to get $\mathbf{PrecTab}[i, l_x, l_y]$ at (71) and (72) from known values of $\mathbf{PrecTab}$ is $O(1)$. Therefore the time complexity to get all values of $\mathbf{PrecTab}[i, j, k]$ is $O(n^2)$. \square

From Proposition 12 and Proposition 13, we have the following theorem.

Theorem 5. *The time complexity of Pruned Parsing is $O(n^2)$.*

7 An Example for a Grammar with Strong Ambiguity

In this section, as an example of this Pruned Graph-Structured Stack Parsing, we show that a grammar $S_m(m \geq 3)$ is parsed in $O(n^2)$. The time complexity for parsing S_m is $O(n^3)$ in [3] and $O(n^{m+1})$ in [11].

Example 19. Let $S3$ be a context free grammar with following productions.

$$(r^0) S' : S \dashv \quad (r^1) S : S S S \quad (r^2) S : S a \quad (r^3) S : a$$

$S3$ is a grammar used in [3] as an example grammar that forces the worst-case behavior for algorithms in [11] and [3]. We call a grammar that has same productions with $S3$ except that r^1 has m nonterminals (S) in the right hand side.

Example 20. For $S3$ in example 19, L^α for inputs 'aa' and 'aaa' is shown in Fig. 17. In Fig. 17, \uparrow reduce and \downarrow reduce is omitted.

Stack tops of parse stacks in L^{a^n} and $L^{a^{n+1}}$ for $S3$ are shown in Fig. 18. In Fig. 18, a^n in an input part is denoted as n .

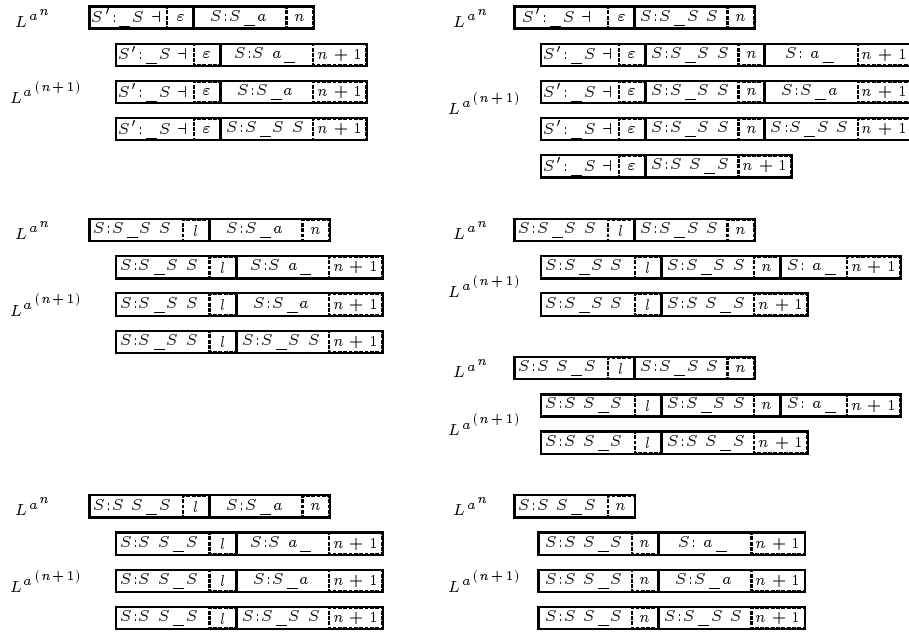


Figure 18: Generation stack top elements of $L^{a^{(n)}}$ from $L^{a^{(n+1)}}$ in S_3

From Fig. 18, we have

Proposition 14.

$$\begin{aligned} \text{Parent}(\langle S : S_a, n+1 \rangle) &= \{e_0, \langle S : S_S S, k \rangle, \langle S : S S_S, l \rangle\} (1 \leq k \leq n, 2 \leq l \leq n), \\ \text{Parent}(\langle S : S_S S, n+1 \rangle) &= \{e_0, \langle S : S_S S, k \rangle, \langle S : S S_S, l \rangle\} (1 \leq k \leq n, 2 \leq l \leq n), \\ \text{Parent}(\langle S : S S_S, n+1 \rangle) &= \{e_0, \langle S : S_S S, k \rangle, \langle S : S S_S, l \rangle\} (1 \leq k \leq n-1, 2 \leq l \leq n-1). \end{aligned}$$

From the definition of \prec and above proposition, we have

Proposition 15. For k and l that satisfy $k < l$,

$$\begin{aligned} \langle S : S_a, k \rangle &\prec \langle S : S_a, l \rangle, \\ \langle S : S_S S, k \rangle &\prec \langle S : S_S S, l \rangle, \\ \langle S : S S_S, k \rangle &\prec \langle S : S S_S, l \rangle. \end{aligned}$$

From the definition of Prune and above proposition, we have

Proposition 16. For case of $n \geq 2$

$$\begin{aligned} \text{Prune}(\text{Parent}(\langle S : S_a, n+1 \rangle)) &= \{e_0, \langle S : S_S S, n \rangle, \langle S : S S_S, n \rangle\}, \\ \text{Prune}(\text{Parent}(\langle S : S_S S, n+1 \rangle)) &= \{e_0, \langle S : S_S S, n \rangle, \langle S : S S_S, n \rangle\}, \\ \text{Prune}(\text{Parent}(\langle S : S S_S, n+1 \rangle)) &= \{e_0, \langle S : S_S S, n-1 \rangle, \langle S : S S_S, n-1 \rangle\}. \end{aligned}$$

From above proposition, S_3 is prunable and we have

Proposition 17. The Time complexity for parsing S_3 by Pruned Graph-Structured Stack Parsing is $O(n^2)$.

Similar in proposition 17, S_m ($m > 3$) is also prunable and the time complexity for parsing S_m by Pruned Graph-Structured Stack Parsing is $O(n^2)$.

Example 21. Relations between $\langle S : S_S S, k \rangle$ ($1 \leq k \leq 4$) and their parent sets in S_3 in Structured Graph Stack Parsing and Pruned Structured Graph Stack Parsing are shown in Fig. 19 in which the upper part indicates the case for Structured Graph Stack Parsing and the

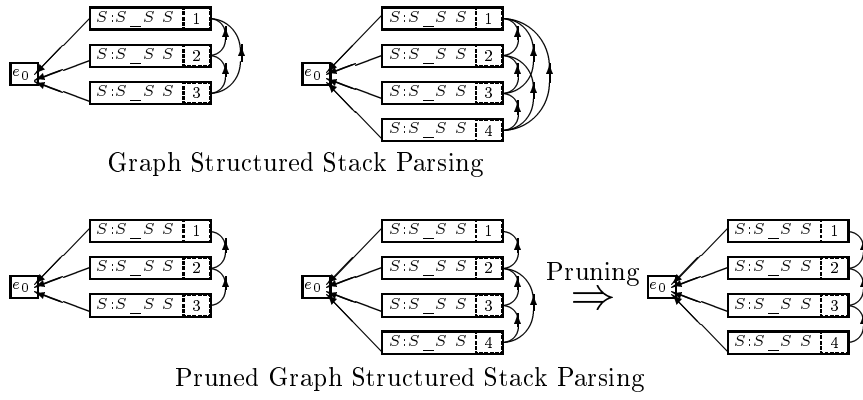


Figure 19: Generation of $\langle S : S_S S, 4 \rangle$ in $S3$ (pruned and unpruned version)

lower part indicates the case for Pruned Structured Graph Stack Parsing. In this figure, a^n in an input part is denoted as n and $S : S_S S$ is denoted as i .

The left part shows relations between $\langle i, k \rangle (1 \leq k \leq 3)$ and their parent sets after reading "aaa". In unpruned case (upper left part), $\langle i, j \rangle \in \text{Parent}(\langle i, k \rangle) (1 \leq j < k \leq 3)$, $e_0 \in \text{Parent}(\langle i, k \rangle)$ hold from proposition 15. In pruned case (lower left part), $\langle i, k-1 \rangle \in \text{Parent}(\langle i, k \rangle) (1 < k \leq 4)$, $e_0 \in \text{Parent}(\langle i, k \rangle)$ hold from proposition 16.

The middle part shows relations between $\langle i, k \rangle (1 \leq k \leq 4)$ and their parent sets after reading "aaaa". In unpruned case (upper middle part), $\langle i, j \rangle \in \text{Parent}(\langle i, 4 \rangle)$ and $(1 \leq j < 4)$, $e_0 \in \text{Parent}(\langle i, 4 \rangle)$ hold. In pruned case (lower middle part), $\langle i, j \rangle \in \text{Parent}(\langle i, 4 \rangle) (j = 2, j = 3)$ and $e_0 \in \text{Parent}(\langle i, 4 \rangle)$ hold.

The right part shows relations between $\langle i, k \rangle (1 \leq k \leq 4)$ and their parent sets after pruning. As $\langle i, 2 \rangle \prec \langle i, 3 \rangle$ holds, $\langle i, 2 \rangle$ is deleted from $\text{Parent}(\langle i, 4 \rangle)$ by pruning. Therefore $\langle i, 3 \rangle \in \text{Parent}(\langle i, 4 \rangle)$, $e_0 \in \text{Parent}$ holds. In this case, $I^{\langle i, k \rangle} (1 \leq k \leq 4)$ is conserved after pruning.

8 Applicable and Inapplicable grammars

In this section, we show examples of grammars that can be parsed or can not be parsed in $O(n^2)$ by Graph Structured Stack Parsing and Pruned Graph Structured Stack Parsing.

From proposition 13, for a grammar G , if $\text{Parent}(x)$ is separable for any $x \in E$, then G can be parsed in $O(n^2)$ by Graph Structured Stack Parsing. From proposition 13 and 14, if G is prunable, then G can be parsed in $O(n^2)$ by Pruned Graph Structured Stack Parsing. Furthermore, similar for proposition 13 and 14, we can prove that for a grammar G , if $|\text{Parent}(x)|$ is constant (not depend on n), then G can be parsed in $O(n^2)$.

The condition for not parsing in $O(n^2)$ by Graph Structured Stack Parsing and Pruned Graph Structured Stack Parsing is that there exists $x \in E$ whose parent set has elements y_1, \dots, y_l that satisfy following conditions.

1. $l = O(n)$,
2. $y_j \prec y_k$ does not hold for any $j, k (1 \leq j, k \leq l)$.

Example 22. Let G_{RL} be context free grammar with following productions.

$$(r^0) S' : S \dashv \quad (r^1) S : a S \quad (r^2) S : X \quad (r^3) X : X a \quad (r^4) X : a$$

Hereafter we denote $S : a_S$ as i and $X : X_a$ as i' .

$|\text{Parent}(\langle i', n \rangle)| = O(n)$ holds because $\text{Parent}(\langle i', n \rangle) = \{\langle i, k \rangle\} (2 \leq k \leq n)$ and \prec relation does not hold between $\langle i, k \rangle (1 \leq k \leq n)$ because $\text{Parent}(\langle i, 1 \rangle) = \{e_0\}$, $\text{Parent}(\langle i, n \rangle) = \{\langle i, n-1 \rangle\} (2 \leq n)$. G_{RL} can NOT be parsed in $O(n^2)$ by neither Graph Structured Stack Parsing nor Pruned Graph Structured Stack Parsing.

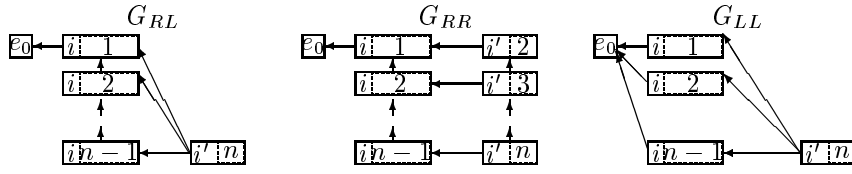


Figure 20: Relations among Parent sets

Example 23. Let G_{RR} be context free grammar with following productions.

$$(r^0) S' : S \dashv \quad (r^1) S : a S \quad (r^2) S : X \quad (r^3) X : a X \quad (r^4) X : a$$

Hereafter we denote $S : a _ S$ as i and $X : a _ X$ as i' .

Parent($\langle i, n \rangle$) and Parent($\langle i', n \rangle$) is separable because Parent($\langle i, 1 \rangle$) = $\{e_0\}$, Parent($\langle i, n \rangle$) = $\{\langle i, n-1 \rangle\}$, Parent($\langle i', n \rangle$) = $\{\langle i, n-1 \rangle, \langle i', n-1 \rangle\}$ ($2 \leq n$) hold. G_{RR} can be parsed in $O(n^2)$ by Graph Structured Stack Parsing.

Example 24. Let G_{LL} be context free grammar with following productions.

$$(r^0) S' : S \dashv \quad (r^1) S : S X \quad (r^2) S : X \quad (r^3) X : X a \quad (r^4) X : a$$

Hereafter we denote $S : a _ S$ as i and $X : X _ a$ as i' .

G_{LL} is prunable because $\langle i, j \rangle \prec \langle i, k \rangle$ ($1 \leq j < k < n$) hold from Parent($\langle i, n \rangle$) = $\{e_0\}$, Parent($\langle i', n \rangle$) = $\{\langle i, k \rangle\}$ ($1 \leq k < n$). G_{LL} can NOT be parsed in $O(n^2)$ by Graph Structured Stack Parsing but can be parsed by Pruned Graph Structured Stack Parsing .

Parent relation between $\langle i, j \rangle, \langle i', k \rangle$ ($1 \leq j, k \leq n$) in G_{RL}, G_{RR}, G_{LL} are shown in Fig. 20.

9 Application for Natural Language Processing

In this section, application of this algorithm for natural language processing is stated. In natural language processing, large-scale grammars are prerequisite for parsing a great variety of sentences. Such grammars are derived from a large-scale syntactically annotated corpus because it is difficult to build such grammars by hand [6]. As these grammars are derived automatically, they are highly ambiguous and some input sentence generates many parse trees.

Current approaches to reduce ambiguity are

- (A) Use of probabilistic context free grammars [2].
- (B) Modifying the language grammar [6].

Probabilistic context free grammar is a context free grammar in which probabilities are assigned to each productions or shift/reduce actions in an LR table. Therefore probabilistic context free grammar can be regarded to assign total order(probabilities) between productions or shift/reduce actions to determine the production or action to be selected. On the other hand, as \prec relation of definition 15 can be regarded as partial order over e-items (to be exact, \prec is preorder because it is reflexive and transitive but not antisymmetric), this algorithm assigns partial order between e-items to be selected. As probabilities are assigned to each productions, it is difficult to assign consistent probabilities over productions. In this algorithm, \prec relation is consistent over e-items because it is defined recursively over e-items.

According to [6], the major ambiguities of a grammar for Japanese are the following.

- (B-1) Compound noun structure
- (B-2) Adnominal phrase attachment
- (B-3) Conjunctive structure

The production for the compound noun structure is "CompoundNoun : CompoundNoun CompoundNoun" and it is modified to "CompoundNoun : Noun CompoundNoun" to reduce ambiguity.

The productions for the adnominal phrase attachment are "NounPhrase : AdnominalPhrase NounPhrase" and "AdnominalPhrase : NounPhrase Particle" . The latter production is modified to "AdnominalPhrase : NounPhrase' Particle" where NounPhrase' is a nonterminal symbol that derives same symbols as NounPhrase except AdnominalPhrase.

The productions for the conjugate structure is "NounPhrase : NounPhrase Particle NounPhrase" and it is modified to "AdnominalPhrase : NounPhrase' Particle".

The situations of above modifications are shown in Fig.21, Fig.22 and Fig.23. These figures are modified from figures in [6].

As these modifications are done manually, it is difficult to make consistent modifications over large-scale grammar. Productions mentioned above are essentially similar to the production(r^1) of S_2 , therefore ambiguities are reduced by this algorithm without modifying productions.

10 Conclusion

In this report, I introduce a parsing algorithm that improves the time complexity by pruning the edges of the graph structured stack. Many approaches such as use of tables and memo functions have been proposed to improve the time complexity of parsing. These approaches improve the complexity by suppressing redundant process that are already executed. This approach(pruning) improves the complexity by suppressing redundant process that will be executed. This approach uses stacks and data in stack can be regarded as data processed in future because all data in stack will be processed and popped in future. Therefore it is possible to find redundant process that will be executed in future in this algorithm.

References

- [1] Earley, J. : An efficient context free parsing algorithm, *Comm. ACM*, **13** (1970) 94–102.
- [2] Inui, K. , Sornlertlamvanich, V. , Tanaka, H. and Tokunaga, T. : Probabilistic GLR Parsing: A New Formalization and Its Impact on Parsing Performance, *Journal of Natural Language Processing*, **5** (1998) 33–52.
- [3] Kipps, J. R. : GLR parsing in time $O(n^3)$, in [10], Chap. 4, 43–59.
- [4] Leemakers, R. : A recursive ascent Earley parser, *Info. Process. Lett.* , **41** (1992) 87–91.
- [5] Nederhof, M. J. : Generalized left-corner parsing, *Proc. 6th Conf. on European Chapter of the ACL*, 1993, 305–314.
- [6] Noro, T. , Hashimoto, T. , Tokunaga, T. and Tanaka, H. : Building a Large-Scale Japanese Grammar. *Journal of Natural Language Processing*, **12** (2005) 3–31. (in Japanese)
- [7] Schabes, Y. : Polynomial time and space shift-reduce parsing of arbitrary context-free grammars, *Proc. 29th ACL*, 1991, 106–113.
- [8] Shann, P. : Experiments with GLR and chart parsing, in [10], Chap. 2, 17–34.
- [9] Sippu, S. and Soisalon-Soininen, E. : *Parsing Theory, Vols. I and II*, Springer-Verlag, Berlin, 1988, 1990.
- [10] Tomita, M. , ed. : *Generalized LR Parsing*, Kluwer Academic Publishers, 1991.
- [11] Tomita, M. , Ng, S. K. : The generalized LR parsing algorithm, in [10], Chap. 1, pp. 1–16 (1991).

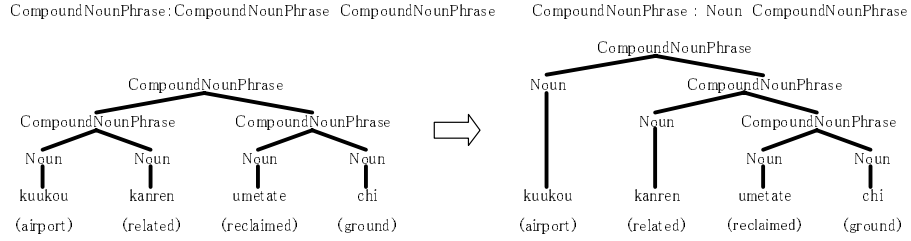


Figure 21: Production modification for Compound Noun Structure

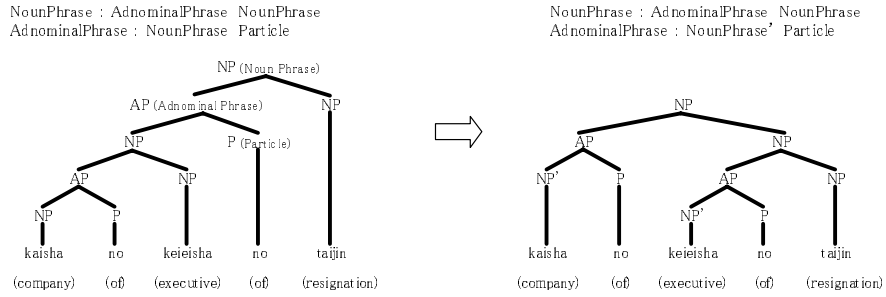


Figure 22: Production modification for Adnominal Phrase Attachment

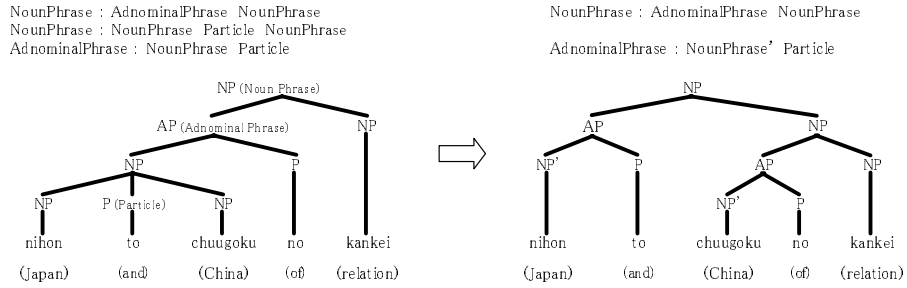


Figure 23: Production modification for Conjunctive Structure

Appendix Proofs of Theorem 1 and Theorem 2

I hereafter abbreviate $\text{Reduce}^{(n)}(L^\alpha, a)$ by $R^{(n)}(L^\alpha, a)$, and $\text{Reduce}_{(n)}(E^\alpha, a)$ by $R_{(n)}(E^\alpha, a)$.

Theorem 1. For $\alpha \in V_T^*$, if $\lambda \in L^\alpha$ then $\lambda \in E^\alpha$.

Proof. From Algorithm 1, $\lambda \in L^\alpha$ is expressed as one of the following five forms:

1. $\alpha = \varepsilon, \lambda = e_0$.
2. $\alpha = \alpha' \cdot a, \lambda = \mu \cdot y$ for $x \in J_G$, where $\lambda' = \mu \cdot x \in L^{\alpha'}, x \in \text{S0item}(a), y \in \text{Shift0}(x, a)$.
3. $\alpha = \alpha' \cdot a, \lambda = \mu \cdot x \cdot y$ for $x \in J_G$, where $\lambda' = \mu \cdot x \in L^{\alpha'}, x \in \text{S1item}(a), y \in \text{Shift1}(x, a)$.
4. $\alpha = \alpha' \cdot a, \lambda = \mu \cdot y$ for $x, x' \in J_G$ where $\lambda' = \mu \cdot x' \cdot x \in R^{(n)}(L^{\alpha'}, a) \subseteq E^\alpha, n \geq 0, x \in \text{Redex}, x' \in \text{R0item}(x), y \in \text{Reduce0}(x', x)$.
5. $\alpha = \alpha' \cdot a, \lambda = \mu \cdot x' \cdot y$ for $x \in J_G$ where $\lambda' = \mu \cdot x' \cdot x \in R^{(n)}(L^{\alpha'}, a) \subseteq E^\alpha, n \geq 0, x \in \text{Redex}, x' \in \text{R1item}(x), y \in \text{Reduce1}(x', x)$.

For Case 1, this theorem holds, because $L^\varepsilon = \{e_0\} = E^\varepsilon$.

For Cases 2, 3, 4 and 5 following relations (A) and (B) hold from the induction hypothesis.

$$\text{Parent}(\lambda_1) = \emptyset, \quad (\text{A})$$

$$\lambda_j \in \text{Parent}(\lambda_{(j+1)})(1 \leq j < \#\lambda - 1). \quad (\text{B})$$

Moreover, (C) and (D) hold in cases 2, 3, 4 and 5.

$$\lambda_{\#} \in E_\alpha, \quad (\text{C})$$

$$\lambda_{(\#\lambda-1)} \in \text{Parent}(\lambda_{\#}). \quad (\text{D})$$

From (B) and (D)

$$\lambda_j \in \text{Parent}(\lambda_{(j+1)})(1 \leq j < \#\lambda) \quad (\text{E})$$

holds. By (A), (C) and (E), λ is in E^α , and this theorem is satisfied.

Proofs for (C) and (D) in cases 2, 3, 4 and 5 are as follows.

(2-C) As $\lambda_{\#} = y \in \text{Shift0}(x, a)$ and $x = \lambda'_{\#} \in E_{\alpha'}$, $\lambda_{\#} = y \in E_{\alpha' \cdot a} = E_\alpha$ holds.

(2-D) As $\mu_{\#} = \lambda'_{(\#\lambda'-1)} \in \text{Parent}(\lambda'_{\#}) = \text{Parent}(x)$ (by induction hypothesis)

and $\text{Parent}(y) = \text{Parent}(x)$ (from lemma 1),

$\lambda_{(\#\lambda-1)} = \mu_{\#} \in \text{Parent}(x) = \text{Parent}(y) = \text{Parent}(\lambda_{\#})$ holds.

(3-C) As $\lambda_{\#} = y \in \text{Shift1}(x, a)$ and $x = \lambda'_{\#} \in E_{\alpha'}$, $\lambda_{\#} = y \in E_{\alpha' \cdot a} = E_\alpha$ holds.

(3-D) As $x = \lambda'_{\#} \in E_{\alpha'}$ (by induction hypothesis) and

$\text{Parent}(y) = \{z \mid z \in E_{\alpha'}, z \in \text{S1item}(a), y \in \text{Shift1}(z, a)\}$ (from case 3 in Algorithm 2),

$\lambda_{(\#\lambda-1)} = x \in \text{Parent}(y) = \text{Parent}(\lambda_{\#})$ holds.

(4-C) As $\lambda_{\#} = y \in \text{Reduce0}(x', x)$ and $x \in E_\alpha, x' \in \text{Parent}(x), \lambda_{\#} = y \in E_\alpha$ holds.

(4-D) As $x' = \lambda'_{(\#\lambda'-1)} \in \text{Parent}(\lambda'_{\#}) = \text{Parent}(x)$ (by induction hypothesis)

and $\text{Parent}(x) \subseteq \text{Parent}(y)$ (from case 4 of Algorithm 2),

$\lambda_{(\#\lambda-1)} = x' \in \text{Parent}(x) \subseteq \text{Parent}(y) = \text{Parent}(\lambda_{\#})$ holds.

(5-C) As $\lambda_{\#\lambda} = y \in \text{Reduce1}(x', x)$ and $x \in E_\alpha, x' \in \text{Parent}(x), \lambda_{\#} = y \in E_\alpha$ holds.

(5-D) As $x' = \lambda'_{(\#\lambda'-1)} \in \text{Parent}(\lambda'_{\#\lambda'}) = \text{Parent}(x)$ (by induction hypothesis) and

$x' \in \text{Parent}(y)$ (from case 5 of Algorithm 2 and $x = \lambda'_{\#} \in E_\alpha \subseteq R_{(n)}(E^\alpha, a)$),

$\lambda_{(\#\lambda-1)} = x' \in \text{Parent}(y) = \text{Parent}(\lambda_{\#})$ holds. □

Theorem 2. For $\alpha \in V_T^*$, if $\sigma \in E^\alpha$ then $\sigma \in L^\alpha$.

Proof. From Algorithm 2, $\sigma \in E^\alpha$ is expressed as one of the following five forms:

1. $\alpha = \varepsilon, \sigma = e_0$.

2. $\alpha = \alpha' \cdot a$, $\sigma = \tau \cdot y$ for $x \in J_G$ where $\tau \cdot x \in E^{\alpha'}$, $x \in \text{S0eitem}(a)$, $y \in \text{Shift0}(x, a)$.
3. $\alpha = \alpha' \cdot a$, $\sigma = \tau \cdot x \cdot y$ where $\tau \cdot x \in E^{\alpha'}$, $x \in \text{S1eitem}(a)$, $y \in \text{Shift1}(x, a)$.
4. $\alpha = \alpha' \cdot a$, $\sigma = \tau \cdot y$ for x' , $x \in J_G$ where $\tau \cdot x' \cdot x \in E^{\alpha'}$,
 $n \geq 0$, $x \in \text{Redex}$, $x' \in \mathbf{R}_{(n)}(E^{\alpha'}, a)$, $x \in \text{R0eitem}(x)$, $y \in \text{Reduce0}(x', x)$.
5. $\alpha = \alpha' \cdot a$, $\sigma = \tau \cdot x' \cdot y$ for $x \in J_G$ where $\tau \cdot x' \cdot x \in E^{\alpha'}$,
 $n \geq 0$, $x \in \text{Redex}$, $x' \in \mathbf{R}_{(n)}(E^{\alpha'}, a)$, $x \in \text{R1eitem}(x)$, $y \in \text{Reduce1}(x', x)$.

For Case 1, this theorem holds, since $E^\varepsilon = \{e_0\} = L^\varepsilon$.

For case 2, as $\tau \cdot x \in L^{\alpha'}$ (by $\tau \cdot x \in E^{\alpha'}$ and induction hypothesis) and $x \in \text{S0eitem}(a)$, $y \in \text{Shift0}(x, a)$ hold, $\sigma = \tau \cdot y \in L^{\alpha' \cdot a} = L^\alpha$ is satisfied from case 2 of Algorithm2.

For case 3, $\sigma \in L^{\alpha' \cdot a} = L^\alpha$ is satisfied similar to case 2.

For case 4, as $\tau \cdot x' \cdot x \in L^{\alpha'}$ (by $\tau \cdot x' \cdot x \in E^{\alpha'}$ and induction hypothesis) and $x \in \text{Redex}$, $x' \in \mathbf{R}_{(n)}(E^{\alpha'}, a)$, $x \in \text{R0eitem}(x)$, $y \in \text{Reduce0}(x', x)$ hold, $\sigma = \tau \cdot y \in L^\alpha$ is satisfied from case 4 of Algorithm2.

For case 5, $\sigma \in L^\alpha$ is satisfied similar to case 4. □