# PATTERN AMBIGUITIES FOR PURE CONTEXT–FREE GRAMMARS

Erkki Mäkinen and Ferucio Laurenţiu Ţiplea

# PATTERN AMBIGUITIES FOR
# PURE CONTEXT–FREE GRAMMARS

Erkki Mäkinen and Ferucio Laurenţiu Ţiplea

# Pattern Ambiguities for
# Pure Context–Free Grammars

**Erkki MÄKINEN**[a]  and  **Ferucio Laurenţiu ŢIPLEA**[b]

[a]Department of Computer Science
University of Tampere, P.O. Box 607
FIN-33101 Tampere, Finland
e-mail: `em@cs.uta.fi`

[b]Faculty of Informatics
"Al. I. Cuza" University of Iaşi
6600 Iaşi, Romania
e-mail: `fltiplea@infoiasi.ro`

## Abstract

We consider here some special forms of ambiguity for pure context-free grammars, namely *pattern avoiding ambiguity*, *pattern preserving ambiguity*, *pattern ambiguity*, and *grammar avoiding ambiguity*. The first two properties are undecidable for arbitrary pure context-free grammars, and in a particular but general enough case we can effectively decide whether or not a pure context-free grammar is pattern preserving ambiguous.

# 1   Introduction and Preliminaries

The notion of ambiguity is one of the key concepts of formal language theory. For results concerning ambiguity in standard Chomsky type context-free grammars, the reader is referred e.g. to [2]. In this paper we study some special forms of ambiguity for pure context-free grammars. This study is motivated by an encryption method proposed in [5]. We start with some terminology and notation (for details concerning formal language theory the reader is referred to [2], [9]).

A *pure context-free grammar* ($PCF$ grammar, for short) [6] is a 3-tuple $G = (V, V_0, P)$, where $V$ is a finite alphabet, the set of *axioms* $V_0$ is a finite subset of $V^+$, and $P$ is a finite set of productions of the form $a \longrightarrow v$, where $a \in V$ and $v \in V^*$. If $v = \lambda$ then the production $a \longrightarrow v$ is called a $\lambda$-*production*. The relation $\Longrightarrow$ (yields directly) and its reflexive and transitive closure $\overset{*}{\Longrightarrow}$ is defined as usual. $G$ is said to be *production-complete* [5] if for each symbol $a \in V$ there is at least one production with $a$ in its left hand side.

Let $G = (V, V_0, P)$ be a pure context-free grammar. We consider a finite set $Lab$ of *labels* and a bijective function $\rho$ from $P$ into $Lab$. The couple $\gamma = (G, \rho)$ is called a *labelled PCF grammar* ($lPCF$ grammar, for short). If a production $u \longrightarrow v$ is associated with the label $r$ by $\rho$, we will write $r : u \longrightarrow v$. The leftmost derivation

of $y$ from $x$ by using the sequence of productions $\alpha = r_1 \ldots r_n \in Lab^+$ is denoted by $x \overset{\alpha}{\Longrightarrow}_{\gamma,l} y$.

With a labelled grammar $\gamma$ we associate a morphism $g : Lab^* \longrightarrow \{0,1\}^*$ with the property $\mid g(r) \mid \leq 1$ for all $r \in Lab$. A couple $(\gamma, g)$ as above will be called an *arbitrarily labelled grammar*, (*alPCF* grammar, for short). The terminology is justified by the fact that $g \circ \rho$ associates to each production $p \in P$ an element of the set $\{0, 1, \lambda\}$, and thus $g \circ \rho$ can be thought as an arbitrary (not necessarily one-to-one) labelling of $G$, permiting $\lambda$ as a label. We have to remark that the labelling $\rho$ helps us to easily handle the productions, whereas the labelling $g$ is in fact the true labelling permiting different "actions" with the same "name" or without name (labelled by $\lambda$).

An example of an *alPCF* grammar $(\gamma, g)$ is presented in the next table. The last column lists the productions of the grammar, the second one describes the labeling $\rho$, whereas the first one gives the morphism $g$:

| $g$ | $\rho$ | P |
|-----|--------|---|
| $\lambda$ | $r_1$ | $S \longrightarrow ab$ |
| $0$ | $r_2$ | $a \longrightarrow ab$ |
| $0$ | $r_3$ | $a \longrightarrow cb$ |
| $0$ | $r_4$ | $c \longrightarrow cb$ |
| $0$ | $r_5$ | $c \longrightarrow ab$ |
| $1$ | $r_6$ | $b \longrightarrow aa$ |

A compact representation for such grammars is obtained by writing both labels in the front of each production. For example, the first two rules of the above grammar can be written as $\lambda : r_1 : S \longrightarrow ab$, $0 : r_2 : a \longrightarrow ab$.

In the next we shall consider only $lPCF$ grammars with an axiom denoted by $S$. For such a grammar $\gamma$ define the *left Szilard language* of $\gamma$ [7],[8], denoted $Szl(\gamma)$, by $Szl(\gamma) = \{\alpha \in Lab^+ | S \overset{\alpha}{\Longrightarrow}_{\gamma,l} y\}$; the strings $\alpha \in Szl(\gamma)$ are called *control words of $\gamma$*.

## 2　On An Encryption Method

In [1] an encryption method based on control words associated to leftmost derivations in context-free grammars has been proposed. We recall this method in terms of $PCF$ grammars [5].

An $alPCF$ grammar $(\gamma, g)$ defines an encryption method as follows. Assume the plaintext space being the phrases in English. Such a phrase $u$ is encoded, by a one-to-one coding, into a string $x \in \{0, 1\}^+$. From this string we compute a new string $\alpha \in Lab^+$ such that $g(\alpha) = x$ and $\alpha$ is a control word of $\gamma$ (suppose such strings exist[1] and we arbitrarily choose one of them). A possible encryption of $u$ is then the sentential form obtained in $\gamma$ rewriting the axiom with the control word $\alpha$.

---

[1]an $alPCF$ grammar $(\gamma, g)$ satisfying the property $(\forall x \in \{0,1\}^+)(\exists \alpha \in Szl(\gamma))(g(\alpha) = x)$ is called $(0,1)$-*total* [1], [5], [4]

Concerning this encryption method many problems arise [1], [5], [4]. Let us suppose we are legal users of a cryptosystem like that described above:

$$\underbrace{u}_{plaintext} \overset{encode}{\longrightarrow} x \in \{0,1\}^{+} \overset{g^{-1}}{\longrightarrow} \alpha \in Szl(\gamma) \overset{\gamma}{\longrightarrow} S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \underbrace{y}_{cryptotext}$$

(the strings $x \in \{0,1\}^{+}$ will be called *patterns*, and if $g(\alpha) = x$ we say that $\alpha$ *extends* $x$ with respect to $g$).

Starting from $x$ the sender has to look for a control word $\alpha$ extending $x$, and to rewrite the axiom $S$ of $\gamma$ with respect to $\alpha$; the result will be the cryptotext. In general, there could exist more than one control word extending $x$, and therefore the sender has to choose one of them. Then we have to take into consideration the next two cases:

(S1) (easy encryption)

for any two distinct control words $\alpha$ and $\beta$ extending the same pattern ($g(\alpha) = g(\beta)$) we have $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S^2$. That is, the cryptotext associated to $x$ does not depend on the control word extending the pattern $g(\alpha) = g(\beta)$. Clearly, the sender may choose in the easiest way a control word and to encrypt with respect to it; the problem is that the cryptotext could be less safe;

(S2) (supplementary encryption key)

for any two distinct control words $\alpha$ and $\beta$ extending the same pattern ($g(\alpha) = g(\beta)$) we have $\neg(S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$. That is, the cryptotext associated to $x$ heavily depends on the control word extending the pattern $g(\alpha) = g(\beta)$. Such a property permits to the sender the use of an encryption key in order to choose a control word; this could be time consuming but in this way the cryptosystem is strenghtened.

¿From the receiver's point of view we have to take into considerations another two cases:

(R1) (easy decryption)

for any two distinct parsings $\alpha$ and $\beta$ of the same cryptotext we have $g(\alpha) = g(\beta)$. Therefore, the receiver may parse the cryptotext in the easiest way, but the cryptosystem could be less safe;

(R2) (supplementary decryption key)

for any two distinct parsings $\alpha$ and $\beta$ of the same cryptotext we have $g(\alpha) \neq g(\beta)$. Therefore, the receiver needs supplementary information in order to decrypt (a supplementary encryption/decryption key). In fact, this case is similar to the case (S2).

As we will see in the next section, these considerations lead to different forms of ambiguity for $alPCF$'s.

---

[2] "$S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S$" denotes "$\exists y : S \overset{\alpha}{\Longrightarrow}_{\gamma,l} y \ \wedge \ S \overset{\beta}{\Longrightarrow}_{\gamma,l} y$"; sometimes it is necessary to fix $y$, and then we write $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} y \overset{\beta}{\Longleftarrow}_{\gamma,l} S$

# 3 Pattern Ambiguities

An $lPCF$ $\gamma$ is *ambiguous* if $(\exists \alpha, \beta \in Szl(\gamma))(\alpha \neq \beta \ \wedge \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$. This notion of ambiguity may be translated to $alPCF$ $(\gamma, l)$ by means of $\gamma$. But, in this case, the ambiguity can be devided in two classes with respect to $g$:

1. $\exists \alpha, \beta \in Szl(\gamma)$ such that $\alpha \neq \beta$ and $g(\alpha) = g(\beta)$ and $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S$. We will call this one *pattern preserving ambiguity* (pp-ambiguity or ppa, for short), because the computations $\alpha$ and $\beta$ leading to ambiguity have in common the pattern $g(\alpha) = g(\beta)$;

2. $\exists \alpha, \beta \in Szl(\gamma)$ such that $\alpha \neq \beta$ and $g(\alpha) \neq g(\beta)$ and $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S$. We will call this one *pattern avoiding ambiguity* (pa-ambiguity or paa, for short), because the computations $\alpha$ and $\beta$ leading to ambiguity have no pattern in common.

It is clear that the negations of these two kinds of ambiguities correspond to the cases (S2) (or (R2)) and (R1), respectively.

The above two ambiguities are referring to the grammar but we may consider ambiguity with respect to the pattern as well. Therefore, we say that $(\gamma, g)$ is *pattern ambiguous* if there exist $x \in \{0, 1\}^{+}$ and $\alpha, \beta \in Szl(\gamma)$ such that $\alpha \neq \beta$ and $g(\alpha) = g(\beta)$. That is, there is an $x$ which can be extended to two different control words. In other words, $(\gamma, g)$ is pattern ambiguous if $(\exists \alpha, \beta \in Szl(\gamma))(\alpha \neq \beta \ \wedge \ g(\alpha) = g(\beta))$. This ambiguity can be divided in two kinds of ambiguities with respect two the grammar $\gamma$:

1. $\exists \alpha, \beta \in Szl(\gamma)$ such that $\alpha \neq \beta$ and $g(\alpha) = g(\beta)$ and $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S$. In fact this one is the pattern preserving ambiguity;

2. $\exists \alpha, \beta \in Szl(\gamma)$ such that $\alpha \neq \beta$ and $g(\alpha) = g(\beta)$ and $\neg(S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$. This pattern ambiguity avoids grammar ambiguity and we will call it *grammar avoiding ambiguity* (ga-ambiguity or gaa, for short).

It is clear that the negation of ga-ambiguity corresponds to the case (S1). Let us conclude our discussion:

- **ga-ambiguity — easy encryption**

  (gaa) $\quad (\exists \alpha, \beta \in Szl(\gamma))(\alpha \neq \beta \ \wedge \ g(\alpha) = g(\beta) \ \wedge \ \neg(S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S))$

  (non-gaa) $\quad (\forall \alpha, \beta \in Szl(\gamma))(\alpha = \beta \ \vee \ g(\alpha) \neq g(\beta) \ \vee \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$, or

  $\quad (\forall \alpha, \beta \in Szl(\gamma))((\alpha \neq \beta \ \wedge \ g(\alpha) = g(\beta)) \ \Rightarrow \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$;

- **pa-ambiguity — easy decryption**

  (paa) $\quad (\exists \alpha, \beta \in Szl(\gamma))(\alpha \neq \beta \ \wedge \ g(\alpha) \neq g(\beta) \ \wedge \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$

  (non-paa) $\quad (\forall \alpha, \beta \in Szl(\gamma))(\alpha = \beta \ \vee \ g(\alpha) = g(\beta) \ \vee \ \neg(S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S))$, or

  $\quad (\forall \alpha, \beta \in Szl(\gamma))((\alpha \neq \beta \ \wedge \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S) \ \Rightarrow \ g(\alpha) = g(\beta))$;

- **pp-ambiguity — supplementary encryption/decryption key**

(ppa)      $(\exists \alpha, \beta \in Szl(\gamma))(\alpha \neq \beta \ \wedge \ g(\alpha) = g(\beta) \ \wedge \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S)$

(non-ppa)   $(\forall \alpha, \beta \in Szl(\gamma))(\alpha = \beta \ \vee \ g(\alpha) \neq g(\beta) \ \vee \ \neg(S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S))$, or

         $(\forall \alpha, \beta \in Szl(\gamma))((\alpha \neq \beta \ \wedge \ g(\alpha) = g(\beta)) \ \Rightarrow \ \neg(S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S))$, or

         $(\forall \alpha, \beta \in Szl(\gamma))((\alpha \neq \beta \ \wedge \ S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S) \ \Rightarrow \ g(\alpha) \neq g(\beta))$.

**Example 3.1**    *(1) Let $(\gamma, g_1)$ be the alPCF given by:*

$$0 : r_1 : S \longrightarrow ac, \ 0 : r_2 : S \longrightarrow bc, \ \lambda : r_3 : a \longrightarrow aa,$$
$$\lambda : r_4 : b \longrightarrow aa, \ \lambda : r_5 : c \longrightarrow cc.$$

*It is easy to see that this grammar is both pp-ambiguous (take $\alpha = r_1 r_3$ and $\beta = r_2 r_4$) and ga-ambiguous (take $\alpha = r_1$ and $\beta = r_2$) but not pa-ambiguous.*

*(2) Let $(\gamma, g_2)$ be the alPCF given by:*

$$0 : r_1 : S \longrightarrow ac, \ 1 : r_2 : S \longrightarrow bc, \ \lambda : r_3 : a \longrightarrow aa,$$
$$\lambda : r_4 : b \longrightarrow aa, \ \lambda : r_5 : c \longrightarrow cc.$$

*It is easy to see that this grammar is pa-ambiguous (take $\alpha$ and $\beta$ as in (1)) but neither pp-ambiguous nor ga-ambiguous.*

*(3) Let $(\gamma', g_3)$ be the alPCF given by:*

$$0 : r_1 : S \longrightarrow ac, \ 0 : r_2 : S \longrightarrow bc.$$

*It is easy to see that this grammar is ga-ambiguous (take $\alpha = r_1$ and $\beta = r_2$) but not pp-ambiguous. Consider the grammar $(\gamma'', g_4)$ given by*

$$0 : r_1 : S \longrightarrow ac, \ \lambda : r_2 : S \longrightarrow bc, \ 0 : r_3 : b \longrightarrow a.$$

*This grammar is pp-ambiguous (take $\alpha = r_1$ and $\beta = r_2 r_3$) but not ga-ambiguous.*

**Remark 3.1**    *(1) It follows directly from the definitions that pp- or pa-ambiguity of $(\gamma, g)$ implies the ambiguity of $\gamma$. The converse does not hold true as Example 3.1 shows us: the grammar $\gamma$ is ambiguous but $(\gamma, g_1)$ is not pa-ambiguous, and $(\gamma, g_2)$ is not pp-ambiguous.*

*(2) Similarly, pp- or ga-ambiguity implies pattern ambiguity, but the converse does not hold true: $(\gamma, g_2)$ in Example 3.1 is pattern ambiguous (take $x = 0$, $\alpha = r_1$, and $\beta = r_1 r_3$) but neither pp- nor ga-ambiguous.*

*(3) From Example 3.1 we can see that the ga-, pp- and pa-ambiguities are not implying each others.*

We will show that the ambiguity, pp- and pa-ambiguity are in general undecidable.

**Lemma 3.1** *It is undecidable whether an lPCF grammar is ambiguous.*

**Proof** The proof is based on the recursive unsolvability of the Post Correspondence Problem and it follows the same line as for context-free grammars (see for example [3]). $\square$

**Lemma 3.2** *pp-ambiguity is undecidable for alPCF grammars.*

**Proof** We will reduce the ambiguity problem for *lPCF* grammars to the pp-ambiguity problem for *alPCF* grammars.

Let $\gamma$ be an *lPCF* grammar. Consider $(\gamma', g)$, where $\gamma'$ is obtained from $\gamma$ by adding a new symbol $S'$ and a new production $0 : p : S' \longrightarrow S$; for the other productions $r$ of $\gamma'$, we will set $g(r) = \lambda$. We show that $\gamma$ is ambiguous iff $(\gamma', g)$ is pp-ambiguous.

If $\alpha$ and $\beta$ are two distinct control words of $\gamma$ such that $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} y \overset{\beta}{\Longleftarrow}_{\gamma,l} S$ then $\alpha' = p\alpha$ and $\beta' = p\beta$ are two distinct control words of $(\gamma', g)$ with $g(\alpha') = 0 = g(\beta')$ and $S' \overset{\alpha'}{\Longrightarrow}_{\gamma',l} y \overset{\beta'}{\Longleftarrow}_{\gamma',l} S'$. Therefore, the ambiguity of $\gamma$ leads to the pp-ambiguity of $(\gamma', g)$.

The converse can be easily obtained observing that any non-empty control words of $(\gamma', g)$ is of the form $p\alpha$, where $\alpha$ is a control word of $\gamma$. $\square$

Before showing that the pa-ambiguity is undecidable we consider the concept of *2-branched grammar*.

**Definition 3.1** *An lPCF grammar $\gamma$ is said to be 2-branched if, for any symbol $a$ in $\gamma$, there exist at most two productions with $a$ in the left hand side.*

**Lemma 3.3** *pa-ambiguity is undecidable for alPCF grammars.*

**Proof** We will reduce the ambiguity problem for *lPCF* grammars to the pa-ambiguity problem for *alPCF* grammars.

Let $\gamma$ be an *lPCF* grammar. Transform first $\gamma$ into a new 2-branched *lPCF* grammar $\gamma'$ as follows: assume $a \longrightarrow u_1, \ldots, a \longrightarrow u_k$ are all the productions with $a$ in left hand side and $k > 2$. Consider the new symbols $a_1, \ldots, a_{k-2}$ and replace the above productions by the new productions $a \longrightarrow u_1, a \longrightarrow a_1, a_1 \longrightarrow u_2, \ldots, a_{k-3} \longrightarrow u_{k-2}, a_{k-3} \longrightarrow a_{k-2}, a_{k-2} \longrightarrow u_{k-1}, a_{k-2} \longrightarrow u_k$. By iterating this procedure we finally get a 2-branched *lPCF* $\gamma'$.

**Claim** $\gamma$ is ambiguous iff $\gamma'$ is ambiguous.

**Proof of Claim** Assume first $\gamma$ is ambiguous. That is, there exist two distinct control words $\alpha$ and $\beta$ such that $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} y$ and $S \overset{\beta}{\Longrightarrow}_{\gamma,l} y$ for some $y$. If these control words use rules of the form $a \longrightarrow u_i$ with $i > 1$ (as above) then they will be replaced by the sequence of rules $a \longrightarrow a_1, a_1 \longrightarrow a_2, \ldots, a_{i-1} \longrightarrow u_i$, for $i < k$, or by the sequence $a \longrightarrow a_1, a_1 \longrightarrow a_2, \ldots, a_{i-2} \longrightarrow u_k$, for $i = k$. Finally, we will

get two distinct control words of $\gamma'$, $\alpha'$ and $\beta'$, such that $S \overset{\alpha'}{\Longrightarrow}_{\gamma',l} y$ and $S \overset{\beta'}{\Longrightarrow}_{\gamma',l} y$. Hence, $\gamma'$ is ambiguous.

Conversely, let us assume that $\gamma'$ is ambiguous, and let $\alpha$ and $\beta$ be two distinct control words producing the same word $y$ from the axiom $S$. First, we may assume that both $\alpha$ and $\beta$ do not end by rules $x \longrightarrow z$, where $z$ is a new symbol. Indeed, if we suppose $\alpha = \alpha' r_1$, $\beta = \beta' r_2$, $r_1 : x_1 \longrightarrow z_1$, $r_2 : x_2 \longrightarrow z_2$, and $z_1$ and $z_2$ are new symbols, then we may write

$$S \overset{\alpha'}{\Longrightarrow}_{\gamma',l} y_1 = y_1^1 x_1 y_1^2 \Longrightarrow_{\gamma',l} y_1^1 z_1 y_1^2 = y$$
$$S \overset{\beta'}{\Longrightarrow}_{\gamma',l} y_2 = y_2^1 x_2 y_2^2 \Longrightarrow_{\gamma',l} y_2^1 z_2 y_2^2 = y.$$

Both $x_1$ and $x_2$ are the leftmost symbols that are rewritten and therefore $y_1^1 = y_2^1$ and $x_1 = x_2$. Since both derivations produce the same string, $y$, it follows that $z_1 = z_2$, $y_2^1 = y_2^2$, and so $r_1 = r_2$ and $\alpha' \neq \beta'$. Since $z_1$ is a new symbol, there must exist a rule having $z_1$ in its left hand side. As a result, we can add a new rule with $z_1 = z_2$ in the left hand side, to the end of $\alpha$ and $\beta$. Continuing this process as long as it is necessary we finally get two new control words with the desired property.

If $a$ is the leftmost symbol in a string $w = w_1 a w_2$ produced by $\gamma'$, and the rule $a \longrightarrow a'$ is used, where $a'$ is a new symbol, then the leftmost symbol to be rewritten in $w_1 a' w_2$ is $a'$. Consequently, if $\alpha$ and $\beta$ contain rules with new symbols in their right hand sides then these rules come in sequences of the form:

$$a \longrightarrow a_1, a_1 \longrightarrow a_2, \ldots, a_i \longrightarrow u_j,$$

where $i = j - 1$ or $i = j - 2$, $a_1, \ldots, a_i$ are new symbols, and $a \longrightarrow u_j$ is a rule in $\gamma$. This fact permits us to "pack" such sequences of rules into one rule, namely $a \longrightarrow u_j$. Finaly we will get two control words of $\gamma$, $\alpha'$ and $\beta'$, such that $S \overset{\alpha'}{\Longrightarrow}_{\gamma,l} y$ and $S \overset{\beta'}{\Longrightarrow}_{\gamma,l} y$.

The proof of the Claim is complete if we can show that $\alpha' \neq \beta'$. We know that $\alpha = \beta$. So, we can write $\alpha = \theta r_1 \alpha_1$ and $\beta = \theta r_2 \beta_1$, where $r_1 \neq r_2$. As $r_1$ and $r_2$ rewrite the leftmost symbol ($\theta$ being a common prefix of $\alpha$ and $\beta$) it follows that $r_1$ and $r_2$ have the same left hand side, say $x$. The above remarks lead to the fact that $\alpha$ and $\beta$ contain two subsequences $w_1 = \theta_2 r_1 \alpha_1^1$ and $w_2 = \theta_4 r_2 \beta_1^1$ such that:

- $\theta = \theta_1 \theta_2 = \theta_3 \theta_4$ and $\alpha_1 = \alpha_1^1 \alpha_1^2$ and $\beta_1 = \beta_1^1 \beta_1^2$, for some $\theta_1, \theta_3, \alpha_1^2, \beta_1^2$;

- these two sequences have the form:

$$a \longrightarrow a_1, \ldots, r_1, \ldots, a_i \longrightarrow u_j,$$

and respectively

$$a \longrightarrow a_1, \ldots, r_2, \ldots, a_s \longrightarrow u_t,$$

where $a \longrightarrow u_j$ and $a \longrightarrow u_t$ are rules in $\gamma$, and $u_j \neq u_t$ ($r_1$ or $r_2$ may be the first rules of these sequences).

7

The "packing operation" (described above) applied to $\alpha$ and $\beta$ will replace these two subsequences by the rules $a \longrightarrow u_j$ and $a \longrightarrow u_t$, respectively. Consequently, $\alpha'$ and $\beta'$ will be distinct control words.

The Claim is proved.

Consider now the $alPCF$ $(\gamma', g)$, where $\gamma'$ is obtained from $\gamma$ as described above and $g$ is given by the following setting: for any $a$, if there is at most one production with $a$ in the left hand side then label it by 0; otherwise (there are at most two productions with $a$ in the left hand side) label one of them by 0 and the other one by 1. We show that $\gamma'$ is ambiguous iff $(\gamma', g)$ is pa-ambiguous.

If $\alpha$ and $\beta$ are two distinct control words of $\gamma'$ such that $S \stackrel{\alpha}{\Longrightarrow}_{\gamma',l} y \stackrel{\beta}{\Longleftarrow}_{\gamma',l} S$ then $g(\alpha) \neq g(\beta)$ because there must exist a branching symbol for these control words (note that the derivation is leftmost). Therefore, the ambiguity of $\gamma'$ leads to the pa-ambiguity of $(\gamma', g)$.

The converse follows from definitions. $\square$

# 4 A Sufficient Criterion for the Decidability of pp-ambiguity

As mentioned in the previous sections the non-ppa property would be useful both for sender and receiver in order to strenghten the cryptosystem. Unfortunately, we have seen that this property is in general undecidable. In what follows we shall derive a sufficient criterion for the decidability of pp-ambiguity. The PCF grammars that we will consider have no $\lambda$-productions or renamings (that is, productions of the form $a \longrightarrow b$ with $a, b \in V$), and they are production-complete.

**Definition 4.1** *Let $\gamma = (G, \rho)$ be an lPCF grammar, $au \in V^+$ ($a \in V$), $b \in V$ and $\alpha \in Lab^+$. We say that $\alpha$ minimally covers $au$ with respect to $b$ if one of the next two properties is satisfied:*

1. *if $u = \lambda$ then $\alpha \in Lab$ and $b \stackrel{\alpha}{\Longrightarrow}_{\gamma,l} cv$, for some $c \in V$ and $v \in V^*$;*

2. *if $u \neq \lambda$ then $\alpha = r_1 \cdots r_n$ ($n \geq 1$), $b \stackrel{r_1 \cdots r_{n-1}}{\Longrightarrow}_{\gamma,l} cv_1 \stackrel{r_n}{\Longrightarrow}_{\gamma,l} dv_2$, $u$ is a suffix of $v_2$ but not of $v_1$ (for $n = 1$ we have $cv_1 = b\lambda$).*

Hence, we may informally say that $\alpha$ minimally covers $au$ with respect to $b$ if $\alpha$ is "a minimal sequence" such that the string yielded from $b$ by $\alpha$ "covers" the string $au$, excepting the first letter.

**Definition 4.2** *Let $\gamma$ be an lPCF grammar and $\alpha, \beta \in Szl(\gamma)$. We say that $\alpha$ and $\beta$ are compatible control words if $\alpha$ minimally covers $bv$ with respect to $S$ or $\beta$ minimally covers $au$ with respect to $S$, where $S \stackrel{\alpha}{\Longrightarrow}_{\gamma,l} au$ and $S \stackrel{\beta}{\Longrightarrow}_{\gamma,l} bv$.*

For an alphabet $V$ denote by $\leq_{pref}$ the *prefix (partial) order* on $V$, that is $u \leq_{pref} v$ iff $u$ is a prefix of $v$. Let $u$ and $v$ words over $V$. We say that $u$ and $v$ are *comparable* (with respect to the prefix order) if $u \leq_{pref} v$ or $v \leq_{pref} u$. If $u$ and $v$ are comparable

we say that the pair $(u', v')$ is obtained from $(u, v)$ by *left maximal cancellation* if $u'$ and $v'$ are respectively obtained from $u$ and $v$ by deleting the common maximal prefix of them.

Let $T \subseteq V$. By $\#(T, \alpha)$ we denote the number of occurrences of letters $a \in T$ in $\alpha$. Then the *balance* of $\alpha$ and $\beta$ over $T$ is defined by:

$$B_T(\alpha, \beta) = |\#(T, \alpha) - \#(T, \beta)|.$$

**Definition 4.3** *An alPCF grammar $(\gamma, g)$ is balance $k$-bounded, $k \geq 1$, if the next property is satisfied: for all $\alpha, \beta \in Szl(\gamma)$ such that $g(\alpha) = g(\beta)$ and $S \overset{\alpha}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta}{\Longleftarrow}_{\gamma,l} S$, and for all $\alpha' \leq_{pref} \alpha$ and $\beta' \leq_{pref} \beta$ such that $\alpha'$ and $\beta'$ are compatible, it follows that $B_{\{0,1\}}(\alpha', \beta') \leq k$.*

**Theorem 4.1** *It is effectively decidable whether or not a given balance $k$-bounded alPCF grammar $(\gamma, g)$ is pattern preserving ambiguous.*

**Proof**    Let $(\gamma, g)$ be a balance k-bounded *alPCF* grammar. We shall generate a graph $\mathcal{G} = (X, E)$ whose nodes are of the form

$$(au, b, x_1, x_2) \quad \text{or} \quad (a, bu, x_1, x_2)$$

where

- $a, b \in V$;

- $u \in V^*$ such that $|u| \leq max\{|w| \,|\, \exists a \longrightarrow w \in P\}$;

- $x_1, x_2 \in \{0, 1\}^*$.

The arcs of $\mathcal{G}$ will be labelled by couples $(\alpha, \beta)$ where $\alpha, \beta \in Lab^*$ and either $\alpha = \lambda$ or $\beta = \lambda$. The next procedure generates this graph:

1. generate the node $(S, S, \lambda, \lambda)$. This node is called the *initial node* of $\mathcal{G}$ and it is marked as an *open node*;

2. while there are open nodes left choose one of them, say $\delta$, close it and perform one of the next three steps (if both the steps 2.1 and 2.2 can be performed with respect to $\delta$ then choose 2.1):

   2.1 if $\delta = (au, b, x_1, x_2)$ then

        - for any minimal covering sequence $\beta \in Lab^+$ of $au$ with respect to $b$,

   $$b \overset{\beta}{\Longrightarrow}_{\gamma,l} cvu,$$

   generate (provided that it does not yet exist) the node $(a, cv, y_1, y_2)$, where $y_1, y_2$ are obtained from $x_1, x_2 g(\beta)$ by left maximal cancellation. If $y_1 \neq \lambda$ and $y_2 \neq \lambda$, or $B_{\{0,1\}}(y_1, y_2) > k$ then mark the node $(a, cv, y_1, y_2)$ as a *final* node; otherwise, mark it as an *opened* node;

9

- the arc $(\delta, (a, cv, y_1, y_2))$ will be labelled by $(\lambda, \beta)$;

2.2 if $\delta = (a, bv, x_1, x_2)$ then

- for any minimal covering sequence $\alpha \in Lab^+$ of $bv$ with respect to $a$,

$$b \overset{\alpha}{\Longrightarrow}_{\gamma,l} cuv,$$

generate (provided that it does not yet exist) the node $(cu, b, y_1, y_2)$, where $y_1, y_2$ are obtained from $x_1 g(\alpha), x_2$ by left maximal cancellation. If $y_1 \neq \lambda$ and $y_2 \neq \lambda$, or $B_{\{0,1\}}(y_1, y_2) > k$ then mark the node $(cu, b, y_1, y_2)$ as a *final* node; otherwise, mark it as an *opened* node;

- the arc $(\delta, (cu, b, y_1, y_2))$ will be labelled by $(\alpha, \lambda)$;

2.3 if no node can be generated from $\delta$ as described in 2.1 or 2.2 then mark $\delta$ as a *final* node;

The graph $\mathcal{G}$ is finite and can be effectively constructed. Indeed, as the grammar is production-complete and it has no $\lambda$-productions or renamings, it follows that there are finitely many minimal covering sequences of a string $au$ with respect to $b$, and these sequences can be effectively computed. Moreover, the use of such sequences lead to the fact that for any node $(u, v, x_1, x_2)$ the next two properties hold:

(1) $|u|, |v| \leq max\{|w| | \exists a \longrightarrow w \in P\}$;

(2) $|x_1|, |x_2| \leq k + max\{|w| | \exists a \longrightarrow w \in P\}$
(each rule of the grammar add at least one new symbol, and therefore every minimal covering sequence of a string $au'$ with respect to $b$ has the length at most $|u'|$. Now, the statement follows from (1)).

Let $\delta_1, \ldots, \delta_n = (u, v, x_1, x_2)$ $(n \geq 2)$ be a path in $\mathcal{G}$ such that $\delta_1$ is the initial node. Consider the sequence of arc labels determined by this path,

$$(\alpha_1, \beta_1) \cdots (\alpha_{n-1}, \beta_{n-1}),$$

where $\beta_i = \lambda$ for $i$ odd and $\alpha_j = \lambda$ for $j$ even, and let $\alpha = \alpha_1 \cdots \alpha_{n-1}, \beta = \beta_1 \cdots \beta_{n-1}$. The next properties hold true:

(3) $\alpha$ and $\beta$ are compatible control words;

(4) $B_{\{0,1\}}(\alpha, \beta) = B_{\{0,1\}}(x_1, x_2)$;

(5) if $\delta_n$ is a *final* node then there is no $\alpha', \beta' \in Szl(\gamma)$ such that $\alpha \leq_{pref} \alpha'$, $\beta \leq_{pref} \beta'$, $g(\alpha') = g(\beta')$, $S \overset{\alpha'}{\Longrightarrow}_{\gamma,l} \cdot \overset{\beta'}{\Longleftarrow}_{\gamma,l} S$. Indeed, if this is not the case then:

a) if $x_1 \neq \lambda$ and $x_2 \neq \lambda$ then $g(\alpha') \neq g(\beta')$; a contradiction;

b) if $B_{\{0,1\}}(x_1, x_2) > k$ then $(\gamma, g)$ would not have the balance $k$-bounded; a contradiction;

(6) if $u = v \in V$ and $x_1 = x_2 = \lambda$ and $\alpha \neq \beta$ then $(\gamma, g)$ is pattern preserving ambiguous; the converse holds also true (that is, if $(\gamma, g)$ is pp-ambiguous then there is a node $\delta$ satisfying (6)).

Finally, to decide whether $(\gamma, g)$ is pattern preserving ambiguous or not we have to generate the graph $\mathcal{G}$ and to verify whether (4) holds true or not. The theorem is proved. $\square$

One can easily notice that the graph in the proof of the above theorem is a version of a domino graph introduced in [10], [11]. By a rough approximation we can say the complexity of our algorithm is exponential because the graph has an exponential number of nodes (it is enough to count the sequences of 0 and 1's in the last two components of the nodes).

# References

[1] M. Andraşiu, A. Atanasiu, Gh. Păun, A. Salomaa. A New Cryptosystem Based on Formal Language Theory, *Bull. Math. Soc. Sci. Math. Roumanie*, Tome 36(84), nr.1, 1992, pp. 1-15.

[2] M.A. Harrison. *Introduction to Formal Language Theory*, Addison-Wesley 1978.

[3] J.E. Hopcroft, J.D. Ullman. *Introduction to Automata Theory, Languages, and Computation*, Addison-Wesley, 1979.

[4] E. Mäkinen. $(0,1)$-Totality is Undecidable for Arbitrary Context-Free Grammars, *Fundamenta Informaticae* (to appear).

[5] C. Matei, F.L. Ţiplea. $(0,1)$- Total Pure Context-Free Grammars, in *Proc. of 2nd Conference on Development in Language Theory*, Magdeburg, 1995.

[6] H.A. Maurer, A. Salomaa, D. Wood. Pure Grammars, *Inform. and Control* 44 (1), 1980, pp. 47-72.

[7] E. Moriya. Associate Languages and Derivational Complexity of Formal Grammars and Languages, *Inform. and Control* 22, 1973, pp. 139-162.

[8] M. Penttonen. On Derivation Languages Corresponding to Context-free Grammars, *Acta Inform.* 3, 1973, pp. 285-291.

[9] A. Salomaa. *Formal Languages*, Academic Press, New York, London, 1973.

[10] A. Weber, T. Head. The Finest Homophonic Partition and Related Code Concepts, in *Proc. of the 19th Symposium on Mathematical Foundations of Computer Science*, Kosice, Slovakia, 1994.

[11] A. Weber. Computing Deciphering and Synchronization Delays of a Code by means of Dominoes, J.W. Goethe-Universität, Frankfurt am Main, 1995.