



Available at  
[www.ComputerScienceWeb.com](http://www.ComputerScienceWeb.com)  
POWERED BY SCIENCE @ DIRECT®

Theoretical Computer Science 298 (2003) 179–206

Theoretical  
Computer Science

[www.elsevier.com/locate/tcs](http://www.elsevier.com/locate/tcs)

# Polynomial-time identification of very simple grammars from positive data<sup>☆</sup>

Takashi Yokomori

*Department of Mathematics, School of Education, Waseda University, 1-6-1 Nishiwaseda,  
Shinjuku-ku, Tokyo 169-8050, Japan*

---

## Abstract

This paper concerns a subclass of simple deterministic grammars, called *very simple grammars*, and studies the problem of identifying the subclass in the limit from positive data. The class of very simple languages forms a proper subclass of simple deterministic languages and is incomparable to the class of regular languages. This class of languages is also known as the class of left Szilard languages of context-free grammars.

After providing some properties of very simple languages, we show that the class of very simple grammars is polynomial-time identifiable in the limit from positive data in the following sense. That is, we show that there effectively exists an algorithm that, given a target very simple grammar  $G_*$  over alphabet  $\Sigma$ , identifies a very simple grammar  $G$  equivalent to  $G_*$  in the limit from positive data, satisfying the property that the time for updating a conjecture is bounded by  $O(m)$ , and the total number of prediction errors made by the algorithm is bounded by  $O(n)$ , where  $n$  is the size of  $G_*$ ,  $m = \text{Max}\{N^{|\Sigma|+1}, |\Sigma|^3\}$  and  $N$  is the total length of all positive data provided.

© 2002 Elsevier Science B.V. All rights reserved.

*Keywords:* Very simple grammars; Identification in the limit from positive data; Polynomial-time learning

---

## 1. Introduction

Since the class of regular languages has been shown to be efficiently identifiable using deterministic finite-state automata (DFAs) from what is called “minimally adequate teacher” [3], a computational approach to learning theory has been again receiving much attention, and a numerous number of intensive work on grammatical inference

---

<sup>☆</sup> A preliminary version appeared in the Proceedings of COLT'91, Santa Cruz, CA, 1991.

*E-mail address:* [yokomori@mn.waseda.ac.jp](mailto:yokomori@mn.waseda.ac.jp) (T. Yokomori).

has been reported. From the practical point of view, there are, we believe, two major requirements for an inductive inference algorithm: the identification algorithm must have a good *time efficiency* and hopefully run with only *positive data (examples)*.

Angluin [1] has given several conditions for a class of languages to be identifiable in the limit from positive data, and presented some examples of identifiable classes. She has also proposed subclasses of regular languages called *k-reversible languages* (where  $k \geq 0$ ) and shown these classes are identifiable in the limit from positive data, requiring a polynomial time for updating conjectures [2]. Wright [22] proposes a sufficient condition for a language class to be identifiable in the limit from positive data, while Sato and Umayahara [18] introduce another sufficient condition for the identifiability in the limit from positive data and discuss the relationship to Wright's condition.

On the other hand, as for efficient identifiability results on subclasses of regular languages, Tanida and Yokomori [20,23] have introduced a new subclass of regular languages called strictly deterministic regular languages and shown its identifiability from positive data. Further, in [6,24], another subclass of regular languages called strictly *k-testable languages* is studied and shown to be identifiable in the limit from positive data, together with its application results to pattern recognition and DNA sequence analysis, respectively.

Motivated by a question posed by Angluin, however, one natural question has been recognized to be significant: In what sense we should analyse the time complexity of an “in-the-limit” algorithm? One may define the notion of *polynomial-time identification in the limit* in various ways. And, it was not until quite recently that the polynomial-time identifiability in the limit was reasonably defined by Pitt. By making a slight modification of his definition in [17], we propose the following definition for polynomial-time identifiability in the limit from positive data.

Informally, we say a class of grammars  $\mathcal{R}$  is *polynomial-time identifiable in the limit from positive data* if and only if there is an algorithm  $A$  which, given  $r$  in  $\mathcal{R}$ , identifies  $r'$  in  $\mathcal{R}$  equivalent to  $r$  in the limit from positive data, with the property that there exist polynomials  $p$  and  $q$  such that for any  $n$ , for any  $r$  of size  $n$ , the number of times  $A$  makes a wrong conjecture is at most  $p(n)$ , and the time for updating a conjecture is at most  $q(n, N)$ , where  $N$  is the sum of lengths of data provided.

This paper deals with a class of grammars called *very simple grammars* and discusses the identification problem of the class of very simple grammars. To the author's knowledge, the notion of a very simple grammar was originally introduced in [5] in the study of some types of Thue systems and the equivalence problem for simple deterministic grammars by Korenjak and Hopcroft [10]. Informally, a context-free grammar  $G$  is *very simple* if the righthand side of each rule starts with distinct terminal symbol, possibly followed by a nonterminal string. For example, the context-free grammar  $G_1 = (\{E, R\}, \{+, -, *, (, ), a\}, P_1, E)$ , where  $P_1 = \{E \rightarrow +EE, E \rightarrow -EE, E \rightarrow *EE, E \rightarrow (ER, E \rightarrow a, R \rightarrow )\}$  generates a subset of the arithmetic expressions in prefix notation. Further, consider the context-free grammar  $G_2 = (\{S, T, G\}, \{a, c, d, \bar{a}, \bar{c}\}, P_2, S)$ , where  $P_2 = \{S \rightarrow aST, S \rightarrow cSG, S \rightarrow d, T \rightarrow \bar{a}, G \rightarrow \bar{c}\}$  generating  $\{xd\bar{x}^R \mid x \in \{a, c\}^*\}$ . Under the interpretation that  $\bar{a} = t$ ,  $\bar{c} = g$ , this language represents biological “palindromes” to form secondary structures in nucleic acids. These grammars  $G_1, G_2$  are typically very simple, and, thus, very simple grammars can capture the non-regularity feature

necessary for describing a certain class of popular languages. (Also, see Example 1 below.) The class of very simple languages forms a proper subclass of the simple deterministic languages and is incomparable to the class of regular languages. On the other hand, Mäkinen has discussed the left Szilard languages of context-free grammars in great detail, investigating the derivational properties of context-free grammars [12], and it turns out that the class of left Szilard languages of context-free grammars coincides with the class of very simple languages.

After providing some properties of very simple languages, we show that the class of very simple grammars is polynomial-time identifiable in the limit in the sense above. In fact, the identification of the class is achieved using only positive data, and this result is in marked contrast with the fact that the class of regular languages is not at all identifiable in the limit from positive data only [7].

Thus, the main result in this paper provides an interesting instance of a language class containing nonregular languages which is polynomial-time identifiable in the limit. As a corollary, it immediately follows that the class of very simple grammars is polynomial-time identifiable via only equivalence queries, provided that only positive counterexamples are supplied in the identification process.

## 2. Definitions

### 2.1. Basic definitions and notations

We assume the reader to be familiar with the rudiments of automata and formal language theory. (For notions and notations not stated here, see, e.g., [9].)

Let  $\Sigma$  be a finite alphabet and  $\Sigma^*$  be the set of all finite length strings over  $\Sigma$ . Further, let  $\Sigma^+ = \Sigma^* - \{\lambda\}$ , where  $\lambda$  is the null string. By  $len(u)$  we denote the length of a string  $u$ . A *language over  $\Sigma$*  is a subset of  $\Sigma^*$ . For a string  $w$  in  $\Sigma^*$ ,  $\text{alph}(w)$  denotes the set of terminal symbols appearing in  $w$ . For a language  $L$ , let  $\text{alph}(L) = \bigcup_{w \in L} \text{alph}(w)$ .

Let  $h$  be a homomorphism from  $\Sigma^*$  to  $\Delta^*$ . Then,  $h$  is called a *coding* if and only if for all  $a$  in  $\Sigma$ ,  $h(a)$  is in  $\Delta$ . Further, a coding  $h$  is called a *renaming* if and only if it is a bijection.

Let  $G = (V_N, \Sigma, P, S)$  be a context-free grammar (CFG). The language generated by  $G$  is defined by  $\{w \in \Sigma^* \mid S \Rightarrow^* w \text{ in } G\}$ .

By  $\text{Label}(P)$  we denote the set of labels of all rules in  $P$ . Then, define  $D_L(G)$  as  $\{\pi_1 \cdots \pi_n \in \text{Label}(P)^* \mid S \Rightarrow_L^{\pi_1} a_1 \alpha_1 \Rightarrow_L \cdots \Rightarrow_L^{\pi_n} a_1 \cdots a_n \in \Sigma^*\}$ , where each  $\pi_i$  is a label of a rule used in the  $i$ th step, and  $\Rightarrow_L$  indicates the left-most derivation. The language  $D_L(G)$  is called *left Szilard language* of  $G$ .

Let  $G_i = (V_{N_i}, \Sigma, P_i, S_i)$  ( $i = 1, 2$ ) be CFGs. Then,  $G_1$  and  $G_2$  are *isomorphic* if and only if there is a renaming  $h$  from  $(V_{N_1} \cup \Sigma)^*$  to  $(V_{N_2} \cup \Sigma)^*$  such that  $h(S_1) = S_2$ , the restriction of  $h$  to  $\Sigma$  is the identity, and  $P_2 = \{h(X) \rightarrow h(\alpha) \mid X \rightarrow \alpha \in P_1\}$ . A CFG  $G = (V_N, \Sigma, P, S)$  is *linear* if and only if any rule of  $G$  is of the form either  $A \rightarrow uBv$  or  $A \rightarrow v$ , where  $A, B \in V_N$ ,  $u, v \in \Sigma^*$ .

For a language  $L$  over  $\Sigma$  and  $x$  in  $\Sigma^*$ , let  $x \setminus L = \{y \mid xy \in L\}$ . The cardinality of a set  $S$  is denoted by  $|S|$ .

A language  $L$  is said to have *the prefix-free property* if and only if for any  $x \in \Sigma^+$ ,  $y \in \Sigma^*$ , both  $x$  and  $xy$  are in  $L$  implies that  $y = \lambda$ .

A (directed finite) labeled graph is an ordered triple  $(V, E, \Sigma)$ , where  $\Sigma$  is a finite alphabet of labels,  $V$  is a finite set of nodes and  $E = \{(p, a, q) \mid p, q \in V, a \in \Sigma\}$  is a finite set of edges.

A sequence of edges:  $\sigma = (p_1, u_1, q_1), (p_2, u_2, q_2), \dots, (p_n, u_n, q_n)$ , ( $n \geq 1$ ), is a *path* from  $p_1$  to  $q_n$  if and only if  $q_i = p_{i+1}$  for all  $i$  ( $1 \leq i < n$ ). In this case, a string  $x = u_1 u_2 \cdots u_n$  is called a *path string* of  $\sigma$ .

## 2.2. Polynomial-time identification in the limit from positive data

In this paper we adopt a slight modification of the original definition in [17] for the notion of polynomial-time identification in the limit from positive data. (In fact, the definition turns out to be a restricted version of the Pitt's definition where positive data is only provided.)

For any class of languages to be identified, let  $\mathcal{R}$  be a *class of representations* for the class of languages. Given an  $r$  in  $\mathcal{R}$ ,  $L(r)$  denotes the language represented by  $r$ . A *positive presentation* of  $L(r)$  is any infinite sequence of data such that every  $w \in L(r)$  occurs at least once in the sequence and no other strings not in  $L(r)$  appear in the sequence. Each element of  $L(r)$  is called a *positive example* (or simply, *example*) of  $L(r)$ .

Let  $r$  be a representation in  $\mathcal{R}$ . An algorithm  $\mathcal{A}$  is said to *identify  $r$  in the limit from positive data* if and only if  $\mathcal{A}$  takes any positive presentation of  $L(r)$  as input, and outputs an infinite sequence of  $r_i$ s in  $\mathcal{R}$  such that there exist  $r'$  in  $\mathcal{R}$  and  $j > 0$  so that for all  $i \geq j$ , the  $i$ th conjecture (representation)  $r_i$  is identical to  $r'$  and  $L(r') = L(r)$ . A class  $\mathcal{R}$  is *identifiable in the limit from positive data* if and only if there exists an algorithm  $\mathcal{A}$  that, given any  $r$  in  $\mathcal{R}$ , identifies  $r$  in the limit from positive data.

Let  $\mathcal{A}$  be an algorithm for identifying  $\mathcal{R}$  in the limit from positive data. Suppose that after examining  $i$  examples, the algorithm  $\mathcal{A}$  conjectures some  $r_i$ . We say that  $\mathcal{A}$  makes an *implicit error of prediction* at step  $i$  if  $r_i$  is not consistent with the  $(i + 1)$ st example, i.e., if  $L(r_i)$  fails to contain the  $(i + 1)$ st example.

A class  $\mathcal{R}$  is *polynomial-time identifiable in the limit from positive data* if and only if there exists an algorithm  $\mathcal{A}$  for identifying  $\mathcal{R}$  in the limit from positive data with the property that there exist polynomials  $p$  and  $q$  such that for any  $n$ , for any  $r$  of size  $n$ , and for any positive presentation of  $L(r)$ , the number of implicit errors of prediction made by  $\mathcal{A}$  is at most  $p(n)$ , and the time used by  $\mathcal{A}$  between receiving the  $i$ th example  $w_i$  and outputting the  $i$ th conjecture  $r_i$  is at most  $q(n, m_1 + \cdots + m_i)$ , where  $m_j = \text{len}(w_j)$  and the size of  $r$  is the length of a description for  $r$ .

## 3. Very simple grammars and languages

Let  $G = (V_N, \Sigma, P, S)$  be a context-free grammar (CFG) in Greibach normal form (GNF), i.e., each rule of  $P$  is of the form:  $A \rightarrow ax$ , where  $A \in V_N$ ,  $a \in \Sigma$ ,  $x \in V_N^*$ . If for

all  $(a, A) \in \Sigma \times V_N$ , there is at most one rule of the form  $A \rightarrow a\alpha$  (for some  $\alpha \in V_N^*$ ) in  $P$ , then  $G$  is said to be *simple deterministic*.

Let  $G$  be a CFG in GNF. For each terminal symbol  $a \in \Sigma$ , a rule whose right-hand side is of the form  $a\alpha$  (where  $\alpha \in V_N^*$ ) is called an *a-handle rule*. Then,  $G$  is said to be *very simple* if and only if for each  $a$  in  $\Sigma$ , there exists exactly one *a-handle rule* in  $P$ . A language  $L$  is *very simple* if and only if there exists a very simple CFG  $G$  such that  $L = L(G)$  holds. (Note that since every very simple grammar is  $\lambda$ -free, so is every very simple language.)

Let  $S \Rightarrow_L^* x\alpha$  be the left-most derivation in  $G = (V_N, \Sigma, P, S)$ , where  $x \in \Sigma^*$ ,  $\alpha \in V_N^*$ . Then, define  $L(\alpha) = \{w \in \Sigma^* \mid \alpha \Rightarrow_L^* w\}$ . A CFG  $G$  is *reduced* if and only if for any  $\alpha$  such that  $S \Rightarrow_L^* x\alpha$  in  $G$ , it holds that  $L(\alpha) \neq \emptyset$ .

**Example 1** (Laird and Gamble [11]). Consider a CFG  $G = (\{S, C, L, T, R\}, \{p, s, t, x, \wedge, (, ), 0, \bullet\}, P, S)$ , where  $P$  consists of the following:

$$\begin{aligned} S &\rightarrow pLTCTCTR, S \rightarrow \wedge LSCSR, S \rightarrow t, \\ T &\rightarrow sLTR, \quad T \rightarrow x, \quad T \rightarrow 0, \\ C &\rightarrow \bullet, \quad L \rightarrow (, \quad R \rightarrow ). \end{aligned}$$

It is easily seen that the grammar  $G$  is very simple and that, for example, strings

$$“p(s(0)\bullet x\bullet 0)” \quad \text{and} \quad “\wedge(p(0\bullet 0\bullet 0)\bullet p(s(0)\bullet 0\bullet s(0)))”$$

are generated by  $G$ . If we take  $p$ ,  $s$  and  $\wedge$  as a predicate symbol “**plus**”, a function symbol “*successor*”, and a logical “*and*” in prefix notation, respectively, then these two strings are interpreted as logical formulas: **plus**( $s(0), x, 0$ ) and **plus**( $0, 0, 0$ )  $\wedge$  **plus**( $s(0), 0, s(0)$ ). Note that “ $t$ ”, “ $x$ ”, and “ $\bullet$ ” denote a logical constant **true**, a variable, and comma, respectively. Thus, the language  $L(G)$  defines a subset of formulas in first-order predicate logic which is not a regular language. Note that the very simple grammar  $G$  is reduced.

[**Convention**] In what follows, it is assumed that very simple grammars we consider in this paper are all *reduced*. Further, a notation  $\Rightarrow$  (or  $\Rightarrow_G$ ) always denotes the leftmost derivation (by  $G$ ).

### 3.1. Basic properties

The next result immediately follows from the definition.

**Lemma 1.** *Let  $L$  be a very simple language.*

- (i) [Prefix-free Property] *For each  $x$  in  $\Sigma^*$  and  $y$  in  $\Sigma^+$ , if  $x \in L$ , then  $xy \notin L$ .*
- (ii) *For each  $w$  in  $L$ , if  $\text{len}(w) \geq 2$ , then the initial symbol of  $w$  must differ from the last symbol of  $w$ .*

**Example 2.** The followings are not very simple languages:  $\{abba\}$ ,  $\{a^n \mid n \geq 1\}$ , or  $\{c^m a c^n \mid m, n \geq 0\}$ .

Thus, the class of very simple languages forms a proper subclass of simple deterministic languages by Korenjak and Hopcroft [10]. Further, it is easy to see that the class is incomparable to the class of regular languages.

**Lemma 2.** *For any very simple grammar  $G$ , there exists a renaming  $f$  such that  $L(G) = f(D_L(G))$ , where  $D_L(G)$  is the left Szilard language of  $G$ .*

**Proof.** For a given very simple grammar  $G = (V_N, \Sigma, P, S)$ , define a renaming  $f$  by: for each  $\pi_a \in \text{Label}(P)$ ,  $f(\pi_a) = a$ , where  $\pi_a$  is a label of the  $a$ -handle rule. Then, it is obvious that  $S \Rightarrow^{\pi_{a_1} \cdots \pi_{a_n}} a_1 \cdots a_n \in \Sigma^*$  if and only if  $f(\pi_{a_1} \cdots \pi_{a_n}) = a_1 \cdots a_n \in \Sigma^*$ . Thus,  $L(G) = f(D_L(G))$  holds.  $\square$

Since left Szilard languages of CFGs of finite index (i.e., of nonterminal bounded CFGs) are regular [15], it is seen from Example 1 that there exists a very simple language  $L$  which is of infinite index (i.e., not nonterminal bounded).

### 3.2. Closure properties

In this subsection, we will present the closure properties of very simple languages.

**Lemma 3.**  $L_0 = \{a^{3n}b \mid n \geq 0\}$  is not a very simple language.

**Proof.** Suppose that  $L_0 = L(G)$  for some very simple grammar  $G = (V_N, \{a, b\}, P, S)$ , where we may assume that  $V_N = \{S, X\}$ ,  $P = \{S \rightarrow a\alpha, X \rightarrow b\beta\}$  ( $\alpha, \beta \in V_N^*$ ). Since  $b$  is in  $L_0$ , there exists a rule:  $S \rightarrow b$ . Hence, the  $b$ -handle rule is  $S \rightarrow b$  and  $X = S$ ,  $\beta = \lambda$ . Further, since  $aaab$  is in  $L_0$ , there must exist derivations  $S \Rightarrow a\alpha$  and  $\alpha \Rightarrow^* aab$ . Let  $\alpha = Y\alpha'$  ( $Y \in V_N, \alpha' \in V_N^*$ ), then there exist a rule  $Y \rightarrow a\gamma$  and a derivation  $\gamma\alpha' \Rightarrow^* ab$ . From the rules  $S \rightarrow a\alpha$  and  $Y \rightarrow a\gamma$ ,  $S = Y$  and  $\alpha = \gamma$  follow, and hence the  $a$ -handle rule is  $S \rightarrow aS\alpha'$ . Since  $\alpha (= S\alpha') \Rightarrow^* aab\alpha'^3$  and  $\alpha \Rightarrow^* aab$ , it must hold that  $\alpha' = \lambda$ . Thus, we have that  $P = \{S \rightarrow aS, S \rightarrow b\}$  and hence,  $L(G) = \{a^n b \mid n \geq 0\}$  contradicting to that  $L_0 = L(G)$ .  $\square$

**Theorem 4.** *The class of very simple languages is closed under none of the following: union, concatenation, intersection, complement, Kleene closure (+, \*), ( $\lambda$ -free) homomorphism, inverse homomorphism, intersection with regular languages, or reversal.*

**Proof.** (Union:) Consider  $L_1 = \{ab\}$  and  $L_2 = \{abc\}$  which are clearly very simple. Since very simple (or simple deterministic) languages have the prefix-free property [8],  $L_1 \cup L_2 = \{ab, abc\}$  is not very simple.

(Concatenation:) It is obvious that  $L_3 = \{ba\}$  is very simple, while  $L_1 L_2 = \{abba\}$  is not, as seen in Example 2, which is shown from (ii) of Lemma 1.

(Intersection:) Let  $L_4 = \{ab^m cd^n e f^n g | m, n \geq 0\}$  and  $L_5 = \{ab^m cd^m e f^n g | m, n \geq 0\}$  be languages generated by very simple grammars  $G$  and  $G'$  with the sets of rules:

$$\begin{array}{l|l} S \rightarrow aABC & A \rightarrow bA \\ A \rightarrow c & B \rightarrow dBD \\ B \rightarrow e & C \rightarrow g \\ D \rightarrow f & \end{array} \quad \begin{array}{l|l} S' \rightarrow aA'B'C' & A' \rightarrow c \\ A' \rightarrow bA'D' & C' \rightarrow fC' \\ B' \rightarrow e & C' \rightarrow g \\ D' \rightarrow d. & \end{array}$$

Then,  $L_4 \cap L_5 = \{ab^n cd^n e f^n g | n \geq 0\}$  is a language which is proven to be noncontext-free.

(Complement:)  $L_6 = \{a\}$  is clearly a very simple language over  $\Sigma$  (containing  $\{a\}$ ), while its complement  $L_6^c = \Sigma^* - L_6$  is not from (ii) of Lemma 1, because  $L_6^c$  contains  $\{a^n | n \geq 2\}$ .

(Kleene closure:) Consider again  $L_6$ . Then,  $L_6^*$  (or  $L_6^+$ ) is not very simple, which is proved from (ii) of Lemma 1.

(Homomorphism:) Consider a very simple language  $L_7 = \{a^n b | n \geq 0\}$  and a homomorphism  $h_1$  defined by  $h_1(a) = a$ ,  $h_1(b) = \lambda$ . Then,  $h_1(L_7) = \{a^n | n \geq 0\}$  is not very simple. Further, consider a  $\lambda$ -free homomorphism  $h_2$  defined by  $h_2(a) = a$ ,  $h_2(b) = a$ . Then,  $h_2(L_7) = \{a^n | n \geq 1\}$  is not very simple.

(Inverse homomorphism:) For the very simple language  $L_6$ , consider a homomorphism  $h$  defined by  $h(b) = a$ ,  $h(c) = \lambda$ . Then,  $h^{-1}(L_6) = \{c^m bc^n | m, n \geq 0\}$  is not very simple, as seen in Example 2.

(Intersection with regular languages:) For the very simple language  $L_7$ , consider  $L_7 \cap (a^3)^* b = \{a^{3n} b | n \geq 0\}$  which is not very simple from Lemma 1.

(Reversal:) For the very simple language  $L_7$ , the reversal  $L_7^R = \{ba^n | n \geq 0\}$  is not very simple, because it does not have the prefix-free property of (i) of Lemma 1.  $\square$

### 3.3. Other properties of very simple grammars

Let  $G = (V_N, \Sigma, P, S)$  be any very simple grammar. Then, a symbol  $a$  in  $\Sigma$  is said to be *terminating* (in  $G$ ) if the  $a$ -handle rule is of the form  $A \rightarrow a$ . A symbol  $a$  in  $\Sigma$  is said to be *starting* if and only if there is a string  $w \in L(G)$  whose first symbol is  $a$ . A symbol  $a$  in  $\Sigma$  is said to be *final* if there is a string  $w \in L(G)$  whose last symbol is  $a$ . Note that the notion of “terminating” depends on  $G$ , while those of “starting” and “final” do not.

The next lemma almost immediately follows from definitions.

**Lemma 5.** *Let  $w, w_1, w_2$  be in a very simple language  $L$ . Then, for any  $G = (V_N, \Sigma, P, S)$  generating  $L$ , and for each  $a, b, c \in \Sigma$ ,*

- (1) *if  $w = ax$  for some  $x \in \Sigma^*$ , then the  $a$ -handle rule is of the form:  $S \rightarrow ax$  for some  $\alpha \in V_N^*$ ,*
- (2) *if  $w = xa$  for some  $x \in \Sigma^*$ , then the  $a$ -handle rule is of the form:  $A \rightarrow a$  for some  $A \in V_N$ ,*
- (3) *if  $\beta \Rightarrow^* a^n$  for some  $\beta \in V_N^+$ ,  $n \geq 1$ , then  $\beta = A^n$ , where  $A$  is the left-hand side of the  $a$ -handle rule of the form:  $A \rightarrow a$ ,*



- (4) if  $w_1 = xby_1$  and  $w_2 = xc y_2$  for some  $x, y_i \in \Sigma^*$ , then there exist  $X \in V_N$  and  $\alpha, \beta \in V_N^*$  such that  $X \rightarrow b\alpha$  and  $X \rightarrow c\beta$  are in  $P$ .

The following are useful for our later discussion.

**Lemma 6.** Let  $G = (V_N, \Sigma, P, S)$  be any very simple grammar and let  $\alpha$  and  $\alpha'$  be in  $V_N^+$ . Further, let  $x \in \Sigma^+, u, u_1, u_2 \in \Sigma^*$ . Then,

- (1) if  $S \Rightarrow^* x\alpha$  and  $S \Rightarrow^* x\alpha'$ , then  $\alpha = \alpha'$  (Forward determinism),
- (2) if  $\alpha \Rightarrow^* x$  and  $\alpha' \Rightarrow^* x$ , then  $\alpha = \alpha'$  (Backward determinism),
- (3) if  $S \Rightarrow^* u_1\alpha \Rightarrow^* u_1x$  and  $S \Rightarrow^* u_2\alpha' \Rightarrow^* u_2x$ , then  $u_1 \setminus L = u_2 \setminus L$ .

**Lemma 7** (Iteration Lemma). Let  $L$  be a very simple language. If  $w = ux^n v$  and  $w' = uv$  are in  $L$  for some  $n \geq 1$ ,  $u \in \Sigma^*$ ,  $x, v \in \Sigma^+$ , then for each  $i \geq 0$ ,  $ux^i v$  is in  $L$ .

Finally, we can easily make a characterization on the relationships among the classes of regular, zero-reversible, and very simple languages.

A regular language  $L$  is zero-reversible if and only if whenever  $u_1x$  and  $u_2x$  are in  $L$ ,  $u_1 \setminus L = u_2 \setminus L$  holds [2].

From the definitions together with Lemma 6, the next theorem almost immediately follows.

**Theorem 8.** For a language  $L$  over  $\Sigma$ , the following are equivalent:

- (1)  $L$  is the left Szilard language of a linear grammar,
- (2)  $L$  is very simple and regular,
- (3)  $L$  is very simple and zero-reversible.

Thus, the intersection of the class of very simple languages and that of regular languages is exactly the class of (left) Szilard languages of linear grammars. (See Fig. 1.)

### 3.4. Schema representations of very simple grammars

We introduce a new representation method for very simple grammars called *grammar schemata*, which enables us to deal with very simple grammars in a uniform manner. Any very simple grammar can be analysed and decomposed into two constructs: a grammar schema and its interpretation. Informally, the former provides a skeletal structure of the original grammar, while the latter specifies the details of it.

Given a finite alphabet  $\Sigma$ , let  $V_{N,S} = \{X_a | a \in \Sigma\} \cup \{S\}$  and let  $PAR (= \{x_a | a \in \Sigma\})$  be a finite set of *parameters*, where  $S$  is a specific symbol not in  $(\{X_a | a \in \Sigma\} \cup \Sigma \cup PAR)$ , and the value of each parameter ranges over all elements from  $V_{N,S}^*$ . Let  $\Gamma = (V_{N,S} \cup PAR)$ . Then, a construct  $X \rightarrow ax$ , where  $X \in V_{N,S}$ ,  $a \in \Sigma$  and  $x \in \Gamma^*$ , is called a *rule form*. We call a quadruple  $\mathcal{G} = (V_{N,S}, \Sigma, P, S)$  a *grammar schema* if  $P$  is a finite set of rule forms.

An *interpretation*  $I = (f_n, f_p)$  over  $\Sigma$  is an ordered pair of mappings, where  $f_n$  is a coding defined on  $V_{N,S}$  such that  $f_n(S) = S$ , and  $f_p$  is a homomorphism defined on  $PAR$



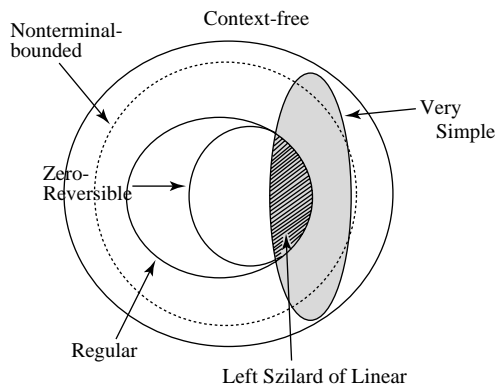


Fig. 1. Language class relations.

such that for all  $x \in PAR$ ,  $f_p(x)$  is in  $\Gamma^*$ . Then, given a grammar schema  $\mathcal{G}$ , let  $I(\mathcal{G})$  be a quadruple defined by  $(f_n(V_{N,S}), \Sigma, I(P), S)$ , where  $I(P) = \{f_n(X) \rightarrow af_p(x) \mid X \rightarrow ax \in P\}$ . Thus,  $I$  induces a mapping on the set of all grammar schemata. An interpretation  $I$  is said to be *ground* if and only if for all  $x \in PAR$ ,  $f_p(x) \in V_{N,S}^*$ . Note that any CFG in Greibach normal form is taken as a grammar schema of special type.

### 3.5. Finding very simple grammars consistent with positive data

In this subsection, we consider the following problem for very simple grammars. Suppose that we are given a finite set of strings  $R = \{w_1, w_2, \dots, w_t\}$  from an unknown very simple language  $L = L(G_*)$  for some very simple grammar  $G_*$ . Starting with constructing the initial grammar schema  $\mathcal{G}_0 = (V_{N,S}, \Sigma, \{X_a \rightarrow ax_a \mid a \in \Sigma\}, S)$ , our goal here is to find (identify) a ground interpretation  $I$  such that  $I(\mathcal{G}_0)$  is consistent with  $R$ , i.e.,  $L(I(\mathcal{G}_0))$  contains  $R$ , where  $\Sigma = \text{alph}(R)$ .

In order to identify such a ground interpretation  $I = (f_n, f_p)$  satisfying the requirements, we have to determine  $f_n$  and  $f_p$ . First, we describe the manner of determining  $f_n$ .

#### [Computation of $f_n$ ]

We will demonstrate the way how to determine a mapping  $f_n$  from a graph called the *structure graph of  $R$*  denoted by  $T_R$ .

Let

$$\Sigma_s(R) = \{a \in \Sigma \mid \exists w \in R, \exists x \in \Sigma^* (w = ax)\}.$$

$$\Sigma_f(R) = \{a \in \Sigma \mid \exists w \in R, \exists x \in \Sigma^* (w = xa)\}.$$

Then, we may define that for all  $a \in \Sigma_s(R)$ ,  $f_n(X_a) = S$ . In order to determine  $f_n(X)$  for other  $X \in V_{N,S} - (\{X_a \mid a \in \Sigma_s(R)\} \cup \{S\})$ , we construct a directed graph called the *structure graph  $T_R$*  of  $R$  as follows.

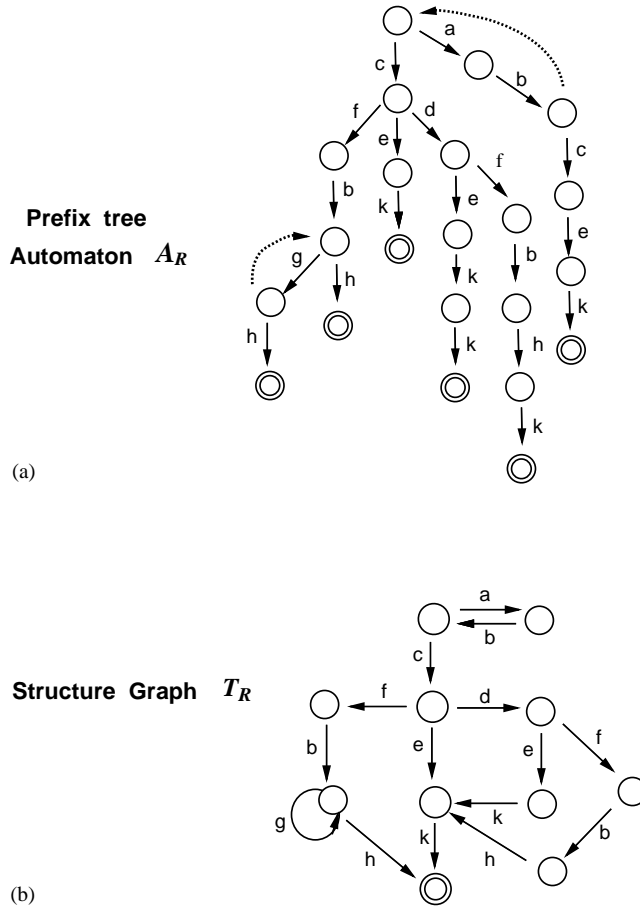


Fig. 2. Prefix tree automaton  $A_R$  and structure graph  $T_R$ . In (a) two nodes connected by a dotted arrow are to be identified.

From a given  $R$ , we first construct the prefix tree automaton  $A_R$  for  $R$  in the sense of the paper [2]. Then, using the properties of very simple languages presented in Lemmas 6 and 7, from  $A_R$  we construct the structure graph  $T_R$ . Once the graph  $T_R$  is obtained, a mapping  $f_n$  is immediately determined from  $T_R$ . Let us illustrate this procedure using an example.

**Example 3.** Suppose that  $R = \{abcek, cdek^2, cek, cdfbhk, cfbh, cfbgh\}$  ( $= \{w_1, \dots, w_6\}$ ) is given as a sample set of  $L(G_*)$  for some unknown very simple grammar  $G_*$ . Then, we have the prefix tree automaton  $A_R$  for  $R$  shown in (a) of Fig. 2. Note that  $\Sigma = \text{alph}(R) = \{a, b, c, d, e, f, g, h, k\}$ .

Then, making use of properties of very simple grammars in Lemma 6, we merge all the nodes (nonterminals) which must be identical in any correct grammar for  $L(G_*)$ .

By Lemma 7, we know that  $ab$  and  $g$  are loopings. Hence, we eventually obtain the structure graph of  $R$ ,  $T_R$ , shown in (b) of Fig. 2.

Now, from 4 of Lemma 5, we know that if terminal symbols  $s$  and  $t$  are the labels of edges that have an identical starting node in  $T_R$ , then  $X_s$  and  $X_t$  must be identical. Note that, for our example,  $\Sigma_s(R) = \{a, c\}$  and  $\Sigma_f(R) = \{h, k\}$ . Hence, we have that  $X_a = X_c (= S)$ ,  $X_d = X_e = X_f$ , and  $X_g = X_h$ . In this example, we define  $f_n$  as follows:

$$f_n(X_a) = f_n(X_c) = S, \quad f_n(X_e) = f_n(X_f) = X_d, \quad f_n(X_h) = X_g,$$

$$\text{and for other } X \in \{X_b, X_d, X_g, X_k\}, \quad f_n(X) = X.$$

Thus, in what follows, we assume that  $f_n$  is defined so that for each  $X$  in  $V_{N,S} - (\{X_a | a \in \Sigma_s\} \cup \{S\})$  the image  $f_n(X)$  is the one with the first index in alphabetical order, among those to be indentified.

**Note 1.** As described above,  $f_n$  is determined from the structure graph  $T_R$  whose construction is solely due to the properties proved valid in Lemmas 5, 6 and 7, which assures that any very simple grammar consistent with  $R$  must satisfy those properties.

#### [Computation of $f_p$ ]

Recall that

$$\Sigma_f(R) = \{a \in \Sigma | \exists w \in R, \exists x \in \Sigma^* (w = xa)\}.$$

We say that a parameter  $x_a$  is *empty* if it is  $\lambda$ . Obviously, for all  $a \in \Sigma_f(R)$ ,  $x_a$  must be empty. Also, for all  $a \in \Sigma_s(R) - \Sigma_f(R)$ ,  $x_a$  is *not empty* (i.e.,  $a$  is not terminating). These facts are used for determining  $f_p$  below.

In what follows, first we are going to present procedures and definitions, and then we will show a sample run to illustrate them.

Determining unknown values of  $f_p$  is performed in the following three steps.

*Step 1: Constructing Length Equation Set Lg(R).* Let  $\text{alph}(R) = \Sigma = \{a_1, \dots, a_r\}$  be an ordered alphabet. For  $w \in \Sigma^*$  and  $a \in \Sigma$ ,  $\#_a(w)$  denotes the number of occurrences of  $a$  in  $w$ .

For all  $a \in \Sigma$ , let  $n_a = \text{len}(f_p(x_a)) - 1$ , where  $f_p(x_a)$  is the value (a string of nonterminals) to be determined for the unknown parameter  $x_a$ . (Note that  $n_a$  indicates the length increase of a sentential form induced by one application of the  $a$ -handle rule. Hence, it holds that  $n_a = -1$  if and only if  $f_p(x_a) = \lambda$ .) Further, let  $R_a = \{w \in R | w \text{ contains } a\}$ . Then, we have that for any  $a \in \Sigma$ ,

(C1)  $-1 \leq n_a \leq B(a, w)$  (where  $B(a, w) = \text{len}(w) / \#_a(w)$  and  $w$  is a string of the minimum length in  $R_a$ ),

(C2) if  $a \in \Sigma_s(R) - \Sigma_f(R)$ , then  $n_a \geq 0$ .

It is easily seen that for each  $w_i \in R$ ,

$$\#_{a_1}(w_i)n_{a_1} + \dots + \#_{a_r}(w_i)n_{a_r} = -1 \quad \dots \quad (\ell.w_i).$$

By constructing the equation  $(\ell.w_i)$  for each  $w_i$  in  $R$  in the manner described above, we have a set of equations  $\text{Lg}(R) = \{(\ell.w_1), \dots, (\ell.w_t)\}$ . Since for  $a \in \Sigma_f(R)$ ,  $f_p(x_a) = \lambda$ , we know that  $n_a = -1$ .

Let  $\text{Lg}(R') = \{(\ell.w_i)' | 1 \leq i \leq t\}$  be the set of equations  $(\ell.w_i)'$  obtained from  $(\ell.w_i)$  by substituting  $n_a = -1$  for each  $a \in \Sigma_f(R)$ .

In order to find a concrete very simple grammar  $G$  which is consistent with the given positive data  $R$ , first we have to solve the set of equations  $\text{Lg}(R')$ , so that we may obtain a possible ground interpretation  $I = (f_n, f_p)$  satisfying the requirements.

*Step 2: Solving Equation Set  $\text{Lg}(R')$ .* Let  $\Sigma' (= \{a_1, \dots, a_m\})$  be the ordered set  $\text{alph}(R) - \Sigma_f(R)$  and let  $\text{LPAR} = \{n_a | a \in \Sigma'\}$  be the set of length parameters to be solved. Further, let  $\text{vec}(w_i)$  be a row vector defined by

$$\text{vec}(w_i) = (\#_{a_1}(w_i), \dots, \#_{a_m}(w_i)),$$

where, each  $\#_{a_j}(w_i)$  is a coefficient of  $n_{a_j}$  in the left-hand side of  $(\ell.w_i)'$ .

Then, construct a  $(t \times m)$ -matrix  $M_R$  associated with  $R$  as follows:

$$(i, j)\text{-entry of } M_R \stackrel{\text{def}}{=} \begin{array}{l} \text{the coefficient of the } j\text{th parameter } n_{a_j} \text{ in} \\ \text{the } i\text{th equation } (\ell.w_i)' \end{array}$$

that is,

$$M_R = \begin{pmatrix} \text{vec}(w_1) \\ \vdots \\ \text{vec}(w_t) \end{pmatrix}.$$

(Note that  $t = |R|$  and  $m = |\text{LPAR}|$ .) Then, the set of equations  $\{(\ell.w_i)' | 1 \leq i \leq t\}$  is reformulated by a matrix equation:

$$M_R \mathcal{X}^T = \mathcal{C}^T \quad \dots \quad (\text{E})$$

where  $\mathcal{X} = (n_{a_1}, \dots, n_{a_m})$  is a vector of length parameters,  $\mathcal{C} = (k'_1, \dots, k'_t)$  is a constant vector and each  $k'_j$  is the constant in the right-hand side of  $(\ell.w_i)'$ .

Solving the matrix equation (E) is performed in a usual manner in linear algebra. In some case, we may have indeterminate parameters in the solution. The next lemma is one of the well-known results in linear algebra.

**Lemma 9.** *The equation (E) has a solution if and only if the rank of  $M_R$  is equal to that of  $M_R \mathcal{C}^T$ .*

Let  $m$  and  $t$  be the numbers of parameters and equations in (E), respectively. Then, we also know the following.

**Lemma 10.** *Suppose that  $m = t$  holds. Then,  $M_R$  is nonsingular if and only if (E) has a unique solution.*

Thus, if  $m = \text{rank}(M_R)$ , then the equation (E) has a unique solution. And, if  $m > t$  or  $m \neq \text{rank}(M_R)$ , then (E) is solved in part. That is, let  $\text{Sol}(\text{Lg}(R'))$  be the set of solutions of  $\text{Lg}(R')$ . Then, we have that

$$\text{Sol}(\text{Lg}(R')) = \text{P-Sol}(\text{Lg}(R')) \cup \text{Uns}(\text{Lg}(R')),$$

where  $\text{P-Sol}(\text{Lg}(R'))$  is the set of solutions partly solved, and  $\text{Uns}(\text{Lg}(R'))$  is the set of equations involved in indeterminate parameters left unsolved.

Then, by assigning an appropriate value to each indeterminate parameter in  $\text{Uns}(\text{Lg}(R'))$ , it is always possible to have a complete solution for  $\text{Lg}(R')$ .

*Step 3: Determining the values of  $f_p$ .*

(3-1) First, we may define  $f_p$  in part by

$$\text{for all } a \text{ in } \Sigma_f(R), \quad f_p(x_a) = \lambda.$$

(3-2) Determining the values of  $f_p$  for other parameters is performed by *simulating* the derivation for each  $w$  in  $R$  as follows:

Let  $\mathcal{X} = (\ell_{a_1}, \dots, \ell_{a_m})$  be a solution vector obtained from (E) in (Step 2). Then, for  $a \in \Sigma'$  such that  $\ell_a \geq 0$ , we may assume that  $f_p(x_a) = Z_{a,1} \cdots Z_{a,\ell_a+1}$ . For  $a \in \Sigma'$  such that  $\ell_a = -1$ , we may trivially define  $f_p(x_a) = \lambda$ . The set  $\{f_n(X_a) \rightarrow a f_p(x_a) \mid a \in \Sigma'\}$  is called the *set of candidate rules* obtained from  $\mathcal{X}$ .

In order to identify each unknown nonterminal  $Z_{a,j}$  introduced above, for each  $w = b_1 \cdots b_k \in R$ , we simulate the derivation for  $w$  using each candidate  $b_i$ -handle rule:  $f_n(X_{b_i}) \rightarrow b_i Z_{b_i,1} \cdots Z_{b_i,\ell_{b_i}+1}$  or  $f_n(X_{b_i}) \rightarrow b_i$  (for  $i = 1, \dots, k$ ).

Starting with  $S \Rightarrow b_1 Z_{b_1,1} \cdots Z_{b_1,\ell_{b_1}+1}$ , since  $Z_{b_1,1}$  must produce  $b_2 u$  (for some  $u \in \Sigma^*$ ), it must hold that  $Z_{b_1,1} = f_n(X_{b_2})$  and that  $f_n(X_{b_2}) \rightarrow b_2 Z_{b_2,1} \cdots Z_{b_2,\ell_{b_2}+1}$ . Thus, we have that

$$\begin{aligned} S &\Rightarrow b_1 Z_{b_1,1} \cdots Z_{b_1,\ell_{b_1}+1} \\ &\Rightarrow b_1 b_2 Z_{b_2,1} \cdots Z_{b_2,\ell_{b_2}+1} Z_{b_1,2} \cdots Z_{b_1,\ell_{b_1}+1} \\ & (= b_1 b_2 Z_{b_1,3} \cdots Z_{b_1,p}). \end{aligned}$$

In the same manner, for each  $j$  ( $3 \leq j \leq p$ ), we can determine the value of  $Z_{b_i,j}$ . That is, for each  $i = 1, \dots, k$ ,  $f_p(x_{b_i})$  is determined. Thus, a *simulation of the derivation* for  $w$  is completed.

Determining complete values of  $f_p$  is performed by combining the results of simulation of the derivations for all  $w \in R$ .

From Lemma 9, a corollary immediately follows.

**Corollary 11.** *It is decidable given any finite subset  $R$  of  $\Sigma^*$  whether or not there exists a very simple grammar consistent with  $R$ .*

(Recall that a finite set is not necessarily very simple.)

[Procedure **Consistent**( $R$ )]

By  $\text{Int}(R)$  we denote the set of all ground interpretations  $I = (f_n, f_p)$  obtained from  $R$  in the manner described above. Let  $\text{CR}(\mathcal{X})$  be the set of candidate rules determined from a solution vector  $\mathcal{X}$  of (E).

There are, in general, two cases for the simulation process of  $R$  via  $\text{CR}(\mathcal{X})$ : the simulation finishes without any *additional* nonterminal identification (i.e., merging two nonterminals) on  $f_n$ , or it forces us to make a new nonterminal identification. In the former case,  $\text{CR}(\mathcal{X})$  is said to be *good*.

In order to obtain the least general grammar  $I(\mathcal{G}_0)$  consistent with  $R$ , we want to *avoid* choosing an interpretation  $I = (f_n, f_p)$  that causes the latter case, as much as we possibly can.

We say that a ground interpretation  $I = (f_n, f_p)$  in  $\text{Int}(R)$  is *admissible to  $R$*  if and only if either there exists a good  $\text{CR}(\mathcal{X})$  by which  $f_p$  is determined or there exists no good set of candidate rules.

Our goal here is to find a ground interpretation  $I$  admissible to  $R$ .

Now, let us define a *procedure Consistent*( $R$ ) to obtain  $I$  admissible to  $R$  from a given  $R$ . That is, **Consistent**( $R$ ) is defined as follows:

*Input*:  $R$  (a finite subset of the target language  $L_*$ );

*Output*:  $I = (f_n, f_p)$  admissible to  $R$ ;

*Procedure*:

initialize  $F_v = \emptyset$ ;

determine  $f_n$  by constructing the structure graph of  $R$ ;

construct  $\text{Lg}(R)$  and solve it, where for each  $a \in \Sigma_f(R)$ ,  $n_a = -1$ ;

let  $S_v$  be the set of all solution vectors  $\mathcal{X}$  satisfying (C1) and (C2);

**while**  $S_v \neq \emptyset$  **do**

**begin**

    choose any solution  $\mathcal{X}$  from  $S_v$  to obtain the values for the rest  
    of unspecified parameters of  $f_p$ ;

    let  $S_v := S_v - \{\mathcal{X}\}$  and  $F_v := F_v \cup \{\mathcal{X}\}$ ;

    determine  $\text{CR}(\mathcal{X})$  (the set of candidate rules) from  $\mathcal{X}$ ;

    simulate  $R$  with  $\text{CR}(\mathcal{X})$  to determine  $f_p$ ;

**if**  $\text{CR}(\mathcal{X})$  is good, **then** output  $I = (f_n, f_p)$  and halt;

**end**

take any solution  $\mathcal{X}$  from  $F_v$ ;

let  $f_p$  be the one obtained by simulating  $R$  via  $\text{CR}(\mathcal{X})$  (where additional nonterminal identifications are allowed);

renew  $f_n$  by making additional identifications obtained through the simulation of  $R$  via  $\text{CR}(\mathcal{X})$ ;

output  $I = (f_n, f_p)$  and halt.

We write the input–output relation as  $I = \mathbf{Consistent}(R)$ . (Note that from the manner of construction,  $I(= \mathbf{Consistent}(R))$  can obviously provide a very simple grammar  $I(\mathcal{G}_0)$  that is consistent with  $R$ .)

**Note 2.** In the actual implementation of the procedure above, one can make **Consistent** ( $R$ ) more efficient in an incremental fashion so that it may produce and try a new solution vector  $\mathcal{X}$  only after the current  $\text{CR}(\mathcal{X})$  turned out to be not good. We chose the current presentation for **Consistent**( $R$ ), because of its simplicity and understandability.

**Example 4.** Returning to the previous example, let  $R$  be the sample set from Example 3. Let  $\mathcal{G}_0 = (V_{N,S}, \Sigma, P_0, S)$ , where  $V_{N,S} = \{S, X_a, X_b, X_c, X_d, X_e, X_f, X_g, X_h, X_k\}$ . Since we have already obtained  $f_n$ , all we have to do is to construct  $\text{Lg}(R)$  and solve it for finding  $f_p$  which satisfies the requirements.

Firstly, since  $\Sigma_f(R) = \{h, k\}$ , we have already obtained that

$$f_p(x_h) = f_p(x_k) = \lambda, \quad (\text{i.e., } n_h = n_k = -1).$$

From  $R$ , we have that  $\text{Lg}(R) = \{(\ell.w_1), (\ell.w_2), \dots, (\ell.w_6)\}$ , where

$$\begin{aligned} n_a + n_b + n_c + n_e + n_k &= -1 & \cdots & (\ell.w_1) \\ n_c + n_d + n_e + 2n_k &= -1 & \cdots & (\ell.w_2) \\ n_c + n_e + n_k &= -1 & \cdots & (\ell.w_3) \\ n_b + n_c + n_d + n_f + n_h + n_k &= -1 & \cdots & (\ell.w_4) \\ n_b + n_c + n_f + n_h &= -1 & \cdots & (\ell.w_5) \\ n_b + n_c + n_f + n_g + n_h &= -1 & \cdots & (\ell.w_6). \end{aligned}$$

That is, the set of equations  $\text{Lg}(R')$  comprising

$$\begin{aligned} n_a + n_b + n_c + n_e &= 0 & \cdots & (\ell.w_1)' \\ n_c + n_d + n_e &= 1 & \cdots & (\ell.w_2)' \\ n_c + n_e &= 0 & \cdots & (\ell.w_3)' \\ n_b + n_c + n_d + n_f &= 1 & \cdots & (\ell.w_4)' \\ n_b + n_c + n_f &= 0 & \cdots & (\ell.w_5)' \\ n_b + n_c + n_f + n_g &= 0 & \cdots & (\ell.w_6)' \end{aligned}$$

is obtained, where  $\Sigma' (= \text{alph}(R) - \Sigma_f(R)) = \{a, b, c, d, e, f, g\}$ . Then, construct a matrix equation

$$M_R \mathcal{X}^T = \mathcal{C}^T \dots \quad (\text{E})$$

where  $\mathcal{X} = (n_a, n_b, n_c, n_d, n_e, n_f, n_g)$ ,  $\mathcal{C} = (0, 1, 0, 1, 0, 0)$ , and

$$M_R = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 \end{pmatrix}.$$



Computing (E) is performed in a usual manner:

$$\left( \begin{array}{cccccc|c} 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 \end{array} \right) \Rightarrow^* \left( \begin{array}{cccccc|c} 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{array} \right).$$

Thus, we have that  $\text{Sol}(\text{Lg}(R')) = \{n_d = 1, n_g = 0\} \cup \text{Uns}(\text{Lg}(R'))$ , where  $\text{Uns}(\text{Lg}(R')) = \{n_a + n_b = 0, n_c + n_e = 0, n_b + n_c + n_f = 0\}$ .

By taking (C1) and (C2) into consideration, we may choose a set of solutions  $\{n_a = 1, n_b = -1, n_c = 0, n_e = 0, n_f = 1\}$  that produces a solution vector:

$$\mathcal{X} = (1, -1, 0, 1, 0, 1, 0) = (n_a, n_b, n_c, n_d, n_e, n_f, n_g)$$

In order to determine the values of  $f_p$ , the set of candidate rules  $\text{CR}(\mathcal{X})$  is constructed:

$$\begin{array}{ll} S \rightarrow aZ_{a,1}Z_{a,2}, & X_b \rightarrow b \\ S \rightarrow cZ_{c,1}, & X_{td} \rightarrow dZ_{d,1}Z_{d,2} \\ X_d \rightarrow eZ_{e,1}, & X_d \rightarrow fZ_{f,1}Z_{f,2} \\ X_g \rightarrow gZ_{g,1}. & \end{array}$$

(Note that  $f_n(X_a) = f_n(X_c) = S$ ,  $f_n(X_e) = f_n(X_f) = X_d$ .)

In simulating the derivation for  $w_1 = abcck$ , since  $S \Rightarrow aZ_{a,1}Z_{a,2} \Rightarrow^* abcck$ , it must hold that  $Z_{a,1} = X_b$  and that  $Z_{a,2} = X_c (= S)$ . Further, from  $S \Rightarrow cZ_{c,1} \Rightarrow^* cek$ , it must hold that  $Z_{c,1} = X_e (= X_d)$  and that  $Z_{e,1} = X_k$ . Thus, we have that

$$f_p(x_a) = Z_{a,1}Z_{a,2} = X_bS, \quad f_p(x_c) = X_d \quad \text{and} \quad f_p(x_e) = X_k.$$

(With these assignments for length parameters, we succeed in simulating the derivations for  $w_2$  and  $w_3$ .)

In a similar manner, from the simulation of  $w_4$ ,

$$f_p(x_d) = X_dX_k, \quad f_p(x_f) = X_bX_g \quad \text{and} \quad f_p(x_g) = X_g$$

are obtained. Thus, we complete a solution for  $f_p$  using  $\text{CR}(\mathcal{X})$ . (Note that  $\text{CR}(\mathcal{X})$  is good.) Further, we can see that any other solution fails in simulating  $R$  without any additional nonterminal identification on  $f_n$  in this sample run.

This eventually provides the unique interpretation  $I = (f_n, f_p)$  admissible to  $R$ , where

$$\begin{aligned} f_n(X_a) &= S, & f_p(x_a) &= X_b S \\ f_n(X_b) &= X_b, & f_p(x_b) &= \lambda \\ f_n(X_c) &= S, & f_p(x_c) &= X_d \\ f_n(X_d) &= X_d, & f_p(x_d) &= X_d X_k \\ f_n(X_e) &= X_d, & f_p(x_e) &= X_k \\ f_n(X_f) &= X_d, & f_p(x_f) &= X_b X_g \\ f_n(X_g) &= X_g, & f_p(x_g) &= X_g \\ f_n(X_h) &= X_g, & f_p(x_h) &= \lambda \\ f_n(X_k) &= X_k, & f_p(x_k) &= \lambda. \end{aligned}$$

The rule set of  $I(\mathcal{G}_0)$  consists of

$$\begin{aligned} S &\rightarrow aX_b S, & X_b &\rightarrow b, & S &\rightarrow cX_d, \\ X_d &\rightarrow dX_d X_k, & X_d &\rightarrow eX_k, & X_d &\rightarrow fX_b X_g, \\ X_g &\rightarrow gX_g, & X_g &\rightarrow h, & X_k &\rightarrow k. \end{aligned}$$

#### 4. Identifying very simple grammars from positive data

##### 4.1. Characterizing the structure of $\mathcal{I}_\Sigma$

Throughout this section, we assume that  $G_* = (V_{N_*}, \Sigma, P_*, S_*)$  is a target very simple grammar and that  $R$  is a finite subset of  $L(G_*)$ . Further, let  $\mathcal{I}_\Sigma = \{I = (f_n, f_p) \mid I: \text{ground interpretation over } \Sigma\}$ , where  $|\Sigma| \geq 2$ . (It is obvious that any (nonempty) very simple language over the unary alphabet trivially becomes  $\{a\}$ .)

Recall that  $\mathcal{G}_0 = (V_{N,S}, \Sigma, \{X_a \rightarrow ax_a \mid a \in \Sigma\}, S)$  is the initial grammar schema. Now, for any  $I = (f_n, f_p)$ ,  $I' = (f'_n, f'_p) \in \mathcal{I}_\Sigma$ , we define a binary relation  $\preceq$  as follows:  $f_n \preceq f'_n$  if and only if for every  $X_a, X_b \in V_{N,S} (a \neq b)$ ,  $f_n(X_a) = f_n(X_b)$  implies that  $f'_n(X_a) = f'_n(X_b)$ . When  $f_n \preceq f'_n$  and  $f_n \neq f'_n$ , we denote it by  $f_n \prec f'_n$ . Further, we say that  $f_n$  is *incomparable* to  $f'_n$  if and only if  $f_n \not\preceq f'_n$  and  $f'_n \not\preceq f_n$  holds, where the symbol  $\not\preceq$  denotes the negation of  $\preceq$ . (Recall that we assume all very simple grammars to be reduced.)

We now need to show the following series of technical lemmas.

**Lemma 12.** *Let  $I = (f_n, f_p)$  and  $I' = (f'_n, f'_p)$  be in  $\mathcal{I}_\Sigma$  such that  $L(I(\mathcal{G}_0)) \subseteq L(I'(\mathcal{G}_0))$ . Then, it holds that  $f_n \preceq f'_n$ .*

**Proof.** Let  $I(\mathcal{G}_0) = G = (V_N, \Sigma, P, S)$  and  $I'(\mathcal{G}_0) = G' = (V'_N, \Sigma, P', S')$  be very simple grammars such that  $L(G) \subseteq L(G')$ . Suppose that  $f_n \not\preceq f'_n$ , then we may assume that there

exist  $X_a, X_b$  ( $a \neq b$ ) in  $V_{N,S}$  of  $\mathcal{G}_0$  such that  $f_n(X_a) = f_n(X_b)$  and  $f'_n(X_a) \neq f'_n(X_b)$ . There exists a string  $w = xay$  in  $L(G)$  such that

$$S \Rightarrow^* xA\alpha, A \Rightarrow a\gamma \quad \text{and} \quad \gamma\alpha \Rightarrow^* y \quad \text{for some } \alpha, \gamma \in V_N^*, \quad \text{where } A = f_n(X_a).$$

Further, since  $f_n(X_a) = f_n(X_b)$ , there also exists a string  $w' = xbz$  in  $L(G)$  such that

$$S \Rightarrow^* xA\alpha, A \Rightarrow b\beta \quad \text{and} \quad \beta\alpha \Rightarrow^* z \quad \text{for some } \beta \in V_N^*, \quad \text{where } A = f_n(X_b).$$

Since  $L(G) \subseteq L(G')$ , there must exist derivations in  $G'$  such that

$$S' \Rightarrow^* x\alpha', \alpha' \Rightarrow^* ay \quad \text{and} \quad \alpha' \Rightarrow^* bz \quad \text{for some } \alpha' \in V_N'^+.$$

Let  $X$  be the left-most nonterminal of  $\alpha'$ , then there must exist rules  $X \rightarrow a\beta'$  and  $X \rightarrow b\gamma'$  in  $P'$  of  $G'$ , which leads to a contradiction that  $X = f'_n(X_a) \neq f'_n(X_b) = X$  in  $V_N'$  of  $G'$ .  $\square$

**Corollary 13.** *Let  $I = (f_n, f_p)$  and  $I' = (f'_n, f'_p)$  be in  $\mathcal{I}_\Sigma$ . Then,  $L(I(\mathcal{G}_0)) = L(I'(\mathcal{G}_0))$  implies  $f_n = f'_n$ .*

**Lemma 14.** *Let  $I = (f_n, f_p)$  and  $I' = (f'_n, f'_p)$  be in  $\mathcal{I}_\Sigma$ . Then,*

- (i) *if  $L(I(\mathcal{G}_0)) \subset L(I'(\mathcal{G}_0))$ , then  $f_n \prec f'_n$ ,*
- (ii) *if  $f_n = f'_n$  and  $L(I(\mathcal{G}_0)) \neq L(I'(\mathcal{G}_0))$ , then  $L(I(\mathcal{G}_0))$  and  $L(I'(\mathcal{G}_0))$  are incomparable.*

**Proof.** Let  $G = I(\mathcal{G}_0) = (V_N, \Sigma, P, S)$  and  $G' = I'(\mathcal{G}_0) = (V_N', \Sigma, P', S')$ . (i) Assume that  $L(G) \subset L(G')$ . Then, by Lemma 12, it holds that  $f_n \preceq f'_n$ . Now, by assuming that  $f_n = f'_n$ , we will derive a contradiction.

From the assumption that  $L(G) \subset L(G')$ , there exists a string  $w_0 = a_1 a_2 \cdots a_k \in L(G') - L(G)$ . Let

$$S \Rightarrow a_1 \alpha_1 \Rightarrow^* a_1 \cdots a_j \alpha_j \Rightarrow^* a_1 \cdots a_{k-1} \alpha_{k-1} \Rightarrow a_1 \cdots a_k \quad \cdots \quad (*)$$

be a derivation for  $w_0$  in  $G'$ , where each  $\alpha_j \in V_N'^+$ . Since  $G$  fails to generate  $w_0$ , there are three possible cases in  $G$ .

*Case 1.* There exists  $j$  ( $1 \leq j \leq k$ ) such that  $S \Rightarrow^* a_1 \cdots a_{j-1} f_n(X_b) \beta$  and  $f_n(X_b) \neq f_n(X_{a_j})$  ( $b \neq a_j$ ), for some  $b \in \Sigma, \beta \in V_N^*$  in  $G$ : Let  $f_n(X_b) \rightarrow b\gamma$  be in  $P$  and  $\gamma\beta \Rightarrow^* u$  ( $u \in \Sigma^*$ ). (Recall our initial assumption that very simple grammars considered are all reduced.) Then,  $S \Rightarrow^* a_1 \cdots a_{j-1} bu (= w) \in L(G)$ . On the other hand, since  $L(G) \subset L(G')$ , there also exists a derivation for  $w$  in  $G'$ . Let  $S' \Rightarrow^* a_1 \cdots a_{j-1} \alpha'$  and  $\alpha' \Rightarrow^* bu$  in  $G'$ , where  $\alpha' \in V_N'^+$ . Then, this together with (\*) implies that  $\alpha' = \alpha_{j-1}$ ,  $\alpha' \Rightarrow^* bu$  and  $\alpha' \Rightarrow^* a_j \cdots a_k$ . Therefore, we have that  $f'_n(X_b) = f'_n(X_{a_j})$  (and  $b \neq a_j$ ). Since  $f_n = f'_n$ , this also holds for  $G$ , i.e.,  $f_n(X_b) = f_n(X_{a_j})$ , which is a contradiction.

*Case 2.* There exists  $\beta \in V_N^+$  such that  $S \Rightarrow^* w_0 \beta$  and  $\beta \Rightarrow^* x$  ( $x \in \Sigma^+$ ) in  $G$ . Then,  $w_0 x$  is in  $L(G)$ . Since  $L(G) \subset L(G')$ , it holds that  $w_0 x$  is in  $L(G')$ . This contradicts the prefix-free property of a very simple language  $L(G')$ .

*Case 3.* There exists  $j$  ( $1 \leq j < k$ ) such that  $S \Rightarrow^* a_1 \cdots a_j (= y)$  in  $G$ : Since  $L(G) \subset L(G')$ , it holds that  $y$  is in  $L(G')$ . This together with the fact that  $ya_{j+1} \cdots a_k (= w_0)$  is in  $L(G')$  contradicts the prefix-free property. Thus, we have that  $f_n \prec f'_n$ .

(ii) Suppose that  $L(G) \subset L(G')$ . Then, from (i) we have that  $f_n \prec f'_n$ , contradicting  $f_n = f'_n$ . Hence, it holds that  $L(G) \not\subset L(G')$ . Similarly, we also have that  $L(G') \not\subset L(G)$ . These, together with  $L(G) \neq L(G')$ , imply that  $L(G)$  and  $L(G')$  are incomparable. Thus, the proof is completed.  $\square$

**Lemma 15.** For a finite subset  $R$  of  $L(G_*)$ , let  $I = (f_n, f_p) = \mathbf{Consistent}(R)$  and  $G = I(\mathcal{G}_0)$ . Then, either  $L(G_*) = L(G)$  or  $L(G_*) - L(G) \neq \emptyset$ .

**Proof.** Suppose otherwise, that is,  $L(G_*) \subset L(G)$ . (We prove by contradiction.) Then, from (i) of Lemma 14, we have that  $f_{n*} \prec f_n$ . (Note that  $R \subseteq L(G_*)$  and  $I = (f_n, f_p) = \mathbf{Consistent}(R)$ .) On the other hand, from the manner of constructing  $\mathbf{Consistent}(R)$ , there are two cases to be considered.

*Case 1.* A good  $\text{CR}(\mathcal{X})$  was found in the **while** loop:  $f_n$  is determined from the structure graph  $T_R$  whose construction is solely due to the properties common to all very simple grammars consistent with  $R$  and proved valid in Lemmas 5–7 (see Note 1 immediately after Example 3). This implies that  $f_n \preceq f_{n*}$ .

*Case 2.* Otherwise: In this case, since no good  $\text{CR}(\mathcal{X})$  was found for determining  $f_p$ , there must occur additional nonterminal identifications in simulating  $R$ . Let  $\tilde{f}_n$  be the one, constructed from the structure graph  $T_R$  alone, from which the renewed version  $f_n$  is obtained by making such additions. Then, it holds that  $\tilde{f}_n \prec f_{n*}$  and  $\tilde{f}_n \prec f_n$ , and we see that  $\tilde{f}_n$  is the greatest lower bound of  $\{f_n, f_{n*}\}$ , i.e.,  $\tilde{f}_n = \inf\{f_n, f_{n*}\}$  in the lattice with  $\preceq$  consisting of all possible  $f_n$ s with the maximum  $\tilde{f}_n$  (where for any  $X \in V_{N,S}$ ,  $\tilde{f}_n(X) = S$ ) and the minimum  $\underline{f}_n$  (where  $\underline{f}_n$  is an identity on  $V_{N,S}$ ). Since  $f_{n*} \prec f_n$ , it holds that  $\tilde{f}_n \prec f_{n*} \prec f_n$ , which contradicts that  $\tilde{f}_n$  is the greatest lower bound of  $\{f_n\}$ .  $\square$

#### 4.2. Identification algorithm: IA

We now present an identification algorithm *IA* which is consistent, responsive and conservative [1], that is, intuitively, it always produces a guess consistent with the sample set  $R$  and makes at least one guess in response to each input before requesting another input, and only changes its guess when it conflicts with the sample set read in so far. The algorithm *IA* is given in Fig. 3, in which  $\mathcal{G}_{0,\Sigma}$  denotes  $(\{S\} \cup V_{N,\Sigma}, \Sigma, P_\Sigma, S)$ , where  $V_{N,\Sigma} = \{X_a \mid a \in \Sigma\}$ ,  $P_\Sigma = \{X_a \rightarrow ax_a \mid a \in \Sigma\}$ .

**Lemma 16.** Let  $G_{R_0}, G_{R_1}, \dots, G_{R_i}, \dots$  be the sequence of conjectured grammars produced by *IA*, where  $G_{R_i} = I_i(\mathcal{G}_{0,\Sigma})$ . Then, there exists  $r \geq 0$  such that for all  $i \geq 0$ ,  $G_{R_r} = G_{R_{r+i}}$  and  $L(G_{R_r}) = L(G_*)$ .

**Proof.** From the property of *IA*, in particular, of the procedure  $\mathbf{consistent}(R)$ , there is an upper bound  $B$  (depending on the size of  $G_*$ ) for which for each  $i \geq 1$ , the number of candidate interpretations  $I_i (= \mathbf{consistent}(R_i))$  for the  $i$ -th conjecture  $G_{R_i}$  is no more than  $B$ . (Note that for each  $i \geq 1$ , and interpretation  $I_*$  for  $G_*$  is potentially included in the set of those candidate interpretations  $I_i$ .) Thus, there exists  $r \geq 0$  such that for

---

**Input:** a positive presentation of a very simple language  $L(G_*)$   
**Output:** a sequence of very simple grammars  $G_{R_0}, G_{R_1}, \dots$   
**Procedure**  
 initialize  $R_0 = \emptyset$ ;  
 initialize the grammar schema  $\mathcal{G}_{0, \emptyset}$ ;  
 let  $G_{R_0} = (\{S\}, \emptyset, \emptyset, S)$ ;  
 let  $i = 1$ ;  
**repeat** (forever)  
   read the next positive example  $w_i$ ;  
   let  $R_i = R_{i-1} \cup \{w_i\}$ ;  
   let  $\text{alph}(R_i) = \text{alph}(R_{i-1}) \cup \{\text{alph}(w_i)\}$ ;  
   **if**  $w_i \in L(G_{R_{i-1}})$ , **then** let  $G_{R_i} = G_{R_{i-1}}$ ;  
   output  $G_{R_i}$ ;  
   **else**  
     augment  $\mathcal{G}_{0, \Sigma}$  using  $\Sigma = \text{alph}(R_i)$ ;  
     let  $I_i = \mathbf{Consistent}(R_i)$ ;  
     output  $G_{R_i} = I_i(\mathcal{G}_{0, \Sigma})$ ;  
   let  $i := i + 1$ ;

---

Fig. 3. The identification algorithm  $IA$ .

all  $i \geq 0$ ,  $G_{R_r} = G_{R_{r+i}}$ . Suppose that  $L(G_*) \neq L(G_{R_r})$ , then by Lemma 15, there exists a string  $w \in L(G_*) - L(G_{R_r})$  such that  $w$  is not yet provided as a positive example. This implies that  $IA$  produces a conjecture distinct from  $G_{R_r}$ , a contradiction.  $\square$

Thus, we have the following.

**Theorem 17.** *The class of very simple grammars is identifiable in the limit from positive data.*

#### 4.3. Time complexity analysis

We analyse the time complexity of  $IA$  in two respects: time for updating a conjecture and a bound on the number of implicit errors of prediction.

We analyse the time complexity of  $IA$  as follows:

(1) *Time for updating a conjecture.* Let  $N = \sum_{w_j \in R} \text{len}(w_j)$  and  $\ell$  be the maximum length of positive data in  $R$ . The time for updating a conjecture is obviously dominated by the time for the procedure  $\mathbf{Consistent}(R)$ .

In performing  $\mathbf{Consistent}(R)$ , determining  $f_n$  requires at most  $O(N)$  time. It takes at most  $O(N)$  time to construct  $\text{Lg}(R)$ . Solving  $\text{Lg}(R)$  requires at most  $O(|\Sigma|^3)$  time, because it is reduced to the computation of an inverse matrix with at most  $|\Sigma|$ -dimension. From (C1) of (Step 1) for the computation of  $f_p$  in Section 3.5, the value  $n_a$  is bounded by  $\ell$ , the number of all solution vectors of  $\text{Lg}(R)$  is bounded

by  $\ell^{|\Sigma|}$ . Hence, **while** loops are repeatedly performed at most  $\ell^{|\Sigma|}$  times. Each **while** loop requires at most  $O(N)$  time.

Thus, time for updating a conjecture is bounded by

$$O(|\Sigma|^3) + O(\ell^{|\Sigma|}) \times O(N) \leq O(\text{Max}\{N^{|\Sigma|+1}, |\Sigma|^3\}).$$

**Note 3.** The time complexity result above gives only an upper bound and, as mentioned in Note 2, an actual implementation could involve more efficient techniques.

(2) *A bound on the number of implicit errors of prediction.* Given a target very simple grammar  $G_*$ , let  $R$  be the current sample set of positive examples of  $L(G_*)$  over  $\Sigma$ . Further, let  $\text{Vec}(R) = \{\text{vec}(w) \mid w \in R\}$ , where  $\text{vec}(w) = (\#_{a_1}(w), \dots, \#_{a_m}(w))$ , and  $\Sigma' = \{a_1, \dots, a_m\}$  is an ordered set  $\text{alph}(R) - \Sigma_f(R)$  (see Section 3.5).

The number of implicit errors of prediction  $IA$  makes is analysed as follows.

**Lemma 18.** *The number of implicit errors of prediction  $IA$  makes is bounded by  $O(|\Sigma|)$ .*

**Proof.** Let  $I = (f_n, f_p)$  be the output of  $\text{Consistent}(R)$  such that  $I(\mathcal{G}_{0,\Sigma})$  is the current conjectured grammar from  $IA$ . Each time the conjecture  $I(\mathcal{G}_{0,\Sigma}) (= G_R)$  is refuted by the next new example  $w$ , we examine what changes happen on the current conjecture, that is, on  $I = (f_n, f_p)$ .

Each time  $I(\mathcal{G}_{0,\Sigma})$  fails to generate  $w$ ,  $IA$  produces a new conjecture  $I'(\mathcal{G}_{0,\Sigma})$  which can generate  $w$ . Let  $I' = (f'_n, f'_p)$  and  $R' = R \cup \{w\}$ , where  $I' = \text{Consistent}(R')$ . We analyse the following cases:

*Case (i).*  $f'_n$  is distinct from  $f_n$ : From the definition of **Consistent**( $R$ ) in  $IA$ , this change entails either at least one new nonterminal identification in  $\mathcal{G}_{0,\Sigma}$  or  $\text{alph}(R) \subset \text{alph}(R')$ , because of the monotonic incremental nature of nonterminal identification. Therefore, this case may happen in total at most  $2|\Sigma|$  times.

*Case (ii).*  $f'_p$  is distinct from  $f_p$ : Let  $\text{Vec}(R)$  and  $\text{Vec}(R')$  be the sets of vectors for  $R$  and  $R'$ , respectively. We note that since  $R \subset R'$ ,  $\text{Vec}(R) \subseteq \text{Vec}(R')$ . There are two possible cases: (1) a vector  $\text{vec}(w)$  is a linear combination of  $\text{Vec}(R)$  and (2)  $\text{Vec}(R') (= \text{Vec}(R) \cup \{\text{vec}(w)\})$  is linearly independent.

In case (1),  $\text{Sol}(\text{Lg}(R))$  remains unchanged. If  $\text{alph}(R) \subset \text{alph}(R')$ , then Case (i) must occur. Hence, suppose that  $\text{alph}(R) = \text{alph}(R')$ , and therefore, it must hold that for each  $a \in \text{alph}(R')$ ,  $\text{len}(f_p(x_a)) = \text{len}(f'_p(x_a))$ . However, since  $f'_p$  is distinct from  $f_p$ ,  $f_p(x_a) \neq f'_p(x_a)$  for some  $a \in \text{alph}(R')$ . Suppose that  $f_n = f'_n$ . Then, since the simulation of  $R$  by the set of candidate rules for  $I'$  must succeed using only the set of candidate rules for  $I$ , at that moment, for each  $a \in \text{alph}(R')$ , the value of  $f'_p(x_a)$  is completely determined through the simulation of  $R (= R' - \{w\})$ . Hence, without additional nonterminal identifications, we cannot have a situation where  $f_p(x_a) \neq f'_p(x_a)$  for some  $a \in \text{alph}(R')$  by simulating the derivation for  $w$ , because  $f_n = f'_n$  and  $\text{alph}(R) = \text{alph}(R')$ . Thus, as far as  $f_p \neq f'_p$ , we conclude that  $f_n \neq f'_n$ , that is, Case (i) must occur.

In case (2), for each  $R$  the number of vectors in  $\text{Vec}(R)$  linearly independent is bounded by the number of parameters to be solved. Hence, this can happen in total at most  $|\Sigma|$  times.

Each time the conjecture fails, and hence,  $I'$  is distinct from  $I$ , at least one of the above cases must occur. Therefore, the number of implicit errors of prediction is bounded by  $O(|\Sigma|)$ .  $\square$

One problem to be discussed here is how we should define the size of a grammar  $|G|$  for a very simple grammar  $G$  with the terminal alphabet  $\Sigma$ . In a usual setting, the size of a grammar may often be defined as the number of rules. However, for a fixed  $\Sigma$ , the number of rules of any very simple grammar is equal to the alphabet size, which implies that every very simple grammar with  $\Sigma$  has the same size of  $|\Sigma|$ .

We may define the size of a very simple grammar  $G = (V_N, \Sigma, P, S)$  as the length of its description for  $P$ . That is, let us define  $size(G)$  by  $\sum_{A \rightarrow \alpha \alpha \in P} (len(\alpha) + 3)$ . Then, it is obvious that  $|\Sigma| \leq size(G)$ .

Summing up the above observations together with Lemma 18, we have the following theorem.

**Theorem 19.** *The algorithm IA may be implemented to run in time polynomial in  $O(m)$  for updating a conjecture, and the number of implicit errors of prediction is bounded by  $O(n)$ , where  $n = size(G_*)$ , and  $m = \text{Max}\{N^{|\Sigma|+1}, |\Sigma|^3\}$ ,  $N = \sum_{w_j \in R} len(w_j)$ , and  $R$  is the set of positive data provided.*

In a similar manner discussed in [4,17], under the condition that only positive counterexamples are supplied in the identification process, the polynomial-time identifiability in the limit from positive data implies the polynomial-time identifiability via equivalence queries. That is, we have the following corollary.

**Corollary 20.** *The class of very simple grammars is polynomial-time identifiable via only equivalence queries, provided that only positive counterexamples are supplied in the identification process.*

## 5. An example run

**Example 5.** Consider  $G_* = (\{S, A, B, C\}, \{a, b, c, d, e, k\}, P, S)$  as the target grammar, where  $P = \{S \rightarrow aAS, S \rightarrow cBC, A \rightarrow b, B \rightarrow dBC, B \rightarrow e, C \rightarrow k\}$ . After all initial procedures, IA is ready to read the input example.

(1) Let  $w_1 = cd^2ek^3$  be the first example string from  $L(G_*)$ .  $R_1 = \{w_1\}$  and  $\text{alph}(R_1) = \{c, d, e, k\}$ . Then, since  $w_1$  is trivially not in  $L(G_{R_0})$ , the algorithm IA constructs the augmented grammar schema  $\mathcal{G}_{0, \Sigma}$  whose set of rule forms is  $\{X_c \rightarrow cx_c, X_d \rightarrow dx_d, X_e \rightarrow ex_e, X_k \rightarrow kx_k\}$ . Let  $V_N = \{X_c, X_d, X_e, X_k, S\}$ . Since a symbol  $k$  is final, the  $k$ -handle rule is  $X_k \rightarrow k$  and  $f_p(x_k) = \lambda(n_k = -1)$ . Further,  $f_n(X_c) = S$  is immediately obtained together with  $f_p(x_c) \neq \lambda(n_c \geq 0)$ . The length equation for  $w_1$  is

$$\begin{aligned} n_c + 2n_d + n_e + 3n_k &= -1 \quad \cdots (\ell.w_1), \text{ and hence,} \\ n_c + 2n_d + n_e &= 2 \quad \cdots (\ell.w_1)'. \end{aligned}$$



Thus, we have that  $\Sigma_s(R_1) = \{c\}$ ,  $\Sigma_f(R_1) = \{k\}$ , and  $\text{Sol}(\text{Lg}((R_1))) = \{n_c + 2n_d + n_e = 2\}$ , where  $\Sigma' = \{c, d, e\}$ . From the structure graph of  $R_1$ , it follows that  $f_n(X_c) = S$ ,  $f_n(X) = X (\forall X \in \{X_d, X_e, X_k\})$ . (See (a) of Fig. 4 for the structure graph  $T_{R_1}$ .)

Among possible solution vectors of  $\{n_c + 2n_d + n_e = 2\}$ , suppose we choose a solution  $\mathcal{X}_1 = (4, -1, 0) = (n_c, n_d, n_e)$ . That is, the set of candidate rules  $\text{CR}(\mathcal{X}_1)$  is

$$\{S \rightarrow cZ_{c,1}Z_{c,2}Z_{c,3}Z_{c,4}Z_{c,5}, X_d \rightarrow d, X_e \rightarrow eZ_{e,1}\}.$$

By simulating the derivation for  $w_1$  via these rules, we have that:  $Z_{c,1} = X_d$ ,  $Z_{c,2} = X_d$ ,  $Z_{c,3} = X_e$ ,  $Z_{c,4} = Z_{c,5} = Z_{e,1} = X_k$ . We see that  $\text{CR}(\mathcal{X}_1)$  is good and, as a result, a ground interpretation  $I_1 = (f_n, f_p)$  admissible to  $R_1$  is obtained, where

$$\begin{aligned} f_n(X_c) &= S, & f_n(X) &= X \text{ (} X \text{: otherwise),} \\ f_p(x_c) &= X_d^2 X_e X_k^2, & f_p(x_d) &= \lambda, \\ f_p(x_e) &= X_k, & f_p(x_k) &= \lambda. \end{aligned}$$

The first conjectured grammar  $G_{R_1} = I_1(\mathcal{G}_{0,\Sigma})$  is  $(\{S, X_d, X_e, X_k\}, \{c, d, e, k\}, P_1, S)$ , where  $P_1$  is

$$\begin{aligned} S &\rightarrow cX_d^2 X_e X_k^2, & X_d &\rightarrow d \\ X_e &\rightarrow eX_k, & X_k &\rightarrow k. \end{aligned}$$

(2) Suppose that  $w_2 = cd^4ek^5$  is given as the second example string. Then,  $R_2 = \{w_1, w_2\}$  and  $\text{alph}(R_2) = \text{alph}(R_1) = \{c, d, e, k\}$ . Since  $w_2$  is not in  $L(G_{R_1})$ ,  $\mathcal{G}_{0,\Sigma}$  is to be augmented but, in fact, unchanged.

First, by constructing the structure graph  $T_{R_2}$  pictured in (b) of Fig. 4, we observe that  $f_n(X_e) = X_d$ .

Then,  $IA$  computes a length equation for  $w_2$  and constructs  $\text{Lg}(R_2) = \{(\ell.w_1), (\ell.w_2)\}$ , where

$$\begin{aligned} n_c + 4n_d + n_e + 5n_k &= -1 \cdots (\ell.w_2), \text{ and hence,} \\ n_c + 4n_d + n_e &= 4 \cdots (\ell.w_2)'. \end{aligned}$$

To solve  $\text{Lg}(R'_2) = \{(\ell.w_1)', (\ell.w_2)'\}$ , we construct an associated matrix  $M_{R_2}$  and a matrix equation:

$$M_{R_2} = \begin{pmatrix} 1 & 2 & 1 \\ 1 & 4 & 1 \end{pmatrix} M_{R_2} \mathcal{X}_2^T = (2, 4)^T \\ \mathcal{X}_2 = (n_c, n_d, n_e).$$

The matrix computation is done as follows:

$$\left( \begin{array}{ccc|c} 1 & 2 & 1 & 2 \\ 1 & 4 & 1 & 4 \end{array} \right) \Rightarrow^* \left( \begin{array}{ccc|c} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{array} \right).$$

Thus,  $\text{Sol}(\text{Lg}(R'_2)) = \{n_d = 1\} \cup \{n_c + n_e = 0\}$ . From  $\text{Sol}(\text{Lg}(R'_2))$ , for example, we may choose a solution vector  $\mathcal{X}_2 = (0, 1, 0) = (n_c, n_d, n_e)$ . Hence, the set of candidate rules

$$\text{CR}(\mathcal{X}_2) = \{S \rightarrow cZ_{c,1}, X_d \rightarrow dZ_{d,1}Z_{d,2}, X_e \rightarrow eZ_{e,1}\}$$

is obtained. Simulating  $R_2$  via  $\text{CR}(\mathcal{X}_2)$  succeeds without any additional nonterminal identification, and produces the identification results that  $Z_{c,1}=X_d$ ,  $Z_{d,1}=X_d$ ,  $Z_{d,2}=X_k$ ,  $Z_{e,1}=X_k$ . From this, a ground interpretation  $I_2=(f_n, f_p)$  is obtained, where

$$\begin{aligned} f_n(X_c) &= S, & f_n(X_e) &= X_d, \\ f_n(X) &= X & (\forall X \in V_N - \{X_c, X_e\}), \\ f_p(x_c) &= X_d, & f_p(x_d) &= X_d X_k, \\ f_p(x_e) &= X_k, & f_p(x_k) &= \lambda. \end{aligned}$$

Thus, the second conjectured grammar  $G_{R_2}=I_2(\mathcal{G}_{0,\Sigma})$  is the following:  $(\{S, X_d, X_k\}, \{c, d, e, k\}, P_2, S)$ , where  $P_2$  is

$$\begin{aligned} S &\rightarrow cX_d, & X_d &\rightarrow dX_d X_k, \\ X_d &\rightarrow eX_k, & X_k &\rightarrow k. \end{aligned}$$

(3) Suppose that  $w_3=abababcek$  is given as the third example string. Then,  $R_3=\{w_1, w_2, w_3\}$  and  $\text{alph}(R_3)=\text{alph}(R_2)\cup\{a, b\}$ . Since  $w_3$  is not in  $L(G_{R_2})$ , in order to augment  $\mathcal{G}_{0,\Sigma}$ , two rule forms  $X_a \rightarrow ax_a$  and  $X_b \rightarrow bx_b$  are added to the set of rule forms of  $\mathcal{G}_{0,\Sigma}$ . Note that  $V_N=\{X_a, X_b, X_d, X_k, S\}$ ,  $\Sigma_s(R_3)=\{a, c\}$ . The structure graph  $T_{R_3}$  is shown in (c) of Fig. 4. Hence, we have that  $f_n(X_a)=S$  and  $f_p(x_a)\neq\lambda(n_a\geq 0)$ .

IA computes a length equation for  $w_3$  and constructs  $\text{Lg}(R'_3)=\{(\ell.w_1)', (\ell.w_2)', (\ell.w_3)'\}$ , where

$$\begin{aligned} 3n_a + 3n_b + n_c + n_e + n_k &= -1 \cdots (\ell.w_3), \text{ and hence,} \\ 3n_a + 3n_b + n_c + n_e &= 0 \cdots (\ell.w_3)'. \end{aligned}$$

Then, an associated matrix  $M_{R_3}$  and its matrix equation are:

$$M_{R_3} = \begin{pmatrix} 0 & 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & 4 & 1 \\ 3 & 3 & 1 & 0 & 1 \end{pmatrix} \begin{matrix} M_{R_3} \mathcal{X}_3^T = (2, 4, 0)^T \\ \mathcal{X}_3 = (n_a, n_b, n_c, n_d, n_e). \end{matrix}$$

The matrix computation is done as follows:

$$\left( \begin{array}{ccccc|c} 0 & 0 & 1 & 2 & 1 & 2 \\ 0 & 0 & 1 & 4 & 1 & 4 \\ 3 & 3 & 1 & 0 & 1 & 0 \end{array} \right) \Rightarrow^* \left( \begin{array}{ccccc|c} 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 \end{array} \right).$$

Thus,  $\text{Sol}(\text{Lg}(R'_3))=\{n_d=1\}\cup\{n_a+n_b=0, n_c+n_e=0\}$  is obtained from which, for example, we may choose a solution vector  $\mathcal{X}_3=(0,0,0,1,0)=(n_a, n_b, n_c, n_d, n_e)$ . Hence, the set of candidate rules  $\text{CR}(\mathcal{X}_3)$  is

$$\begin{aligned} S &\rightarrow aZ_{a,1}, & X_b &\rightarrow bZ_{b,1}, & S &\rightarrow cZ_{c,1}, \\ X_d &\rightarrow dZ_{d,1}Z_{d,2}, & X_d &\rightarrow eZ_{e,1}. \end{aligned}$$

Then, without any additional nonterminal identification, simulating  $R_3$  via  $\text{CR}(\mathcal{X}_3)$  leads to the identification results that  $Z_{a,1}=X_b$ ,  $Z_{b,1}=S$ ,  $Z_{c,1}=X_d$ ,  $Z_{d,1}=X_d$ ,  $Z_{d,2}=X_k$ ,

$Z_{e,1} = X_k$ . From this, a ground interpretation  $I_3 = (f_n, f_p)$  is obtained, where

$$\begin{array}{ll} f_n(X_a) = f_n(X_c) = S, & f_n(X_e) = X_d, \\ f_n(X_b) = X_b, & f_n(X_k) = X_k, \\ f_p(x_a) = X_b, & f_p(x_b) = S, \\ f_p(x_c) = X_d, & f_p(x_d) = X_d X_k, \\ f_p(x_e) = X_k, & f_p(x_k) = \lambda. \end{array}$$

Thus, the third conjectured grammar  $G_{R_3} = I_3(\mathcal{G}_{0,\Sigma})$  is the following one:  $(\{S, X_b, X_d, X_k\}, \{a, b, c, d, e, k\}, P_3, S)$ , where  $P_3$  consists of

$$\begin{array}{lll} S \rightarrow aX_b, & X_b \rightarrow bS, & S \rightarrow cX_d, \\ X_d \rightarrow dX_d X_k, & X_d \rightarrow eX_k, & X_k \rightarrow k. \end{array}$$

The conjecture  $G_{R_3}$  is equivalent (but not isomorphic) to  $G_*$ . Hence,  $G_{R_3}$  is always output as a conjecture for all input data afterwards.

## 6. Discussion

### 6.1. Related works

Mäkinen [13] discusses the problem of learning Szilard languages of linear grammars and gives a linear update-time algorithm for solving the problem, however he gives no consideration on the polynomial-time identifiability of the class in the sense discussed in this paper. From the definition, the class of Szilard languages of *linear* grammars is clearly a proper subclass of the class of very simple languages, and is also properly included in the class of zero-reversible languages [2]. (In fact, the class of very simple languages coincides with the left Szilard languages of *context-free* grammars [2].) Further, the class of very simple languages is incomparable to the class of zero-reversible languages. (See Fig. 1.) Thus, the main result in the present paper is in contrast with the fact that the class of regular languages is not polynomial-time identifiable in the limit using *DFAs* in the sense of Pitt [17].

Wakatsuki and Tomita [21] study the inclusion problem for the class of very simple grammars and show the problem is decidable by giving an efficient algorithm.

### 6.2. Conclusions

By adopting a slight modification of Pitt's definition, we have shown that the class of very simple grammars is identifiable in the limit from positive data, and presented an algorithm which identifies any very simple grammar in polynomial time in the new definition. One of the interesting features of the learning algorithm developed in this paper is the property of "reliability" in the sense that the algorithm *IA* has the ability of "self-diagnosing" which can detect input data inconsistent with any very simple grammar (see Corollary 11). Hence, besides its ability to identify every very simple language, the algorithm *IA* is additionally able to *refute* (in the sense of Mukouchi

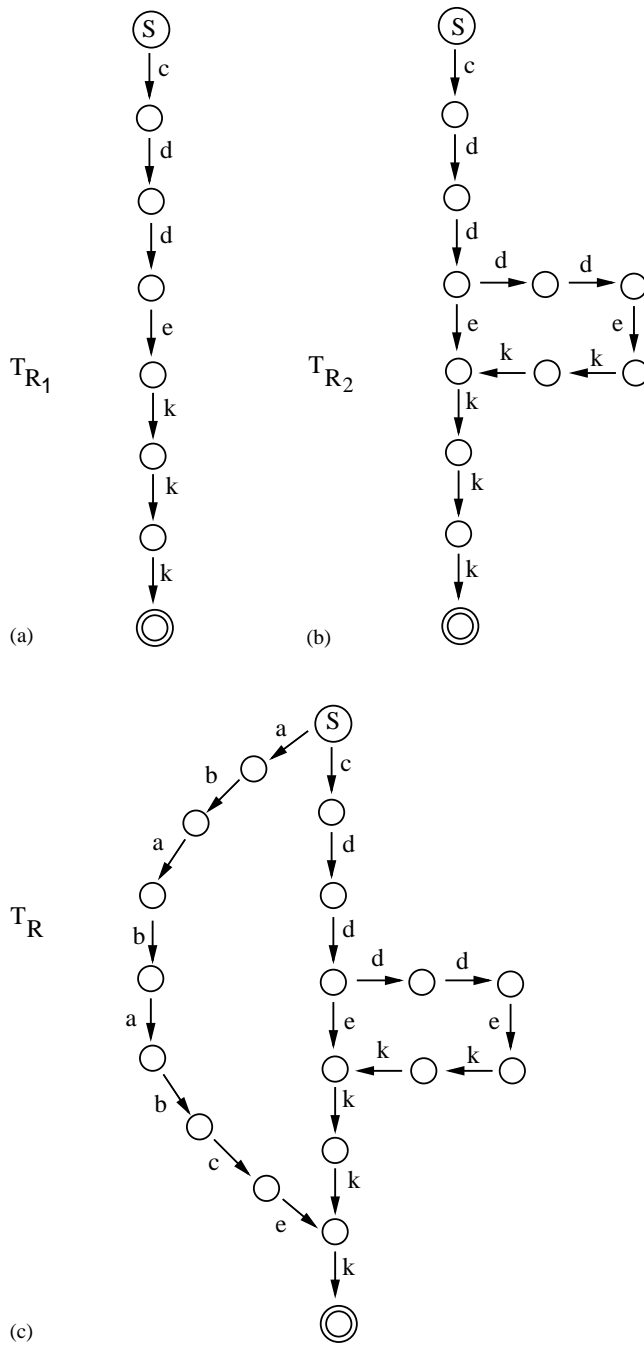


Fig. 4.

and Arikawa [16]) every positive presentation that does *not* describe a very simple language.

As we have shown in Section 3, the class of very simple grammars is moderately powerful in language generative capability, and it bears some of the essential properties of *CFGs* except for the nondeterministic features.

One of the results by Shinohara [19] shows that the class of context-sensitive grammars having a given fixed number of rules is identifiable in the limit from positive data, which implies the identifiability of the class of very simple grammars with a fixed size of terminal alphabet. In fact, for a fixed alphabet  $\Sigma$ , one may argue a simple and enumerative algorithm for identifying the class in polynomial time, because the number of rules in a grammar  $G$  is bounded by  $|\Sigma|$  and hence there are no more than  $O(n^{|\Sigma|})$  feasible hypotheses of what the grammar  $G$  could be, where  $n$  is the size of  $G$ . Compared to such a naive enumerative algorithm, our algorithm has some advantages. First, from the viewpoint of an upper-bound of the time complexity, the latter is only *potentially* exponential in  $|\Sigma|$  in the sense that we do not know at present whether or not the worst case really occurs, while the former must require exponential time if the target grammar is the last one up to equivalence in the enumeration process due to the “size” order of very simple grammars. Second, the number of updating conjectures is crucially important for a learning algorithm, and our algorithm requires at most  $O(|\Sigma|)$  times, while the other needs at the worst case  $O(n^{|\Sigma|})$  times. (Thus, these two algorithms make a great difference as  $|\Sigma|$  grows up.)

For future study, it seems useful to pursue applications of an inference algorithm for very simple grammars to other domains of research interests. In fact, making use of the fact that the left Szilard language of any *CFG* is very simple, Mäkinen [14] discusses a theoretical application of the inference algorithm for very simple grammars to the inference problem for *CFGs* from the structural data in which a Szilard word together with its corresponding terminal string is given as an example. In this problem setting, he shows the structural inference problem for *CFGs* is efficiently solvable. In relation to this topic, it would be a very interesting open problem to know whether or not there is an algorithm for identifying the class of very simple grammars with  $\Sigma$  in polynomial time in  $n$ ,  $N$  (total length of data provided) and  $|\Sigma|$ .

## Acknowledgements

The author would like to express his gratitude to E. Tomita and M. Wakatsuki who not only worked through an earlier draft of the paper but gave invaluable comments and suggestions. In particular the author is very much grateful to M. Kanazawa and R. Yoshinaka for their crucial comments which pointed out flaws in argument concerning the characteristic samples in the original version of this paper. Also, he would like to thank E. Mäkinen for many useful reference papers. Last but not least, the author is deeply indebted to referees for their many useful comments.

This work is supported in part by Grants-in-Aid for Scientific Research Nos. 03245104 and 04229105 from the Ministry of Education, Science, Sports and Culture, Japan.

## References

- [1] D. Angluin, Inductive inference of formal languages from positive data, *Inform. and Control* 45 (1980) 117–135.
- [2] D. Angluin, Inference of reversible languages, *J. ACM* 29 (1982) 741–765.
- [3] D. Angluin, Learning regular sets from queries and counterexamples, *Inform. and Comput.* 75 (1987) 87–106.
- [4] D. Angluin, Negative results for equivalence queries, *Mach. Learning* 5 (1990) 121–150.
- [5] P. Butzbach, Une famille de congruences de thue pour lesquelles le probleme de l'equivalence est decidable. application a l'equivalence des grammaires separees, in: M. Nivat (Ed.), *Automata, Languages and Programming*, North-Holland/American Elsevier, Amsterdam, 1973, pp. 3–12.
- [6] P. Garcia, E. Vidal, Inference of k-testable languages in the strict sense and application to syntactic pattern recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 12 (9) (1990) 920–925.
- [7] E.M. Gold, Language identification in the limit, *Inform. and Control* 10 (1967) 447–474.
- [8] M. Harrison, *Introduction to Formal Language Theory*, Addison-Wesley, Reading, MA, 1978.
- [9] J.E. Hopcroft, J.D. Ullman, *Introduction to Automata Theory, Languages, and Computation*, Addison-Wesley, Reading, MA, 1979.
- [10] A. Korenjak, J.E. Hopcroft, Simple deterministic languages, *Proc. 7th Annu. IEEE Conf. on Switching and Automata Theory*, 1966, pp. 36–46.
- [11] P. Laird, E. Gamble, *Analytical learning and term-rewriting systems*, Report RIA-90-06-17-7, Ames Research Center, NASA, 1990.
- [12] E. Mäkinen, On context-free derivations, Ser A 198, Act Universitatis Tamperensis, 1985.
- [13] E. Mäkinen, The grammatical inference problem for the szilard languages of linear grammars, *Inform. Process. Lett.* 36 (1990) 203–206.
- [14] E. Mäkinen, Remarks on the structural grammatical inference problem for context-free grammars, *Inform. Process. Lett.* 44 (1992) 125–127.
- [15] E. Moriya, The associate language and the derivation complexity of formal grammars, *Inform. and Control* 22 (1973) 139–162.
- [16] Y. Mukouchi, S. Arikawa, Towards a mathematical theory of machine discovery from facts, *Theoret. Comput. Sci.* 137 (1995) 53–84.
- [17] L. Pitt, Inductive inference, DFAs, and computational complexity, *Proc. 2nd Workshop on Analogical and Inductive Inference*, Lecture Notes in Artificial Intelligence, Vol. 397, Springer, Berlin, 1989, pp. 18–44.
- [18] M. Sato, K. Umayahara, Inductive inferability for formal languages from positive data, *IEICE Trans. Inform. Systems* E 75-D (4) (1992) 84–92.
- [19] T. Shinohara, Rich classes inferable from positive data: length bounded elementary formal systems, *Inform. Comput.* 108 (1994) 175–186.
- [20] N. Tanida, T. Yokomori, Polynomial-time identification of strictly regular languages in the limit, *IEICE Trans. Inform. Systems* E 75-D (1) (1992) 125–132.
- [21] M. Wakatsuki, E. Tomita, A fast algorithm for checking the inclusion for very simple deterministic pushdown automata, *IEICE Trans. Inform. Systems* E 76-D (10) (1993) 1224–1233.
- [22] K. Wright, Identification of unions of languages drawn from an identifiable class, *Proc. 2nd Workshop on Computational Learning Theory*, 1989, pp. 328–333, T. Motoki, T. Shinohara, K. Wright, The correct definition of finite elasticity: corrigendum to identification of unions, *Proc. 4th Workshop on Computational Learning Theory*, 1991, pp. 375.
- [23] T. Yokomori, On polynomial-time learnability in the limit of strictly deterministic automata, *Machine Learning* 19 (2) (1995) 153–179.
- [24] T. Yokomori, S. Kobayashi, Learning local languages and their application to DNA sequence analysis, *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (10) (1998) 1067–1079.