

## COMPLEXITY OF NORMAL FORM GRAMMARS\*

Alica KELEMENOVÁ

*Mathematical Institute, Slovak Academy of Sciences, 814 73 Bratislava, Czechoslovakia*

Communicated by A. Salomaa

Received October 1982

Revised March 1983

**Abstract.** Various types of grammars can be used to describe context-free languages. Such are context-free grammars and their normal form restrictions. Rewriting of a context-free grammar to an equivalent grammar in required (normal) form can cause a change of parameters of the grammar such as the number of rules, the number of nonterminals, etc. Greibach normal form grammars and position restricted grammars will be investigated from the point of view of descriptiveness of context-free languages.

### 1. Introduction

Various types of grammars can be used to describe context-free languages. Such are context-free grammars and their normal form restrictions (e.g., Chomsky normal form, Greibach normal form, etc.). Also a special subclass of context-sensitive grammars—terminal bounded grammars—has that property [2].

Rewriting of a context-free grammar to an equivalent grammar in required (normal) form can cause, in general, a change of such parameters of the grammar as the number of rules, the number of nonterminals, etc. For instance, the  $\epsilon$ -rule removal substantially increases the number of rules of context-free grammars, while the total length of the grammar can increase at most ten-times [6].

In the present paper Greibach normal form grammars and position restricted grammars will be investigated from the point of view of descriptiveness of context-free languages. We shall prove that the number of nonterminals needed to describe a language by a grammar in Greibach normal form does not exceed twice the number of nonterminals needed to describe it by an  $\epsilon$ -free chain-free reduced grammar. Increasing of productions and total length for grammars in Greibach normal form is at most cubic in the number of productions or in the total length for  $\epsilon$ -free chain-free reduced grammars, respectively. We shall give an example, where increasing is exactly square. (For the results discussed above, terminals are allowed in any position on right-hand sides of rules of grammars in Greibach normal form.)

\* Part of the results of this paper were previously reported in [11].

For the case of language description by position restricted grammars, the increase of values for the number of nonterminals, the number of rules, the number of grammatical levels and the height of grammatical levels cannot be bounded by any function.

For the total length of grammars we shall give optimal linear increase both for position restricted grammars of type  $(0, 0, 0)$  and for grammars in Chomsky normal form with respect to the total length of  $\varepsilon$ -free chain-free reduced grammars.

## 2. Preliminaries

We assume the reader to be familiar with the basic formal language theory in a scope of [8] or [12]. We briefly review some well-known notations and definitions from descriptonal complexity of formal languages but refer to [7] for further details.

In this paper only context-free grammars and languages will be considered.

We shall use  $G = (T, N, P, S)$  for a context-free grammar, where  $N$  and  $T$  denote the sets of nonterminals and terminals, respectively,  $P$  is the set of rules and  $S$  is the initial nonterminal.

To evaluate the complexity of the grammar  $G = (N, T, P, S)$ , we shall use five parameters known from literature. We shall call them by common name complexity measure. Three of them reflect in some way the size of the grammar. They are:

$$\text{Var}(G) = |N| \quad (\text{the number of nonterminals}).$$

$$\text{Prod}(G) = |P| \quad (\text{the number of productions}).$$

$$\text{Symb}(G) = \sum (2 + |\alpha|) \quad (\text{the total length of } G)$$

where the sum goes over all productions  $A \rightarrow \alpha$  in  $P$ .

Two further parameters, based on the notion of grammatical level, reflect the structure of grammars. To introduce them, first we define the relation  $\triangleright$  on the nonterminals of the grammar  $G = (N, T, P, S)$ .

For  $A, B$  in  $N$ ,  $A \triangleright B$  if there is a production  $A \rightarrow xBy$  in  $P$ . The relation  $\triangleright^*$  is the reflexive and transitive closure of  $\triangleright$ .

We say that nonterminals  $A$  and  $B$  are structurally equivalent and we write  $A \equiv B$  if  $A \triangleright^* B$  and  $B \triangleright^* A$ .

Every class of partition of  $N$  by the equivalence relation  $\equiv$  is called a grammatical level of  $G$ . For two different grammatical levels  $Q_1$  and  $Q_2$  we write  $Q_1 > Q_2$  if there are nonterminals  $A$  in  $Q_1$  and  $B$  in  $Q_2$  such that  $A \triangleright B$ . By the graph of grammatical levels of  $G$  we shall understand the digraph, the nodes of which are grammatical levels of  $G$  and edges are ordered pairs  $(Q_1, Q_2)$  for  $Q_1 > Q_2$ .

Now we are ready to introduce complexity measures Lev and Hei.

$$\text{Lev}(G) = \text{“the number of levels in } G\text{”}.$$

$$\text{Hei}(G) = \max\{\text{Hei}(Q) : Q \text{ is a grammatical level of } G\},$$

where  $\text{Hei}(Q) = 1$  iff  $S \in Q$  and  $\text{Hei}(Q_i) = 1 + \max\{\text{Hei}(Q_i): Q_i \supset Q\}$ , i.e.,  $\text{Hei}(G)$  denotes the length of the longest way in the acyclic graph of grammatical levels, starting with the initial level (i.e., the level containing  $S$ ).

Any complexity measure  $K$  from Var, Prod, Symb, Lev and Hei allows us to introduce the grammatical complexity of the language  $L$  generated by grammars from  $\mathcal{G}$  in the following manner:

$$K_{\mathcal{G}}(L) = \min\{K(G): L(G) = L, G \in \mathcal{G}\}.$$

For  $\mathcal{G}$  being the class of all context-free grammars we shall omit the subscript  $\mathcal{G}$  in  $K_{\mathcal{G}}(L)$ .

By the complexity of class  $\mathcal{L}$  of languages we shall mean

$$K_{\mathcal{G}}(\mathcal{L}) = \sup\{K_{\mathcal{G}}(L): L \in \mathcal{L}\}.$$

A grammar  $G = (N, T, P, S)$  is reduced if

- (1) for any  $A \in N$ ,  $S \Rightarrow^* \alpha A \beta$  for some  $\alpha, \beta$  in  $(N \cup T)^*$ , and
- (2) for any  $A \in N$  there is a  $w$  in  $T^*$  such that  $A \Rightarrow^* w$ .

We shall use  $\mathcal{R}$  to denote the class of all completely reduced grammars, i.e., reduced grammars with no  $\varepsilon$ -rules (i.e., that of  $A \rightarrow \varepsilon$ ) and no chain rules (i.e., that of type  $A \rightarrow B$ ,  $A, B$  in  $N$ ).

In this paper the problem to compare the values  $K_{\mathcal{N}}(L)$  and  $K_{\mathcal{R}}(L)$  will be investigated, where  $\mathcal{R}$  is the class of completely reduced grammars and  $\mathcal{N}$  is the class of grammars in Greibach normal form or alternatively the class of position restricted grammars of type  $t$ . We prefer to compare normal form complexity with that on  $\mathcal{R}$  instead of that on the whole of context-free grammars since comparison for  $K$  and  $K_{\mathcal{R}}$  gives nontrivial results [6, 7]. Also, known techniques transforming a grammar to an equivalent one in Greibach normal form or to an equivalent one in position restricted grammar (especially to a grammar in Chomsky normal form) use as input a grammar in completely reduced form or produce a completely reduced grammar as an intermediate result of the algorithm.

For our purpose it is sufficient to take into consideration  $\varepsilon$ -free languages only. We shall use the following  $\varepsilon$ -free definitions of normal form grammars.

**Definition.** A grammar  $G = (N, T, P, S)$  is in Greibach normal form if all its productions are of type  $A \rightarrow a\alpha$  for  $A \in N$ ,  $a \in T$ ,  $\alpha \in N^*$ .

The class of grammars in Greibach normal form will be denoted by  $\mathcal{G}_t$ .

**Definition.** A grammar  $G = (N, T, P, S)$  is in weak Greibach normal form if all its productions are of type  $A \rightarrow a\alpha$  for  $A \in N$ ,  $a \in T$  and  $\alpha \in (N \cup T)^*$ .

A class of grammars in weak Greibach normal form will be denoted by  $\mathcal{G}'_t$ .

**Definition.** A grammar  $G = (N, T, P, S)$  is a position restricted grammar of type

$t = (m_1, m_2, \dots, m_k, m_{k+1})$ , where  $m_i$  is a nonnegative integer for  $1 \leq i \leq k+1$  if every production of the grammar is either of the form

(i)  $A \rightarrow w_1 A_1 w_2 A_2 \cdots w_k A_k w_{k+1}$ , where  $A, A_i \in N$ ,  $1 \leq i \leq k$  and  $w_j \in T^{m_j}$  for  $1 \leq j \leq k+1$ , or

(ii)  $A \rightarrow aB$  or  $A \rightarrow a$ , where  $A, B \in N$  and  $a \in T$ .

The collection of position restricted grammars of type  $t = (m_1, m_2, \dots, m_k, m_{k+1})$  will be denoted by  $\ell$  if  $t$  is understood or by  $(m_1, \dots, m_k, m_{k+1})$ .

**Remark.** Grammars in Chomsky normal form forms the proper subclass of position restricted grammars of type  $(0, 0, 0)$ . Position restricted grammars are treated in [3, 4, 10].

Algorithms for transformations of grammars to equivalent grammars in weak Greibach normal form are well known and they are described in Appendix A, together with the transformations of grammars to equivalent position restricted grammars of type  $(0, 0, 0)$  and to equivalent grammars in Chomsky normal form.

In order to formulate our results, Knuth's symbols  $O$ ,  $\Omega$  and  $\Theta$  will be used. For  $f, g: \mathbb{Z} \rightarrow \mathbb{Z}$  ( $\mathbb{Z}$  is the set of integers) the meaning of  $g(n) = O(f(n))$ ,  $g(n) = \Omega(f(n))$  and  $g(n) = \Theta(f(n))$  is as in [9], i.e.,

$g(n) = O(f(n))$  if there are positive numbers  $c$  and  $n_0$  such that  $|g(n)| \leq cf(n)$  for all  $n \geq n_0$ ,

$g(n) = \Omega(f(n))$  if there are positive numbers  $c$  and  $n_0$  such that  $g(n) \geq cf(n)$  for all  $n \geq n_0$ ,

$g(n) = \Theta(f(n))$  if there are positive numbers  $c, c'$  and  $n_0$  such that  $cf(n) \leq g(n) \leq c'f(n)$  for all  $n \geq n_0$ .

For a complexity measure  $K$  we shall define the set of languages  $\mathcal{L}_K^{(n)} = \{L: K(L) = n\}$ , where  $n$  is a natural number, and the function  $\mathcal{L}_K$ , which associates with  $n$  the set of languages  $\mathcal{L}_K^{(n)}$ , i.e., for every  $n$ ,  $\mathcal{L}_K(n) = \mathcal{L}_K^{(n)}$ .

**Lemma 2.1.** *Let  $K$  be a complexity measure and  $\mathcal{N}$  be a set of normal form grammars.*

*If for any natural number  $m \geq k_0$  there is a language  $L_m$  and a constant  $c$  such that  $K_B(L_m) = c$  and  $K_N(L_m) = m$ , then  $K_N(\mathcal{L}_K(n))$  cannot be bounded by any function  $\varphi(n)$ .*

**Notation.** (1) Let  $A$  be a nonterminal of a grammar  $G = (N, T, P, S)$ . Then

$$L(A) = \{w: w \in T^*, A \Rightarrow^* w\}.$$

(2) For a letter  $a \in T$  and for a word  $w = x_1 x_2 \cdots x_n$ ,  $x_i \in T$  for  $1 \leq i \leq n$ , we shall write  $a \in \text{alph}(w)$  if  $a = x_j$  for some  $j$ .

For a language  $L$  and for a letter  $a$ ,  $a \in \text{alph}(L)$  if  $a \in \text{alph}(w)$  for some  $w$  in  $L$ .

**Lemma 2.2.** *Let  $G = (N, T, P, S)$  be a completely reduced grammar and  $Z_1, Z_2$  be nonterminals in the same grammatical level. Then*

$$\text{alph}(L(Z_1)) = \text{alph}(L(Z_2)).$$

**Proof.**  $Z_1$  and  $Z_2$  are in the same grammatical level so there are words  $\alpha_1, \alpha_2, \beta_1$  and  $\beta_2$  such that  $Z_2 \Rightarrow^* \alpha_1 Z_1 \beta_1 \Rightarrow^* w_1$  and  $Z_1 \Rightarrow^* \alpha_2 Z_2 \beta_2 \Rightarrow^* w_2$ . The words generated by  $Z_1$  appear as subwords of words generated by  $Z_2$  and vice-versa.  $\square$

**Notation.** Let  $T_1 \subseteq T$ . By  $e_{T_1}$  we denote the homomorphism

$$e_{T_1}(x) = \begin{cases} \varepsilon & \text{if } x \in T_1, \\ x & \text{if } x \in T - T_1. \end{cases}$$

For  $L \subseteq T^*$ ,  $e_{T_1}(L) = \{e_{T_1}(w) : w \in L\}$ .

**Lemma 2.3.** *Let  $G = (N, T, P, S)$  be a reduced grammar. A language  $e_{\Sigma - \{a\}}(L(G))$  is infinite iff  $G$  contains a nonterminal  $A$  such that  $A \Rightarrow^* \alpha A \beta$  and  $a \in \text{alph}(\alpha \beta)$ .*

**Proof.** Obvious.  $\square$

### 3. Greibach normal form complexity

In this section we shall investigate the complexity of the weak Greibach normal form. An obvious correspondence between  $K_{\mathcal{G}_i}$  and  $K_{\mathcal{G}_i'}$  is given by the following theorem.

**Theorem 3.1.** *Let  $L$  be a language and  $|\text{alph}(L)| = k$ . Then*

$$K_{\mathcal{G}_i}(L) \leq k + K_{\mathcal{G}_i'}(L) \quad \text{for } K \in \{\text{Var}, \text{Prod}, \text{Lev}\},$$

$$\text{Symb}_{\mathcal{G}_i}(L) \leq 3k + \text{Symb}_{\mathcal{G}_i'}(L) < 4 \text{Symb}_{\mathcal{G}_i'}(L),$$

$$\text{Hei}_{\mathcal{G}_i}(L) \leq 1 + \text{Hei}_{\mathcal{G}_i'}(L).$$

*All these estimates are definitive (i.e., a sequence of languages  $\{L_k\}_{k=1}^{\infty}$  can be given such that  $|\text{alph}(L_k)| = k$  and in the estimates above equalities hold).*

**Proof.** The last part of algorithm  $\mathcal{A}_1$  (see Appendix A) beginning with label (L3) produces a grammar in Greibach normal form for input a grammar in weak Greibach normal form. The inequalities can be shown by a straightforward analysis of that part of algorithm  $\mathcal{A}_1$ .

The equalities are fulfilled, e.g., for the languages  $L_k = \{a_1^2 a_2 \cdots a_k\}$ .  $\square$

**Theorem 3.2.** *For every language  $L$ ,  $\text{Var}_{\mathcal{G}_i'}(L) \leq 2 \text{Var}_{\mathcal{G}_i}(L)$ . The coefficient 2 in the*

estimate is the best possible [i.e., for every  $\delta > 0$  there is a language  $L_\delta$  such that  $\text{Var}_{\mathcal{G}_i}(L_\delta) > (2 - \delta) \text{Var}_{\mathcal{R}}(L_\delta)$ ].

**Proof.** The inequality  $\text{Var}_{\mathcal{G}_i}(L) \leq 2 \text{Var}_{\mathcal{R}}(L)$  can be proved immediately by analyzing the algorithm  $\mathcal{A}_{1,m}$  (see Appendix A). Using the algorithm  $\mathcal{A}_{1,m}$ , the number of nonterminals of the grammar can grow in the part between the labels (L1) and (L2) only. The number of nonterminals increases to the value at most twice greater than the original one.

Let  $L_n$  be the language generated by the grammar  $\tilde{G}_n$ :

$$S \rightarrow a_1 A_1 b_1 | \cdots | a_n A_n b_n | A_1 c_1 | \cdots | A_n c_n,$$

$$A_i \rightarrow a_i A_i b_i | A_i c_i | a_i b_i | c_i \quad \text{for } 1 \leq i \leq n.$$

We shall prove that  $\text{Var}_{\mathcal{R}}(L_n) = n + 1$  and  $\text{Var}_{\mathcal{G}_i}(L_n) = 2n + 1$ , i.e., for any  $\delta > 0$  and for  $n > 1/\delta - 1$  the language  $L_n$  satisfies the inequality  $\text{Var}_{\mathcal{G}_i}(L_n) > (2 - \delta) \text{Var}_{\mathcal{R}}(L_n)$ .

Let  $G_n$  be a grammar generating the language  $L_n$ . According to Lemma 2.3,  $G_n$  contains nonterminals  $A_i$  such that  $A_i \Rightarrow^* \alpha A_i \beta$  and  $c_i \in \text{alph}(\alpha\beta)$ . Since for any  $u_0 \in L$  there is a fixed number  $i_0$  for which

$$\{i: a_i \in \text{alph}(u_0), b_i \in \text{alph}(u_0), c_i \in \text{alph}(u_0)\} = \{i_0\},$$

$A_i \neq A_j$  for  $i \neq j$  and every  $A_i$  differs from the initial nonterminal of  $G_n$ . So for  $G_n$  generating  $L_n$ ,  $\text{Var}(G_n) \geq n + 1$ . Because  $\text{Var}(\tilde{G}_n) = n + 1$  we have  $\text{Var}_{\mathcal{R}}(L_n) = n + 1$ .

For  $\tilde{G}_n = \mathcal{A}_{1,m}(\tilde{G}_n)$  being the grammar in weak Greibach normal form, we have  $\text{Var}_{\mathcal{G}_i}(L_n) \leq \text{Var}(\tilde{G}_n) = 2n + 1$ .

Let  $G'_n$  be a grammar in weak Greibach normal form generating the language  $L_n$ . For  $u \in L_n$ ,  $u = u_1 u_2$  and  $u_1 \in a_i^*$ ,  $u_2 \in (b_i \cup c_i)^*$ , the number of occurrences of  $a$  in  $u_1$  equals the number of occurrences of  $b_i$  in  $u_2$ . According to Lemma 2.3,  $G'_n$  contains nonterminals  $A_i$  and  $B_i$  iterating the letters  $a_i$  and  $c_i$ , respectively, i.e.,  $A_i \Rightarrow^* x \alpha A_i \beta$ , where  $x \in T$ ,  $a_i \in \text{alph}(x\alpha\beta)$  and  $B_i \Rightarrow^* y \gamma B_i \delta$ , where  $y \in T$  and  $c_i \in \text{alph}(y\gamma\delta)$ .

From the structure of the words in  $L_n$  we can see that  $x = a_i$  and  $L(\alpha) \subseteq a_i^*$ . Since  $c_i c_i' \in L_n$  we also have  $y = c_i$ .  $A_i \neq A_j$  for  $i \neq j$  since for  $u$  in  $L_n$ ,  $\{i: c_i \in \text{alph}(u) \text{ or } a_i \in \text{alph}(u)\} = \{i_0\}$ ;  $A_i \neq B_i$  for  $i = 1, 2, \dots, n$  because all occurrences of  $a_i$  in  $u \in L_n$  precede the occurrences of  $c_i$ 's and  $S \neq A_i$  and  $S \neq B_i$  for  $i = 1, 2, \dots, n$ .

We have proved that  $\text{Var}(G_n) \geq 2n + 1$  for any  $G'_n$  in weak Greibach normal form generating  $L_n$ . So we conclude that  $\text{Var}_{\mathcal{G}_i}(L_n) = 2n + 1$ .  $\square$

**Theorem 3.3.** For  $\varepsilon$ -free languages,

$$\text{Prod}_{\mathcal{G}_i}(\mathcal{L}_{\text{Prod}_\varepsilon}(n)) = O(n^3) \quad \text{and} \quad \text{Prod}_{\mathcal{G}_i}(\mathcal{L}_{\text{Prod}_\varepsilon}(n)) = \Omega(n^2).$$

**Proof.** Let  $G$  be a  $\text{Prod}_\varepsilon$ -minimal grammar generating a language  $L$  and let  $A = AR + b$  be its matrix expression. Assume that the number of words in the components

of  $R$  (i.e.,  $R$ -words) is  $r$  and similarly the number of  $b$ -words is  $b$ . Then  $\text{Prod}_{\mathcal{A}}(L) = r + b = n$ .

The grammar  $\mathcal{A}_2(G)$  is defined by the equations

$$A = bH + b, \quad (3.1)$$

$$H = R_0H + R_0, \quad (3.2)$$

where  $R_0$  is a matrix constructed from  $R$  by replacing the first nonterminals in  $R$ -words by corresponding words specified by (3.1).

The number of productions determined by (3.1) equals  $mb + b$  for  $m$  being the total number of nonterminals in  $G$ .

Let  $r_1 \leq r$  be a number of those  $R$ -words, which have nonterminals as their first symbol. Since for each nonterminal the number of  $A$ -words is at most  $2b$ , the number of  $R_0$ -words is at most  $2br_1 + r - r_1$  and the upper bound for the number of  $R_0H$ -words is  $m(2b + r - r_1)$ . Thus,

$$\text{Prod}_{\mathcal{A}_2}(L) \leq \text{Prod}(\mathcal{A}_2(G)) \leq (m + 1)(2br + r - r_1 + b).$$

Since for a  $\text{Prod}_{\mathcal{A}}$ -minimal grammar  $m \leq n/2$  and since  $2br + r - r_1 + b \leq n^2$ , there is a constant  $c$  such that  $\text{Prod}_{\mathcal{A}_2}(L) \leq cn^3$  for all  $L \in \mathcal{L}_{\text{Prod}}^{(n)}$ , i.e.,  $\text{Prod}_{\mathcal{A}_2}(\mathcal{L}_{\text{Prod}_{\mathcal{A}}}(n)) = O(n^3)$ .

It remains to prove that  $\text{Prod}_{\mathcal{A}_2}(\mathcal{L}_{\text{Prod}_{\mathcal{A}}}(n)) = \Omega(n^2)$ , i.e.,  $\sup\{\text{Prod}_{\mathcal{A}_2}(L) : \text{Prod}_{\mathcal{A}}(L) \leq n, L \text{ is context-free}\} \geq cn^2$  for some real number  $c$ . It is enough to prove for  $L_n = \{(b_1 \cup \dots \cup b_n)^{2^i} : 1 \leq i \leq n\}$  and  $n > 1$  that  $\text{Prod}_{\mathcal{A}}(L_n) \leq 2n$  and  $\text{Prod}_{\mathcal{A}_2}(L_n) \geq n^2$ .

The grammar with production rules

$$S \rightarrow A^2 | \dots | A^{2^n}, \quad A \rightarrow b_1 | \dots | b_n,$$

generates the language  $L_n$ , i.e.,  $\text{Prod}_{\mathcal{A}}(L_n) \leq 2n$ .

Let  $G'_n = (N, T, P, S)$  be a  $\text{Prod}_{\mathcal{A}_2}$ -minimal grammar for the language  $L_n$ . Let  $P_i = \{A \rightarrow b_i \omega : \omega \in (\{b_i\} \cup N)^*\}$  and let for  $i_0$ ,  $|P_{i_0}| \leq |P_i|$ ,  $i = 1, 2, \dots, n$ . The grammar  $G''_n = (N, T, P, S)$  generates the language  $L(G''_n) = \{b_i^{2^i} : 1 \leq i \leq n\}$ . According to [5, Theorem 6.3],  $\text{Prod}_{\mathcal{A}_2}(L(G''_n)) = n$ . Evidently, the  $P_i$ 's are pairwise disjoint for  $i = 1, 2, \dots, n$  and so  $\text{Prod}_{\mathcal{A}_2}(L_n) = \text{Prod}(G'_n) \geq n \text{Prod}_{\mathcal{A}_2}(L(G''_n)) = n^2$ .  $\square$

The results for the measure  $\text{Symb}$  given in Theorem 3.4 were previously discussed in [11] and they are also included in [8].

**Theorem 3.4.** For  $\varepsilon$ -free languages,

$$\text{Symb}_{\mathcal{A}_2}(\mathcal{L}_{\text{Symb}_{\mathcal{A}}}(n)) = O(n^3), \quad \text{Symb}_{\mathcal{A}_2}(\mathcal{L}_{\text{Symb}_{\mathcal{A}}}(n)) = \Omega(n^2).$$

**Proof.** The proof proceeds in the same way as for the measure  $\text{Prod}$  (see [8, pp. 129–131] or [11, pp. 348–349]).  $\square$

**Theorem 3.5.**  $\text{Lev}_{\mathcal{A}_2}(\mathcal{L}_{\text{Lev}_{\mathcal{A}}}(n))$  cannot be bounded by any total function  $\varphi(n)$ .

**Proof.** According to Lemma 2.1 it is sufficient to prove Lemma 3.6.  $\square$

**Lemma 3.6.** For any natural number  $n$  there is a language  $L_n$  such that

$$\text{Lev}_{\mathcal{A}}(L_n) = 1 \quad \text{and} \quad \text{Lev}_{\mathcal{G}_i}(L_n) = n + 1.$$

**Proof.** For  $n = 1$ ,  $L_1 = ab^*$  satisfies the conditions of Lemma 3.6.

For  $n > 1$ , let  $L_n$  be a language generated by the grammar  $G_n$  with starting nonterminal  $A_1$  and with production rules

$$A_i \rightarrow a_i A_i b_i \mid a_i A_{i+1} b_i \mid A_i c_i, \quad 1 \leq i \leq n-1,$$

$$A_n \rightarrow a_n A_n b_n \mid a_n A_1 b_n \mid A_n c_n \mid a_n b_n.$$

Obviously,  $\text{Lev}(G_n) = \text{Lev}_{\mathcal{A}}(L_n) = 1$ .

$\text{Lev}_{\mathcal{G}_i}(L_n) \leq n + 1$  since a grammar  $\hat{G}_n = \mathcal{A}_{1..n}(G_n)$  has  $n + 1$  grammatical levels.

Let  $G'_n$  be a grammar in  $\mathcal{G}'_i$  and  $L(G'_n) = L_n$ . According to Lemma 2.3, for any  $c_i$  there is a  $Z_i$  in  $G'_n$ , such that  $Z_i \Rightarrow^* a Z_i \beta$ ,  $c_i \in \text{alph}(\alpha\beta)$ . The number of occurrences of  $c_i$  in a word in  $L_n$  does not depend on the number of occurrences of the other symbols, i.e.,  $\alpha\beta$  can be chosen in such a manner that  $\text{alph}(\alpha\beta) = \{c_i\}$ , i.e.,  $Z_i \Rightarrow^* c_i^{k_i} Z_i c_i^{l_i}$ ,  $k_i \geq 1$  and  $l_i \geq 0$ .  $Z_i \neq Z_j$  for  $i \neq j$  since if  $ac_j\beta \in L_n$  then  $ac_i\beta \notin L_n$ .

An initial nonterminal  $S$  differs from all  $Z_i$  because  $c_i$  is prefix of no word in  $L$ .

To get  $\text{Lev}_{\mathcal{G}_i}(L_n) = n + 1$  it remains to prove that  $Z_1, \dots, Z_n$  and  $S$  are elements of different grammatical levels. Following Lemma 2.2, it is enough to prove that  $\text{alph}(L(Z_i)) = \{c_i\}$ .

From the structure of the language  $L_n$  it follows that:

(a) if  $uc'_i v \in L_n$  for  $t \geq 1$ , then  $u \in \Sigma^*(c_i \cup b_i)$ ,  $v \in (c_i \cup b_i \cup b_{i-1})\Sigma^*$  for  $i \geq 2$  and  $v \in \{\epsilon\} \cup (c_1 \cup b_1 \cup b_n)\Sigma^*$  for  $i = 1$ ;

(b) if  $uc_i v \in L_n$ , then  $v \in \Sigma^*(a_1 \cup \dots \cup a_n)\Sigma^*$ , i.e., all letters  $a_i$  precede the letters  $c_i$  in the word in  $L_n$ ;

(c) the number and positions of  $b_i$ 's in a word of language  $L_n$  uniquely determine the number and position of letters  $a_i$  in the word.

According to (a),  $\text{alph}(L(Z_i)) \subseteq \{c_i, b_i, b_{i-1}\}$  but since (c) and (b) must be also satisfied  $\text{alph}(L(Z_i)) = \{c_i\}$ .  $\square$

#### 4. Complexity of position restricted grammars

**Theorem 4.1.** For  $K$  being Var, Prod, Lev or Hei, and  $t$  being any position restricted type,  $K_t(\mathcal{F}_{K_n}(n))$  cannot be bounded by any function  $\varphi$ .

**Proof.** The theorem will be proved by proving the following lemma.

**Lemma 4.2.** Let  $K$  be Var, Prod, Lev or Hei. Let  $t$  be a class of position restricted grammars of type  $t = (m_1, \dots, m_k, m_{k+1})$ .

For any natural number  $n$  there is a language  $L_n$  such that  $K_{\mathcal{R}}(L_n) = 1$  and  $K_{\ell}(L_n) = n$ .

**Proof.** For the language  $L_n = \{a^{f(n)}\}$ , where  $f(n) = k^n + \sum_{j=1}^{n-1} k^j \sum_{j=1}^{k+1} m_j$ , evidently  $\text{Var}_{\mathcal{R}}(L_n) = \text{Prod}_{\mathcal{R}}(L_n) = 1$ .

Let  $G_n$  be a position restricted grammar of type  $t$  and let  $L(G_n) = L_n$ . The language  $L_n$  is finite so all nodes of any branch of a derivation tree of the word  $a^{f(n)}$  in  $G_n$  are denoted by distinct nonterminals. A position restricted grammar of type  $t$  has at most  $k$  nonterminals in the right-hand side of production rules and moreover  $\sum_{i=1}^{k+1} m_i$  is the total length of terminal words in it, so the derivation tree of the word  $a^{f(n)}$  in  $G_n$  has a branch of length  $l \geq n$ , i.e.,  $\text{Var}(G_n) \geq n$  and  $\text{Prod}(G_n) \geq n$  for any  $G_n \in \mathcal{L}$ .

Let  $G_n$  be a  $\text{Lev}_{\ell}$ -minimal grammar for the language  $L_n$ . Since  $L_n$  is finite,  $G_n$  contains single-letter levels only, i.e.,  $\text{Lev}(G_n) = \text{Var}(G_n) \geq n$ .

Finally,  $\text{Hei}_{\ell}(L_n) \geq n$  since  $G_n$  consists of trivial levels only (i.e., single-letter levels) and the derivation tree of the word  $a^{f(n)}$  in  $G_n$  has a branch of length at least  $n$ .

Moreover,  $K_{\ell}(L_n) = n$  for  $K$  being  $\text{Var}$ ,  $\text{Prod}$ ,  $\text{Lev}$  and  $\text{Hei}$  since for a grammar  $G_n$  with the production rules

$$A_i \rightarrow a^{m_1} A_{i+1} a^{m_2} A_{i+1} \cdots a^{m_k} A_{i+1} a^{m_{k+1}}, \quad 1 \leq i \leq n-1,$$

$$A_n \rightarrow a,$$

we have  $G_n \in \mathcal{L}$ ,  $L(G_n) = L_n$  and  $\text{Var}(G_n) = \text{Prod}(G_n) = \text{Lev}(G_n) = \text{Hei}(G_n) = n$ .  $\square$

The relationships between  $\text{Symb}_{\ell}$  and  $\text{Symb}_{\mathcal{R}}$  depend on the grammatical type  $t$ . We shall give a result only for the case  $t = (0, 0, 0)$ .

**Theorem 4.3.** *We have*

$$\text{Symb}_{(0,0,0)}(L) \leq 4 \text{Symb}_{\mathcal{R}}(L) - 9.$$

Also, there is a sequence of languages  $\{L_n\}_{n=1}^{\infty}$  such that  $\text{Symb}_{(0,0,0)}(L_n) = 4 \text{Symb}_{\mathcal{R}}(L_n) - 9$ .

**Proof.** Let  $G$  be a  $\text{Symb}_{\mathcal{R}}$ -minimal grammar for a language  $L$  (i.e.,  $G$  and  $\text{Symb}(G) = \text{Symb}_{\mathcal{R}}(L)$ ). Let  $\mathcal{A}_3$  be the algorithm, described in Appendix A, which transforms a completely reduced grammar to a grammar of type  $(0, 0, 0)$ . Suppose  $G$  has  $s$  productions (i.e.,  $P = \{p_i = A_i \rightarrow \alpha_i : i = 1, \dots, s\}$ ) and, moreover,  $\{j : |\alpha_j| \geq 2\} = \{j : 1 \leq j \leq r, r \leq s\}$ . Then

$$\text{Symb}_{\mathcal{R}}(L) = \text{Symb}(G) = \sum_{i=1}^r (|\alpha_i| + 2) + 3(s - r) = 3s - r + \sum_{i=1}^r |\alpha_i|,$$

$$\text{Symb}_{(0,0,0)}(L) \leq \text{Symb}(\mathcal{A}_3(G)) = \sum_{i=1}^r 4(|\alpha_i| - 1) + 3k + 3(s - r),$$

where  $k$  is the number of terminals, which occur as the last symbol on the right-hand sides of production rules  $p_1, \dots, p_n$  i.e.,  $k \leq \min(r, |\text{alph}(L)|)$ .

$$\begin{aligned} \text{Symb}_{(0,0,0)}(L) &\leq 4 \sum_{i=1}^r |\alpha_i| - 4r + 3 \min(r, |\text{alph}(L)|) + 3(s-r) \\ &\leq 4 \text{Symb}_{\mathcal{R}}(L) - 9s \leq 4 \text{Symb}_{\mathcal{R}}(L) - 9. \end{aligned}$$

We shall prove the second part of the theorem for  $L_n = \{a_1 a_2 \cdots a_n\}$ ,  $n \geq 1$ . Clearly,  $\text{Symb}_{\mathcal{R}}(L_n) = n + 2$ .

In a grammar of type  $t = (0, 0, 0)$  generating  $L_n$ , all vertices of a derivation tree for  $a_1 \cdots a_n$  are labelled by different nonterminals and terminals. To get  $n$  leaves of a binary tree we need at least  $n - 1$  additional vertices (except of leaves) and for a derivation tree of a grammar of type  $(0, 0, 0)$  we need at least one more vertex to get all leaves labelled by terminals. Therefore,

$$\text{Symb}_{(0,0,0)}(L_n) \geq 4(n-1) + 3 = 4 \text{Symb}_{\mathcal{R}}(L_n) - 9.$$

For a grammar  $G_n$  with the productions

$$A_i \rightarrow a_i A_{i+1}, \quad i = 1, 2, \dots, n-1, \quad A_n \rightarrow a_n,$$

$L(G_n) = \{a_1 \cdots a_n\}$ ,  $G_n$  is in  $(0, 0, 0)$  and  $\text{Symb}(G_n) = 4n - 1 = 4 \text{Symb}_{\mathcal{R}}(L_n) - 9$ .  $\square$

Now we shall give some remarks concerning Chomsky normal form, since it is stronger than position restricted type  $(0, 0, 0)$  in a sense that production rules of a form  $A \rightarrow aB$  are not allowed in Chomsky normal form.

**Remark 4.4.** In Theorem 4.1, position restricted grammars of type  $t$  can be replaced by grammars in Chomsky normal form. The proof of this new theorem can be done exactly in the same way as that of Theorem 4.1 (i.e., using the same language as for grammars of type  $(0, 0, 0)$ ).

**Remark 4.5.** If  $\mathcal{C}$  is the class of grammars in Chomsky normal form, then  $\text{Symb}_{\mathcal{C}}(L) \leq 7 \text{Symb}_{\mathcal{R}}(L) - 15$ .

There is a sequence of languages  $\{L_n\}_{n=1}^{\infty}$  such that  $\text{Symb}_{\mathcal{C}}(L_n) = 7 \text{Symb}_{\mathcal{R}}(L_n) - 15$ .

The proof can be done in the same way as the proof of Theorem 4.3. If  $r, s$  and  $A_i \rightarrow \alpha_i$  have the same meaning as in the proof of Theorem 4.3 we have

$$\text{Symb}_{\mathcal{C}}(L) = 3s - r + \sum_{i=1}^r |\alpha_i|.$$

Let  $G$  be a  $\text{Symb}_{\mathcal{R}}$ -minimal grammar for  $L$  and let algorithm  $\mathcal{A}_4$  be that from Appendix A. We have

$$\text{Symb}_{\mathcal{C}}(L) \leq \text{Symb}(\mathcal{A}_4(G)) \leq 3(s-r) + 4 \sum_{i=1}^r (|\alpha_i| - 1) + 3|\text{alph}(L)|.$$

Since  $s \geq 1$ ,  $r \geq 0$  and  $|\text{alph}(L)| \leq \text{Symb}_{\mathcal{A}}(L) - 2$ , we now have  $\text{Symb}_{\mathcal{A}}(L) \leq 7 \text{Symb}_{\mathcal{B}}(L) - 15$ .

The optimality of the result can be proved in the same way as in Theorem 4.3. For  $L_n = \{a_1 \cdots a_n\}$  one can get  $\text{Symb}_{\mathcal{B}}(L_n) = n + 2$  and  $\text{Symb}_{\mathcal{A}}(L_n) = 7n - 1$ .

**Open problem.** Let  $\mathcal{L}_k$  be the class of languages such that, for  $L \in \mathcal{L}_k$ ,  $|\text{alph}(L)| = k$ . Then by analyzing algorithm  $\mathcal{A}_4$  for  $L \in \mathcal{L}_k$  one can get  $\text{Symb}_{\mathcal{A}}(L) \leq 4 \text{Symb}_{\mathcal{B}}(L) - 3k - 9$ .

Improve the coefficient 4 in the estimation or prove its optimality.

**Remark 4.6.** In connection with the coefficient 4 we shall prove in Example 4.7 the optimality of that coefficient for the case of the sequence  $\{L_n\}_{n=1}^{\infty}$  such that  $\lim_{n \rightarrow \infty} |\text{alph}(L_n)| / \text{Symb}_{\mathcal{B}}(L_n) = 0$ .

**Example 4.7.** Let  $\Sigma_n = \{a_0, \dots, a_{n-1}\}$  be the alphabet of  $L_n$ . Let

$$g(i, j) = (\frac{1}{2}i(i+1) + j) \pmod{n}$$

and

$$L_n = \{a_{g(0,i)} a_{g(1,i)} \cdots a_{g(n-1,i)} : 0 \leq i \leq n-1\}.$$

The optimality will be proved in the sense that for any  $\delta > 0$  there is an  $n_0$  such that for all  $n \geq n_0$ ,

$$\text{Symb}_{\mathcal{A}}(L_n) > (4 - \delta) \text{Symb}_{\mathcal{B}}(L_n). \quad (4.1)$$

We shall prove  $\text{Symb}_{\mathcal{B}}(L_n) = n^2 + 2n$  and  $\text{Symb}_{\mathcal{A}}(L_n) = 4n^2 - n$ , i.e., inequality (4.1) holds for  $n > (4 - 2\delta) / \delta$ .

For a language  $L_n$  the following property can be proved immediately:

(P) All subwords of  $L_n$  of length at least 2 differ from each other.

We shall prove that a grammar  $G_n$  with productions

$$S \rightarrow a_{g(0,j)} \cdots a_{g(n-1,j)}, \quad 0 \leq j \leq n-1,$$

is a  $\text{Symb}_{\mathcal{B}}$ -minimal grammar for  $L_n$  and therefore  $\text{Symb}_{\mathcal{B}}(L_n) = n(n+2)$ .

Suppose that  $G'_n$  contains a nonterminal  $A \neq S$  such that either  $L(A) = \{u\}$  and  $A$  occurs at least twice on the right-hand side of rules in  $G'_n$  or  $\{u_1, u_2\} \subseteq L(A)$  for  $u_1 \neq u_2$ .

If  $L(A) = \{u\}$ , then  $|u| = 1$  because of property (P) and it leads to a contradiction with  $\text{Symb}_{\mathcal{B}}$ -minimality of  $G'_n$ .

For  $\{u_1, u_2\} \subseteq L(A)$  let  $u$ ; consider the derivations

$$S \Rightarrow^+ uAv \Rightarrow^+ uu_1v \quad \text{and} \quad S \Rightarrow^+ uAv \Rightarrow^+ uu_2v.$$

Since no two words in  $L_n$  have same prefix or suffix we have  $u = v = \varepsilon$ . So  $S \Rightarrow^+ A$  in  $G'_n$  and this is a contradiction since  $G'_n$  is completely reduced.

Let  $G'_n$  be a  $\text{Symb}_{\mathcal{A}}$ -minimal grammar for  $L_n$ . Because of property (P) and the finiteness of  $L_n$  a derivation tree of the word  $a_{g(0,i)} \cdots a_{g(n-1,i)}$  in  $G'_n$  has  $2n - 1$

vertices labelled by different nonterminals;  $n-1$  of them are left-hand sides of production rules of length 4 and the last  $n$  of them are on left-hand sides of production of type  $A \rightarrow a$ . According to property (P), the sets of internal nonterminals used in derivations of different words from  $L_n$  are disjoint. ( $X$  is an internal nonterminal if  $X \rightarrow AB$  for some  $A$  and  $B$ .) Therefore,  $\text{Symb}_{\mathcal{A}_4}(L_n) \geq 4(n-1)n + 3n$ . Since for  $\mathcal{A}_4(G_n)$  we have  $\text{Symb}(\mathcal{A}_4(G_n)) = 4n^2 - n$ , and we have completed our proof.

## Appendix A

We shall give the algorithms used in previous parts of the paper. For the transformation of a completely reduced grammar to an equivalent one in Greibach normal form we shall give two algorithms. The classical sequential one will be denoted by  $\mathcal{A}_1$  and the matrix algorithm denoted by  $\mathcal{A}_2$ . For the transformation of a completely reduced grammar to an equivalent position restricted grammar of type  $(0, 0, 0)$  we shall give algorithm  $\mathcal{A}_3$ , and the algorithm  $\mathcal{A}_4$  will transform completely reduced grammars to equivalent grammars in Chomsky normal form.

Let us assume for all algorithms that  $G = (N, T, P, S)$  is a completely reduced grammar,  $N = \{S = A_1, A_2, \dots, A_m\}$ ,  $T = \{t_1, t_2, \dots, t_l\}$  and  $P = \{p_1, p_2, \dots, p_s\}$ .

Algorithm  $\mathcal{A}_1$  works in three steps. Firstly, it modifies the grammar  $G$  so that if  $A_i \rightarrow A_j \alpha$  is its production rule then  $j > i$  (that part of algorithm is between labels (L1) and (L2)); then the algorithm produces a grammar in weak Greibach normal form (just before label (L3)) and in the last part of algorithm (beginning with (L3)) terminals on the right-hand side of productions of grammar which are not on the first position are changed to nonterminals to obtain a grammar in Greibach normal form.

We shall fix the following notations for  $\mathcal{A}_1$ . For  $i = 1, 2, \dots, m$ , we denote  $P_i = \{A_i \rightarrow \alpha : A_i \rightarrow \alpha \in P\}$ . For a grammar with no left recursion of nonterminals

$$P_{i,j} = \{A_i \rightarrow \alpha\beta : A_i \rightarrow A_j\beta \in P, A_j \rightarrow \alpha \in P\}$$

and

$$\bar{P}_{i,j} = \{A_i \rightarrow A_j \alpha : A_i \rightarrow A_j \alpha \in P\}$$

for  $i, j = 1, 2, \dots, m$ .

### Algorithm $\mathcal{A}_1$

procedure GREIBACH NORMAL FORM ( $P, P'$ )  
(algorithm starts with the set of productions  $P$ )

(L1) begin  $i := 1$ ;

```

while  $i \leq m$  do  $j := 1$ ;
    while  $j \leq i - 1$  do  $P := (P \cup P_{i,j}) - \bar{P}_{i,j}$ ;
         $j := j + 1$ ; end
     $R := P_i$ ;
     $P := P - P_i \cup \{A_i \rightarrow \beta \mid \beta Z_i: A_i \rightarrow \beta \in R, \beta \notin A_i(N \cup T)^*\}$ 
         $\cup \{Z_i \rightarrow \alpha \mid \alpha Z_i: A_i \rightarrow A_i, \alpha \in R\}$ ;
     $i := i + 1$ ; end
(L2)  $i := m - 1$ ;
    while  $i > 0$  do  $P := P \cup \{P_{i,j}: i < j \leq m\} - \{\bar{P}_{i,j}: i < j \leq m\}$ ;
         $i := i - 1$ ; end
     $i := 1$ ;
    while  $i \leq m$  do  $P := P \cup \{Z_i \rightarrow \beta \alpha: Z_i \rightarrow A_j, \alpha \in P, A_j \rightarrow \beta \in P\}$ 
         $- \left\{ Z_i \rightarrow \alpha: \alpha \in N \left( T \cup N \cup \bigcup_{i=1}^m Z_i \right)^* \right\}$ ;
         $i := i + 1$ ; end
(L3)  $P' := \{T_j \rightarrow t_j: 1 \leq j \leq 1\}$ 
     $\cup \{X \rightarrow tY_1Y_2 \cdots Y_r:$ 
         $X \rightarrow tX_1X_2 \cdots X_r \in P, \text{ where } Y_i = X_i \text{ for } X_i \in \{A_i, Z_i: 1 \leq i \leq m\}$ 
         $\text{and } Y_i = T_j \text{ for } X_i = t_j\}$ ;
    end

```

(The algorithm produces the set of production rules in  $P'$ )

We denote by  $\mathcal{A}_{1,m}$  the initial part of the algorithm  $\mathcal{A}_1$  ended just before the command labelled by (L3) with the set of productions in  $P$ .

**Proposition A.1.** *Let  $G = (N, T, P, S)$  be a grammar in completely reduced form. Let  $\mathcal{A}_1(G) = G'$  and  $\mathcal{A}_{1,m}(G) = G''$ .*

*Then  $G' \in \mathcal{G}_1$ ,  $G'' \in \mathcal{G}_1'$  and  $L(G) = L(G') = L(G'')$ .*

**Proof.** The proof can be found in [8, p. 113] or [1, pp. 156–158].  $\square$

### Algorithm $\mathcal{A}_2$

procedure GREIBACH NORMAL FORM-MATRIX ( $P, P'$ )  
(the algorithm starts with the set of productions in  $P$ )

**begin** express  $G$  by equations of form  $A = AR + b$ ;

construct a new system of equation  $A = bH + b$

$$H = RH + R;$$

(where  $H$  be an  $m \times m$  matrix of new nonterminals)

$P' :=$  "set of productions corresponding to  $A = bH + b$ ";

$P :=$  "set of productions corresponding to  $H = RH + R$ ";

$P' := P' \cup \{H_{ij} \rightarrow \alpha_1 \cdots \alpha_s \in P: \alpha_1 \in T\}$

$\cup \{H_{ij} \rightarrow \beta \alpha_2 \cdots \alpha_s: \alpha_1 \in N, \alpha_1 \rightarrow \beta \in P', H_{ij} \rightarrow \alpha_1 \alpha_2 \cdots \alpha_s \in P\};$

**end**

(The algorithm produces set of production rules in  $P'$ )

**Proposition A.2.** Let  $G = (N, T, P, S)$  be a completely reduced grammar and  $G' = \mathcal{A}_2(G)$ . Then  $G' \in \mathcal{G}'$  and  $L(G) = L(G')$ .

**Proof.** For the proof, see, e.g., [1, p. 162].

**Algorithm  $\mathcal{A}_3$**

procedure POSITIONAL RESTRICTED-TYPE (0, 0, 0) ( $P, P'$ )

(the algorithm starts with the set of productions in  $P$ )

**begin**  $P' := \emptyset$ ;

$i := 1$ ;

**while**  $i \leq s$  **do** if  $p_i$  has the form  $A \rightarrow a$  **then**  $P' := P' \cup \{p_i\}$  **else**  $p_i$  has a form  $A \rightarrow X_1 X_2 \cdots X_n, n \geq 2$ ,

**if**  $X_n \in T, X_n = a_r$  **then**  $P' := \{X_r' \rightarrow a_r\} \cup P'$

$\cup \{A \rightarrow X_1 A_1, A_1 \rightarrow X_1 A_2, \dots, A_{n-2} \rightarrow X_{n-1} X_r'\};$

**else**  $P' := P' \cup \{A \rightarrow X_1 A_1, A_1 \rightarrow X_2 A_2, \dots, A_{n-2} \rightarrow X_{n-1} X_n\};$

$i := i + 1$ ;

**end**

(The algorithm produces a set of production rules in  $P'$ )

**Remark.** The  $A_i$ 's in the algorithm  $\mathcal{A}_3$  are new nonterminals which differ with each other for all productions  $p_i$ .

**Proposition A.3.** Let  $G$  be a completely reduced grammar. Let  $G' = \mathcal{A}_3(G)$ . Then  $G'$  is a position restricted grammar of type (0, 0, 0) and  $L(G) = L(G')$ .

**Proof.** The proof is obvious. It can be done in the same way as that for Algorithm  $\mathcal{A}_4$ .  $\square$

#### Algorithm $\mathcal{A}_4$

procedure CHOMSKY NORMAL FORM ( $P, P'$ )  
(algorithm starts with the set of productions  $P$ )

begin  $P' := \emptyset$ ;

$i := 1$ ;

    while  $i \leq s$  do if  $p_i$  has a form  $A \rightarrow a$  then  $P' := P' \cup \{p_i\}$  else

$p_i$  has a form  $A \rightarrow X_1 X_2 \cdots X_n, n \geq 2$

$P' := P' \cup \{A \rightarrow X'_1 A_1, A_1 \rightarrow X'_2 A_2, \dots, A_{n-2} \rightarrow X'_{n-1} X'_n\}$ ,

        where  $X'_i = X_i$  if  $X_i \in N$  and  $X'_i = T_r$  if  $X_i = t_r$ ;

$i := i + 1$ ;

$P' := P' \cup \{T_r \rightarrow t_r : 1 \leq k \leq l\}$

end

(The algorithm produces the set of production rules in  $P'$ )

**Remark.** The  $A_i$ 's in the algorithm  $\mathcal{A}_4$  are new nonterminals which differ with each other for all productions  $p_i$ .

**Proposition A.4.** Let  $G = (N, T, P, S)$  be a completely reduced grammar and let  $G' = \mathcal{A}_4(G)$ . Then  $G'$  is in Chomsky normal form and  $L(G) = L(G')$ .

**Proof.** For the proof, see, e.g., [8, p. 104].  $\square$

#### References

- [1] A.V. Aho and J.D. Ullman, *The Theory of Parsing, Translation and Compiling, Vol. 1* (Prentice-Hall, Englewood Cliffs, NJ, 1972).
- [2] B.S. Baker, Non-context-free grammars generating context-free languages, *Inform. and Control* **24** (1974) 231-246.
- [3] M. Blattner and S. Ginsburg, *Canonical Forms of Context-free Grammars and Position Restricted Grammar Forms*, Lecture Notes in Computer Science **56** (Springer, Berlin, 1977) pp. 49-55.
- [4] M. Blattner and S. Ginsburg, Position restricted grammar forms and grammars, *Theoret. Comput. Sci.* **17** (1982) 1-27.
- [5] J. Gruska, Some classifications of context-free languages, *Inform. and Control* **14** (1969) 152-179.
- [6] J. Gruska, A note on  $\epsilon$ -rules in context-free grammars, *Kybernetika* **11** (1975) 26-31.
- [7] J. Gruska, *Descriptive Complexity (of Languages). A Short Survey*, Lecture Notes in Computer Science **45** (Springer, Berlin, 1976) pp. 65-80.
- [8] M.A. Harrison, *Introduction to Formal Language Theory* (Addison-Wesley, Reading, MA, 1978).
- [9] D.E. Knuth, Big omicron and big omega and big theta, *SIGACT News* **8** (1976) 18-24.

- [10] H.A. Maurer, A. Salomaa and D. Wood, A super normal form theorem for context-free grammars, *J. ACM* **30** (1983) 95–102.
- [11] A. Pirická-Kelemenová, *Greibach Normal Form Complexity*, Lecture Notes in Computer Science **32** (Springer, Berlin, 1975) pp. 344–350.
- [12] A. Salomaa, *Jewels of Formal Language Theory* (Computer Science Press, Rockville, MD, 1981).