

## DESCRIPTIONAL COMPLEXITY MEASURES OF CONTEXT-FREE LANGUAGES

ANTON ČERNÝ

The properties of several new descriptonal complexity measures of context-free languages are discussed. Though these measures seem to be very simple the basic algorithmic problems remain to be undecidable.

### 1. INTRODUCTION

In [2], [3] and [4] several measures of descriptonal complexity of context-free grammars (cfg's) and context-free languages (cfl's) have been investigated, most of them having the following properties:

1. The corresponding hierarchy of complexity classes of languages over two-letter alphabets is infinite.
2. The basic algorithmic problems are undecidable. (For example, the problems to determine the complexity of a language generated by a given grammar, to decide whether the given grammar is minimal or to construct an equivalent minimal grammar.)

In the present paper an attempt is made to investigate measures which seem to be simpler from two points of view. Over a fixed alphabet they induce a finite hierarchy of languages, and to determine the complexity of a grammar only a part of the grammar has to be considered. In spite of this for most of these new measures the basic algorithmic problems remain to be undecidable.

### 2. PRELIMINARIES

A survey of the descriptonal complexity theory of formal languages is given in [4]. The basic notions of context-free languages theory to be used here are from [1].

A *cfg* is a quadruple  $G = (V, \Sigma, P, S)$  with  $V$  and  $\Sigma \subseteq V$  being finite sets of symbols.

$P$  a finite set of productions of the form  $A \rightarrow x$  where  $A \in V - \Sigma$ ,  $x \in V^*$ , and  $S \in V - \Sigma$  the start symbol. The elements of  $\Sigma$ , resp.  $V - \Sigma$ , are called *terminals*, resp. *variables*. We write  $w_1Aw_2 \Rightarrow w_1xw_2$  iff  $A \rightarrow x$  is in  $P$  and  $w_1, w_2 \in V^*$ . The relation  $\Rightarrow^*$  is a reflexive and transitive closure of  $\Rightarrow$ . The language defined by a cfg  $G$  is  $L(G) = \{w \in \Sigma^*; S \Rightarrow^* w\}$ .

In Section 4 the  $\varepsilon$ -reduced form ( $\varepsilon$  being the empty word) of cfg's will be often used. A cfg  $G = (V, \Sigma, P, S)$  is said to be in  $\varepsilon$ -reduced form iff

- a) for no variable  $A \neq S$ ,  $A \rightarrow \varepsilon$  is in  $P$
- b) if  $S \rightarrow \varepsilon$  is in  $P$ , then  $S$  does not appear in any of the right sides of productions in  $P$ .

The undecidability results in Sections 3 and 4 will be obtained using a reduction to the Post correspondence problem. In doing that a class of languages, denoted by  $L'_{x,y}$ , will be used. To define  $L'_{x,y}$  the languages  $L_{x,y}$  described in [1] are used.

Let  $x = (x_1, \dots, x_n)$ ,  $y = (y_1, \dots, y_n)$  be ordered  $n$ -tuples ( $n \geq 1$ ) of nonempty words over the alphabet  $\{a, b\}$ . Let

$$L(x) = \{ba^{i_k} \dots ba^{i_1}cx_{i_1} \dots x_{i_k}; k \geq 1, n \geq i_j \geq 1, j = 1, 2, \dots, k\}$$

$$L(x, y) = L(x) \cup L(y)^R$$

$$L_s = \{w_1cw_2cw_2^Rcw_1^R; w_1, w_2 \in \{a, b\}^*\}$$

$$L_{x,y} = \{a, b, c\}^* - L_s \cap L(x, y)$$

According to [1],  $L_{x,y}$  is a cfl and a cfg  $G_{x,y}$  generating  $L_{x,y}$  can be effectively constructed given  $x$  and  $y$ . Moreover,  $L_{x,y} = \{a, b, c\}^*$  iff the Post correspondence problem for the pair  $(x, y)$  has no solution.

The languages  $L'_{x,y}$  can now be defined by  $L'_{x,y} = h(L_{x,y})$  where  $h: \{a, b, c\}^* \rightarrow \{a, b\}^*$  be the homomorphism defined by  $h(a) = ab$ ,  $h(b) = a^2b^2$ ,  $h(c) = a^3b^3$ .

Clearly, a cfg  $G'_{x,y}$  generating  $L'_{x,y}$  can be effectively constructed given  $x$  and  $y$ , and  $L'_{x,y} = h(\{a, b, c\}^*)$  iff the Post correspondence problem for the pair  $(x, y)$  has no solution.

In our last definition, the notion of descriptonal complexity measure ([4]) is introduced. Let  $\mathcal{G}(\mathcal{L})$  denote the class of all cfg's (cfl's) and  $\mathbb{N}$  the set of all nonnegative integers. A *descriptonal complexity measure* of cfg's (cfl's) is an arbitrary mapping

$$\mathbf{K}: \mathcal{G} \rightarrow \mathbb{N} \quad (\mathbf{K}: \mathcal{L} \rightarrow \mathbb{N})$$

Every complexity measure  $\mathbf{K}$  of cfg's induces a complexity measure of cfl's. This measure is also denoted  $\mathbf{K}$  and is defined as follows

$$\mathbf{K}(L) = \min \{ \mathbf{K}(G); G \in \mathcal{G}, L(G) = L \}$$

for every cfl  $L$ .

### 3. PRODUCTIONS WITH TERMINALS

The number of productions of a cfg  $G = \text{Prod}(G)$  is one of the basic complexity measures of cfg's. One way to define a simpler measure seems to be to count some special productions only.

Since variables and terminals play such an important role in the definition of cfg's, it is natural to consider as a complexity measure the number of those productions in which terminals and/or variables satisfy some special condition. For example, productions with at least one (with no) terminal, productions with at least one (with no) variable in the right side and so on. The measures of the former type are investigated in this section. We consider the number of productions with a terminal as a complexity measure.

In the following definition two terminal-based complexity measures are defined. They differ only in the case when a cfg contains  $\varepsilon$ -productions.

**Definition 1.** Let  $G = (V, \Sigma, P, S)$  be a cfg. Then  
 $PT(G)$  = the number of productions in  $P$  with right hand side containing at least one terminal  
 $PT_\varepsilon(G) = PT(G) +$  the number of  $\varepsilon$ -productions in  $P$

The basic relations between these two measures and the number of symbols of the underlying alphabet are summarized in the following lemma.

**Lemma 1.** Let  $\Sigma$  be a finite alphabet and  $L \subseteq \Sigma^*$  a cfl. Then

- (i)  $PT(L) \leq PT_\varepsilon(L) \leq PT(L) + 1$
- (ii)  $PT(L) \leq |\Sigma|$

In spite of property (ii) of Lemma 1 measures  $PT$  and  $PT_\varepsilon$  induce an infinite hierarchy of complexity classes of cfl's with no gaps. Indeed, it is easy to verify that for any integer  $n \geq 0$   $PT(L_n) = n$  ( $PT_\varepsilon(L_n) = n$ ) holds for the languages  $L_0 = \emptyset$ ,  $L_n = \{a_1, \dots, a_n\}$  over the alphabets  $\Sigma_n = \{a_1, \dots, a_n\}$  ( $n \geq 1$ ).

Though measures  $PT$  and  $PT_\varepsilon$  are relatively simple, the basic algorithmic problems for them are undecidable. The proofs of undecidability are based on the following lemma.

**Lemma 2.** There is no algorithm to decide for an arbitrary cfg  $G$  whether or not  $PT(L(G)) = 1$  ( $PT_\varepsilon(L(G)) = 2$ ).

**Proof.** Let  $D$  denote the Dyck language generated by the grammar with two productions

$$S \rightarrow SaSb \mid \varepsilon$$

Let  $h$  be the homomorphism from Section 2. For every pair  $(x, y)$  of  $n$ -tuples of non-empty words over  $\{a, b\}$  denote

$$L''_{x,y} = (D - ah(\{a, b, c\}^*b)) \cup aL'_{x,y}b$$

Since the language  $ah(\{a, b, c\}^*)b$  is regular, a cfg  $G''_{x,y}$  can be constructed with  $L(G''_{x,y}) = L''_{x,y}$ .

Two cases will be considered now.

I. The Post correspondence problem for  $(x, y)$  possesses no solution. In such a case  $L'_{x,y} = h(\{a, b, c\}^*)$  and therefore  $L''_{x,y} = D$ ,  $PT(L'_{x,y}) = 1$ ,  $PT_\varepsilon(L''_{x,y}) = 2$ .

II. The Post correspondence problem for  $(x, y)$  has a solution. We shall prove by a contradiction that  $PT(L''_{x,y}) \geq 2$  and  $PT_\varepsilon(L''_{x,y}) \geq 3$  in this case.

Let  $G = (V, \{a, b\}, P, S)$  be a grammar generating  $L''_{x,y}$ . Suppose  $PT(G) = 1$ . Since  $ab \in aL'_{x,y}b \subseteq L''_{x,y}$ , the only production in  $P$  with a terminal in the right side has to be of the form

$$(1) \quad A \rightarrow uavbz, \quad uvz \in (V - \{a, b\})^*$$

and there exist words  $v_1, v_2 \in (V - \{a, b\})^*$  such that

$$S \Rightarrow^* v_1 A v_2 \Rightarrow v_1 uavbz v_2 \Rightarrow^* ab$$

However, this is possible only if

$$(2) \quad S \Rightarrow^* v_1 A v_2 \Rightarrow^* A, \quad v \Rightarrow^* \varepsilon, \quad uz \Rightarrow^* \varepsilon$$

Let us now denote  $L(v) = \{w \in \{a, b\}^*; v \Rightarrow^* w\}$  and suppose for a moment that the inclusion  $D \subseteq L(v)$  has been proven. Since there exists a solution of the Post correspondence problem for  $(x, y)$ , there is a word

$$(3) \quad t \in h(\{a, b, c\}^*) - L'_{x,y}$$

Then  $t \in D \subseteq L(v)$ , and by (1), (2)

$$(4) \quad S \Rightarrow^* A \Rightarrow uavbz \Rightarrow^* avb \Rightarrow^* atb$$

is a derivation in  $G$ . Consequently,  $atb \in L''_{x,y}$ . On the other hand, by (3) the following relations are simultaneously valid:

$$atb \notin D - ah(\{a, b, c\}^*)b, \quad atb \notin aL'_{x,y}b$$

From here we obtain that  $atb \notin L''_{x,y}$  what contradicts to (4). This means that  $PT(L''_{x,y}) \geq 2$  (and  $PT_\varepsilon(L''_{x,y}) \geq 3$ ) in this case. Therefore to finish the proof of lemma the inclusion  $D \subseteq L(v)$  remains to be shown. Since  $\varepsilon$  is in  $L(v)$ , in order to prove  $D \subseteq L(v)$  it is sufficient to show that  $w_1, w_2 \in D \cap L(v)$  implies  $w_1 a w_2 b \in D \cap L(v)$ . Suppose  $w_1, w_2 \in D \cap L(v)$ . The word  $w = a^4 w_1 a w_2 b b^4$  is from  $D - ah(\{a, b, c\}^*)b$ . Indeed,  $w$  is in  $D$  and none of the words from  $ah(\{a, b, c\}^*)b$  contains the subword  $b^5$ . Consequently,  $w$  is from  $L''_{x,y} = L(G)$ . Then  $w = x_1 x_2$  where  $A \Rightarrow uavbz \Rightarrow^* x_1$  (an initial substring has to be generated using (1)). From (2) we obtain

$$S \Rightarrow^* A \Rightarrow^* x_1 \in L''_{x,y} \subseteq D$$

The equality  $x_1 = w$  can be easily shown, as each initial substring of  $w$  except  $w$  itself contains more  $a$ 's than  $b$ 's what is impossible for a word from  $D$ . Thus each derivation of  $x_1$  from  $uavbz$  has to be of the form

$$(5) \quad uavbz \Rightarrow^* av_3b \Rightarrow^* x_1 = w, \quad \text{where} \quad v \Rightarrow^* v_3 \Rightarrow^* a^3 w_1 a w_2 b b^3, v_3 \in V^*$$

Since in (5) the production (1) has to be used again, the existence of the following derivations can be shown similarly:

$$\begin{aligned} v &\Rightarrow^* a^2 w_1 a w_2 b b^2 \\ v &\Rightarrow^* a w_1 a w_2 b b \\ v &\Rightarrow^* w_1 a w_2 b \end{aligned}$$

Therefore  $D \subseteq L(v)$  and this completes the proof of lemma.  $\square$

The unsolvability of determining  $PT(L(G))$  ( $PT_\epsilon(L(G))$ ) can be formulated now in a stronger form.

**Theorem 1.**

(i) For no integer  $n \geq 1$  there is an algorithm to decide for a given grammar  $G$  whether or not  $PT(L(G)) = n$  ( $PT_\epsilon(L(G)) = n + 1$ ).

(ii) The problems  $PT(L(G)) = 0$ ,  $PT_\epsilon(L(G)) = 0$ ,  $PT_\epsilon(L(G)) = 1$  are decidable.

*Proof.* (i) For  $n = 1$  the theorem follows from Lemma 2. Let now  $n \geq 2$ . Let  $\Sigma_{n-1} = \{a_1, \dots, a_{n-1}\}$  be an alphabet, such that  $\Sigma_{n-1} \cap \{a, b\} = \emptyset$ . Denote  $L_{x,y}^n = L_{x,y}'' \cup \Sigma_{n-1}$ . Clearly,  $PT(\Sigma_{n-1}) = PT_\epsilon(\Sigma_{n-1}) = n - 1$ . By the same reasoning as in the proof of Lemma 2 one can show that  $PT(L_{x,y}^n) = n$  ( $PT_\epsilon(L_{x,y}^n) = n + 1$ ) iff the Post correspondence problem has no solution for  $(x, y)$ . (i) now follows from undecidability of the Post correspondence problem.

(ii) From Definition 1 we easily obtain

$$\begin{aligned} PT(L(G)) = 0 &\quad \text{iff} \quad L(G) = \emptyset \quad \text{or} \quad L(G) = \{\epsilon\} \\ PT_\epsilon(L(G)) = 0 &\quad \text{iff} \quad L(G) = \emptyset \end{aligned}$$

Since  $L(G) = \emptyset$  and  $L(G) = \{\epsilon\}$  is decidable we have decidability of the first two problems in (ii). The decidability of the third problem we get from the decidability of the inclusion  $L(G) \subseteq w^*$  ( $[1]$ ) since

$$(6) \quad PT_\epsilon(L(G)) = 1 \quad \text{iff} \quad L(G) \neq \emptyset \quad \text{and there is a word } w \text{ such that} \\ L(G) \subseteq w^+$$

The only trouble in the proof of (6) seems to be with the if-part for the case  $w \neq \epsilon$ . Let  $L(G) \subseteq w^+ \subseteq \Sigma^*$  with  $\Sigma$  being a finite alphabet,  $w \neq \epsilon$ . Let  $h_1 : a^* \rightarrow \Sigma^*$  be a homomorphism,  $h_1(a) = w$ . By  $[1]$   $h_1^{-1}(L(G))$  is a cfl not containing the empty

word. It is easy to see that

$$PT_e(L(G)) \leq PT_e(h_1^{-1}(L(G)))$$

and by Lemma 1  $PT_e(h_1^{-1}(L(G))) = 1$ .  $\square$

There are two interesting problems concerning the minimality of a cfg with respect to  $PT(PT_e)$ . The undecidability of the first one is an immediate corollary of Lemma 2, the undecidability of the second one needs a short proof.

**Corollary 1.** There is no algorithm to construct for a given cfg  $G$  an equivalent cfg  $G'$  such that  $PT(G') = PT(L(G'))$  ( $PT_e(G') = PT_e(L(G'))$ ).

**Theorem 2.** There is no algorithm to decide for a given cfg  $G$  whether or not  $PT(G) = PT(L(G))$  ( $PT_e(G) = PT_e(L(G))$ ).

*Proof.* It is easily to see that given  $x$  and  $y$  a cfg  $G$  can be constructed such that  $L(G) = L''_{x,y}$ ,  $PT(G) = 2$  and  $PT_e(G) = 3$ . It has been shown that  $PT(L''_{x,y}) < 2$  ( $PT_e(L''_{x,y}) < 3$ ) iff the Post correspondence problem for  $(x, y)$  has no solution. However, this implies that  $G$  is a minimal cfg with respect to the complexity measure  $PT(PT_e)$  iff the Post correspondence problem for  $(x, y)$  possesses a solution. The theorem now follows from undecidability of the Post correspondence problem.

#### 4. END PRODUCTIONS

The structure of leaf-parts of derivation trees in cfg's depends on the number and form of productions with no variable in the right hand side. Thus the number of such "end productions" as the complexity measure for cfg's could be of some importance. Similarly, the minimal number of distinct left sides of the end productions gives us some information about the intrinsic complexity of the language.

**Definition 2.** Let  $G = (V, \Sigma, P, S)$  be a cfg. A production  $A \rightarrow w$  will be called *end production* iff  $w \in \Sigma^*$ . A variable  $A$  will be called *end variable* iff there is an end production  $A \rightarrow w$  in  $P$ .

**Definition 3.** Let  $G = (V, \Sigma, P, S)$  be a cfg. If  $G$  is a grammar in  $\varepsilon$ -reduced form, then

$$EP(G) = \text{the number of end productions in } P$$

$$EV(G) = \text{the number of end variables in } P$$

Otherwise

$$EP(G) = EV(G) = |\Sigma| + 1$$

Definition 3 needs some explanation. For an arbitrary grammar  $G$  an equivalent grammar  $G'$  with no end productions can be constructed. Indeed, by concatenating

the right side of each production with a new symbol  $X$  and adding  $X \rightarrow \varepsilon$  to the productions of  $G$  such a grammar is obtained. The separate definition of the measures  $EP$  and  $EV$  for grammars not being in  $\varepsilon$ -reduced form makes both measures non-trivial, as shown in Lemma 3.

**Lemma 3.** Let  $n \geq m \geq 1$  be integers and  $\Sigma_n = \{a_1, \dots, a_n\}$  an alphabet. Then

- (i) for any cfl  $L \subseteq \Sigma_n^*$   $EV(L) \leq EP(L) \leq n$
- (ii) there is a language  $L_{m,n} \subseteq \Sigma_n^*$  such that  $EV(L_{m,n}) = m$ ,  $EP(L_{m,n}) = n$ .

**Proof.** (i) is obvious. To prove (ii) denote  $L_{m,n} = a_1^+ \cup \dots \cup a_{m-1}^+ \cup a_m \cup \dots \cup a_n$ . This language is defined by the grammar with productions

$$\begin{aligned} S &\rightarrow A_i \mid a_j, \quad i = 1, 2, \dots, m-1 \\ A_i &\rightarrow a_i A_i \mid a_i, \quad j = m, m+1, \dots, n \end{aligned}$$

and therefore

$$(7) \quad EV(L_{m,n}) \leq m, \quad EP(L_{m,n}) \leq n$$

Let now  $G = (V, \Sigma_n, P, S)$  be an arbitrary grammar in  $\varepsilon$ -reduced form generating  $L_{m,n}$ .  $L_{m,n}$  contains at least one word from each of the languages  $a_i^+$ ,  $i = 1, 2, \dots, n$ . Hence for each integer  $i = 1, 2, \dots, n$  there is a production

$$(8) \quad B_i \rightarrow a_i^{k_i}, \quad k_i \geq 1$$

in  $P$ . Clearly, we may assume that  $k_i = 1$  for  $i = m, \dots, n$  and that (8) is used in the derivation of some word  $w_i \in a_i^+$ ,  $|w_i| > k_i$  for  $i = 1, 2, \dots, m-1$ .

From (8) we immediately get  $EP(L_{m,n}) \geq n$ , and consequently,  $EP(L_{m,n}) = n$  by (7).

Suppose there are integers  $r, s$ ,  $1 \leq r < s \leq m$ , such that  $B_r = B_s$ . Then the following two derivations exist in  $G$ :

$$\begin{aligned} S &\Rightarrow^* b_r^p B_r b_r^q \Rightarrow b_r^{p+q+k_r} = w_r \in L_{m,n}, \quad p+q > 0 \\ S &\Rightarrow^* b_r^p B_r b_r^q \Rightarrow b_r^p B_s b_r^q \Rightarrow b_r^p b_s^{k_s} b_r^q \in L_{m,n} \end{aligned}$$

This contradicts to the fact that no word in  $L_{m,n}$  contains two different symbols. Thus  $EV(L_{m,n}) \geq m$  and by (7)  $EV(L_{m,n}) = m$ .  $\square$

The decidability results for  $EP$ ,  $EV$  and  $PT$ ,  $PT_c$  complexity measures are quite similar and the following lemma corresponds to Lemma 2.

**Lemma 4.** There is no algorithm to decide for an arbitrary cfg  $G$  whether or not  $EP(L(G)) = 1$  ( $EV(L(G)) = 1$ ).

**Proof.** Denote  $L_{x,y}'' = L_{x,y}'c$ . In order to prove the lemma it is sufficient to show that  $EP(L_{x,y}''') = 1$  ( $EV(L_{x,y}''') = 1$ ) iff the Post correspondence problem for  $(x, y)$  possesses no solution. To this end two cases will be considered.

I. If the Post correspondence problem for  $(x, y)$  has no solution, then  $L_{x,y}''' =$

$= h(\{a, b, c\}^*)c$  ( $h$  is the homomorphism defined in Section 2). This language is generated by the grammar  $S \rightarrow abS \mid a^2b^2S \mid a^3b^3S \mid c$  and therefore  $EP(L''_{x,y}) = EV(L''_{x,y}) = 1$ .

II. Let the Post correspondence problem for  $(x, y)$  have a solution, i.e. there exist words  $u, v, z, t \in \{a, b\}^*$ ,  $z = v^R$ ,  $t = u^R$ , such that for no integer  $m \geq 1$  the word  $h(u^m cv^m cz^m ct^m)c$  is in  $L''_{x,y}$ . Suppose there exists a cfg  $G = (V, \{a, b, c\}, P, S)$  in  $\varepsilon$ -reduced form which generates  $L''_{x,y}$ , and  $EV(G) = 1$ .

Since  $c \in L''_{x,y}$ , the production  $A \rightarrow c$  where  $A$  is the only end variable of  $G$  has to be in  $P$ .

At first we shall show that  $EV(G) = 1$  implies that  $L''_{x,y}$  is a regular language. Let

$$(9) \quad B \rightarrow w_1 C w_2$$

be an arbitrary production from  $P$  used in some derivation of a terminal word, with  $w_1 \in \{a, b, c\}^*$ , i.e.  $C$  is the first variable from left in the right hand side of (9). Then there are words  $w_3, \dots, w_8 \in V^*$  and the derivation

$$\begin{aligned} S &\Rightarrow^* w_3 B w_4 \Rightarrow w_3 w_1 C w_2 w_4 \Rightarrow^* w_3 w_1 w_5 A w_6 w_2 w_4 \Rightarrow \\ &\Rightarrow w_3 w_1 w_5 c w_6 w_2 w_4 \Rightarrow^* w_7 c w_8 \in L''_{x,y} \end{aligned}$$

( $A$  is the only end variable of  $G$ !). Since the only possible position for a symbol  $c$  in  $L''_{x,y}$  is at the end of the words, we get  $w_8 = w_2 = \varepsilon$ . Hence every non-superfluous production of  $G$  is right-linear and  $L''_{x,y}$  is a regular language.

However, from the existence of a solution of the Post correspondence problem for  $(x, y)$  the opposite follows. Indeed, by Theorem 5.6 of [5] if  $L''_{x,y}$  is regular, then the equivalence relation  $\mathbb{E}$  on  $\{a, b, c\}^*$  induced by  $L''_{x,y}$  is of finite index. Thus there exist integers  $j, k \geq 1$ ,  $j \neq k$ , such that  $h(u^j) \mathbb{E} h(u^k)$ . This contradicts to the fact that  $h(u^j) h(cv^j cz^j ct^j)c$  is from  $L''_{x,y}$  while  $h(u^k) h(cv^k cz^k ct^k)c$  does not belong to  $L''_{x,y}$ . Consequently, grammar  $G$  with the described properties cannot exist, and  $EP(L''_{x,y}) \geq EV(L''_{x,y}) \geq 2$ . The lemma follows now from the undecidability of the Post correspondence problem.  $\square$

Using Lemma 4 the following theorem can be proven.

**Theorem 3.**

- (i) For no integer  $n \geq 1$  there is an algorithm to decide for an arbitrary cfg  $G$  whether or not  $EP(L(G)) = n$  ( $EV(L(G)) = n$ ).
- (ii) The problem  $EP(L(G)) = 0$  ( $EV(L(G)) = 0$ ) is decidable.

Proof. (i) For  $n = 1$  the theorem follows from Lemma 4. Let  $n \geq 2$ . Let  $\Sigma_{n-1} = \{a_1, \dots, a_{n-1}\}$  be an alphabet,  $\Sigma_{n-1} \cap \{a, b, c\} = \emptyset$ . Denote  $L''_{x,y} = L''_{x,y} \cup L_{n-1}$  where  $L''_{x,y}$  is the language from the proof of Lemma 4 and  $L_{n-1} = a_1^+ \cup \dots \cup a_{n-1}^+$ . Since the languages  $L''_{x,y}$  and  $L_{n-1}$  are over distinct alphabets we get easily

$$K(L''_{x,y}) = K(L''_{x,y}) + K(L_{n-1}) \quad \text{for } K \in \{EP, EV\}.$$



By a similar reasoning as in the proofs of Lemma 3 and Lemma 4 we get  $EP(L_{x,y}^n) = n$  ( $EV(L_{x,y}^n) = n$ ) iff the Post correspondence problem for  $(x, y)$  has no solution. The theorem follows now from the undecidability of the Post correspondence problem.

(ii) Clearly,  $EP(L(G)) = 0$  ( $EV(L(G)) = 0$ ) iff  $L(G) = \emptyset$  or  $L(G) = \{\varepsilon\}$ . The last two problems are known to be decidable.  $\square$

**Corollary 2.** There is no algorithm to construct for a given cfg  $G$  an equivalent cfg  $G'$  minimal with respect to  $EP$  ( $EV$ ).

The problem of minimality is undecidable for  $EP$ ,  $EV$ , too:

**Theorem 4.** There is no algorithm to decide for a given cfg  $G$  whether or not  $EP(G) = EP(L(G))$  ( $EV(G) = EV(L(G))$ ).

*Proof.* Obviously, a cfg  $G$  in  $\varepsilon$ -reduced form generating  $L_{x,y}$  can be constructed, with  $A \rightarrow a$ ,  $B \rightarrow b$ ,  $C \rightarrow c$  as the only end productions. If these productions are replaced by the productions  $A \rightarrow Db$ ,  $B \rightarrow Dab^2$ ,  $C \rightarrow Da^2b^3$ ,  $D \rightarrow a$ , a new grammar generating  $L'_{x,y}$  is obtained. From this grammar by a small modification a grammar  $G'$  generating  $L''_{x,y}$  can be obtained such that  $EP(G') = EV(G') = 2$ . From the proof of Lemma 4 it follows that  $G'$  is minimal with respect to  $EP(EV)$  iff the Post correspondence problem for  $(x, y)$  has a solution. The theorem now follows from the undecidability of the Post correspondence problem.  $\square$

To finish this section three another complexity measures similar to the measures  $EP$  and  $EV$  are defined and some results concerning these measures are summarized.

**Definition 4.** Let  $G = (V, \Sigma, P, S)$  be a cfg. If  $G$  is in  $\varepsilon$ -reduced form, then

$PTB(G)$  = the number of productions in  $P$  with a terminal at the beginning of the right hand side

$PTE(G)$  = the number of productions in  $P$  with a terminal at the end of the right hand side

$PTBE(G)$  = the number of productions in  $P$  with a terminal at the beginning and at the end of the right hand side.

Otherwise

$$PTB(G) = PTE(G) = PTBE(G) = |\Sigma| + 1$$

**Theorem 5.**

(i) Let  $\Sigma$  be a finite alphabet,  $L \leq \Sigma^*$  a cfl. Then

$$EP(L) \leq PTBE(L) \leq PTB(L) \leq |\Sigma|$$

$$PTBE(L) \leq PTE(L) \leq |\Sigma|$$

(ii) Let  $K \in \{PTB, PTE, PTBE\}$ . There is no algorithm

1. to decide for an arbitrary cfg  $G$  and integer  $n \geq 1$  whether or not  $K(L(G)) = n$
2. to decide for a given cfg  $G$  whether or not  $K(G) = K(L(G))$
3. to construct for a given cfg  $G$  and equivalent cfg  $G'$  minimal with respect to  $K$ .

## 5. COMPLEMENTARY MEASURES

Measures of complexity in Sections 2 and 3 have been defined according to the following schema:

For each cfg  $G = (V, \Sigma, P, S)$  and each production  $p \in P$  an integer  $k(G, p) \in \{0, 1\}$  has been defined. The complexity of a cfg  $G$  has then been defined by

$$K(G) = \sum_{p \in P} k(G, p)$$

In other words, only a part of productions of a cfg  $G$  contributes to the overall complexity of  $G$ . Once this point of view is accepted it is quite natural to investigate the so-called complementary complexity measures to which the remaining part of productions contributes. They are defined by

$$CK(G) = \sum_{p \in P} (1 - k(G, p))$$

and will be discussed in this section.

**Definition 5.** For a cfg  $G = (V, \Sigma, P, S)$  let us define

$CPT(G)$  = the number of productions in  $P$  without terminals in the right side

$CPT_\epsilon(G)$  =  $CPT(G)$  – the number of  $\epsilon$ -productions in  $P$

$CEP(G)$  = the number of productions in  $P$  not being end productions

$CEV(G)$  = the number of variables in  $V$  not being end variables

$CPTB(G)$  = the number of productions in  $P$  with the right hand side not beginning with a terminal

$CPTE(G)$  = the number of productions in  $P$  without terminal at the end of the right hand side

$CPTBE(G)$  = the number of productions in  $P$  with the right hand side not beginning or not ending by a terminal

The main results of this section are summarized in two theorems.

**Theorem 6.** There is no algorithm to determine  $CEP(L(G))$  for an arbitrary cfg  $G$ .

**Theorem 7.** Let  $L$  be a cfl. Then

(i) If  $L = \emptyset$  or  $L = \{\epsilon\}$  then  $CEV(L) = 1$ , otherwise  $CEV(L) = 0$

(ii)  $CPT_\epsilon(L) = 0$

(iii) If  $\epsilon \in L$ , then

$$CPT(L) = CPTB(L) = CPTE(L) = CPTBE(L) = 1$$

Otherwise

$$CPT(L) = CPTB(L) = CPTE(L) = CPTBE(L) = 0$$

Hence, the basic algorithmic problems are decidable for all complementary measures of Definition 5 with the exception of the measure  $CEP$ .

**Proofs.** The proof of Theorem 6 is similar to that of Lemma 2. Let us denote  $\bar{L}_{x,y} = L_{x,y} - \{e\}$  and let us consider two cases.

I. If the Post correspondence problem for  $(x, y)$  possesses no solution, then  $\bar{L}_{x,y} = \{a, b, c\}^+$ . This language is generated by the grammar  $S \rightarrow SS \mid a \mid b \mid c$  and therefore  $CEP(\bar{L}_{x,y}) = 1$ .

II. Let the Post correspondence problem for  $(x, y)$  have a solution. Suppose there is a cfg  $G$  with  $L(G) = \bar{L}_{x,y}$ ,  $CEP(G) = 1$ . Then the start symbol is the only recursive variable in  $G$ . Now a grammar  $G'$  equivalent to  $G$  can be constructed containing just one variable. However, this contradicts to Lemma 3.2 in [3], which says that at least two variables are necessary to generate  $\bar{L}_{x,y}$ . It means that there is no grammar  $G$  for  $\bar{L}_{x,y}$  with  $CEP(G) = 1$  and therefore  $CEP(\bar{L}_{x,y}) \geq 2$ . Theorem 6 follows now from the undecidability of the Post correspondence problem.

If  $G$  is a cfg and  $A$  a variable of  $G$  such that  $A \Rightarrow^* w \in L(G)$ , then by adding  $A \rightarrow w$  to the productions of  $G$  an equivalent grammar is obtained. From that (i) of Theorem 7 follows easily. The parts (ii) and (iii) of Theorem 7 follow from the well known fact that to every cfg there exists an equivalent cfg  $G'$  in  $\varepsilon$ -reduced form all productions of which are of the  $A \rightarrow a$  or  $A \rightarrow axb$  with  $a, b$  being terminals, or  $S' \rightarrow \varepsilon$  where  $S'$  is the start symbol of  $G'$ .  $\square$

#### ACKNOWLEDGEMENT

The author is indebted to Dr. J. Gruska for valuable and stimulating discussions and for reading the manuscript.

(Received July 15, 1977.)

#### REFERENCES

- [1] S. Ginsburg: The Mathematical Theory of Context-Free Languages. McGraw-Hill, New York 1966.
- [2] J. Gruska: Some classifications of context-free languages. *Information and Control* 14 (1969), 2, 152–179.
- [3] J. Gruska: Complexity and unambiguity of context-free grammars and languages. *Information and Control* 18 (1971), 5, 502–519.
- [4] J. Gruska: Descriptive complexity of context-free languages. *Proceedings MFCS' 73*, High Tatras, 71–83.
- [5] A. Salomaa: Formal Languages. Academic Press, New York and London 1973.

*RNDr. Anton Černý, Katedra teoretickej kybernetiky Matematicko-fyzikálnej fakulty Univerzity Komenského (Department of Theoretical Cybernetics, Faculty of Mathematics and Physics – Comenius University), Mlynská dolina, 842 15 Bratislava. Czechoslovakia.*