ON *n***-PATH-CONTROLLED GRAMMARS**

Jiří Koutný

Doctoral Degree Programme (2), FIT BUT E-mail: ikoutny@fit.vutbr.cz

Supervised by: Alexander Meduna E-mail: meduna@fit.vutbr.cz

ABSTRACT

This paper discusses context-free grammars with some root-to-leaf paths in derivation trees restricted by control languages. It demonstrates that if these control languages are linear, then there are several families of generated languages depending on the common part of all restricted paths. The paper deals with the investigation of several properties of these families.

1 INTRODUCTION

Indisputably, the investigation of context-free grammars with restricted derivation trees represents an important trend in today's formal language theory (see [1], [3], [4], [5], [6], [8]). In essence, these grammas generate their languages just like ordinary context-free grammars do; in addition, however, their derivation trees have to satisfy some simple prescribed conditions. The present paper continues with the investigation of these grammars.

The authors in [5] have studied a new type of restriction in derivation: a derivation tree in a context-free grammar is accepted only if it contains a path described by a string generated by another context-free grammar. They take two context-free grammars, G and G', where G' generates a language over the total alphabet of G. A string w generated by G is accepted only if there is a derivation tree t of w with respect to G such that there exists a path in t which is marked by a string from L(G'). Based on this restriction, they have introduced *path-controlled grammars*, *PC* grammars for short, and they have found many properties of this formal model. As they have noticed in the final remarks, however, still many modifications remains unsolved.

The regular paths in a derivation tree of context-free grammars do not increase the generative power (see [4] and see [5], Prop. 2). In this paper, therefore, a control language is supposed to be linear (see page 597 in [7]). Here, we deal with a generalization of path-controlled grammars where the string *w* generated by *G* is accepted only if there is a derivation tree *t* of *w* with respect to *G* such that there exists $n \ge 0$ paths in *t* that are marked by the strings from linear language L(G'). Based on path-controlled grammars introduced in [5] and pumping lemma for linear languages (see page 120 in [5]), we establish several types of *n*-path-controlled grammars. We show that if the restricted paths satisfy the certain property, then there exist an infinite hierarchy of *n*-path-controlled grammars.

In conclusion, we formulate some open problems and suggest some new trends in the investigation of *n*-path-controlled grammars.

2 DEFINITIONS

This paper assumes that the reader is familiar with the theory of formal languages (see [7]) and the theory of regulated rewriting (see [2]). In this section, we introduce some terminology used in the sequel.

For an alphabet *V*, *V*^{*} denotes the letter monoid (generated by *V* under the operation of concatenation), ε is the unit of *V*^{*}, *V*⁺ = *V*^{*} - { ε }. For a word $x \in V^*$, |x| denotes its length.

A context-free grammar is a quadruple G = (V, T, P, S), where, as usual, V is a (finite) alphabet, $T \subseteq V$ is a terminal alphabet, P is a finite set of production rules of the form $A \to x$, where $A \in V - T$, $x \in V^*$, and $S \in V - T$ is the starting symbol. Let N = V - T denote a set of nonterminals. In the standard manner, we introduce the relations $\Rightarrow_G, \Rightarrow_G^i, \Rightarrow_G^+, \Rightarrow_G^*$ (see [7]).

Let $_G \triangle(x)$ denote a set of the derivation trees with frontier *x* with respect to the grammar *G*. Let $t \in _G \triangle(x)$ be a derivation tree. Let path(s) denote the word obtained by concatenating all symbols of the path *s* (in order from the top, i.e. from the root of *t* to a leaf of *t*).

Borrowing the notation from [5], we generalize path-controlled grammar and introduce a *n*-path-controlled grammar, *nPC* grammar for short. An *nPC*-grammar is a pair (G, G'), where G = (V, T, P, S) is a context-free grammar and G' = (V', V, P', S') is a linear grammar (see page 597 in [7]). The language generated by (G, G') is $L(G, G') = \{w \in L(G) | \text{ there is a set } C \text{ of } n \text{ different paths in } t \in_G \triangle(w) \text{ such that for all } p \in C \text{ it holds } path(p) \in L(G') \text{ and all } p \in C \text{ are divided in the common node of } t\}.$

Clearly, each two paths of a derivation tree contain at least one common node. Thus, for a $_nPC$ grammar (G, G'), there is some $m_C \in \mathbb{N}$ that denotes a number of common nodes for all $p \in C$. Hence, for each two $p_1, p_2 \in C$ it holds that $path(p_1) = rDs_1$, $path(p_2) = rDs_2$, where $r \in N^*$, $D \in N$, $s_1, s_2 \in N^*T$ and $|rD| = m_C$.

Consider the pumping lemma for linear languages (see page 120 in [5]). Hence, each path(p), where $p \in C$, such that $|path(p)| \ge k$, for some $k \ge 0$, can be written in the form path(p) = uvwxy and for each $m \ge 0$ it holds that $uv^m wx^m y \in L(G')$.

We distinguish five types of $_nPC$ grammars depending on the value of m_C in relation to pumping lemma for L(G'): The types are $_n^IPC$ if C satisfies $0 \le m_C \le |u|$, $_n^{II}PC$ if C satisfies $|u| < m_C \le |uv|$, $_n^{III}PC$ if C satisfies $|uv| < m_C \le |uvw|$, $_n^{IV}PC$ if C satisfies $|uvw| < m_C \le |uvwx|$, and $_n^VPC$ if C satisfies $|uvwx| < m_C \le |uvwxy|$, where uvwxy is the shortest path from C.

We denote the family of the languages generated by LIN, CF, PC, $_nPC$, $_n^{IPC}$, $_n^{II}PC$, $_n^{III}PC$, $_n^{IV}PC$, $_n^{V}PC$ grammars by LIN, CF, PC, n-PC, I-n-PC, II-n-PC, III-n-PC, IV-n-PC, V-n-PC, respectivelly.

3 RESULTS

Theorem 1. PC = 1-PC = I-1-PC = II-1-PC = III-1-PC = IV-1-PC = V-1-PC.

Proof. The equality clearly follows from the definitions of *PC*, $_nPC$, $_n^iPC$, for i = I, II, III, IV, V, grammars.

Theorem 2. If $L \in III$ -n-PC, for $n = card(C) \ge 0$, then there are $p,q \in \mathbb{N}$ such that each $z \in L$ with |z| > p can be written in the form $z = u_1v_1u_2v_2\dots u_{2n+2}v_{2n+2}u_{2n+3}$, such that $0 < |v_1v_2\dots v_{2n+2}| \le q$ and $u_1v_1^iu_2v_2^i\dots u_{2n+2}v_{2n+2}^iu_{2n+3} \in L$ for all $i \ge 1$.

Proof. Let (G, G') be a ${}_n^{III}PC$ -grammar, where G = (V, T, P, S) and G' = (V', V, P', S'). We deal directly only with grammar G—in the case of G', we deal only with the language L(G').

Consider $t \in_{(G,G')} \triangle(z)$. For each $path(s) = SA_1 \dots A_k a$ of t, where $s \in C$, $S, A_1, \dots, A_k \in N$, $a \in T$, consider the rules $A_i \rightarrow x_i A_{i+1} y_i$ used when passing from A_i to A_{i+1} on this path and the rule $A_k \rightarrow x_k a y_k$ used in the last step of the derivation in G corresponding to the path s.

Consider that any $x_i y_i$, i = 1, ..., k, contains a nonterminal *B* that do not belong on any path $s \in C$. Clearly, there is substring z' of z derived from *B*. Since *G* is context-free, it follows that if $|z'| \ge k_1$, for some $k_1 \ge 0$, then there are two substrings z'_1, z'_2 of z' that can be pumped and by pumping lemma for context-free languages, z'_1, z'_2 are bounded in length.

If L(G) is infinite, the string $path(s) \in L(G')$ is arbitrarily long. Thus, if path(s) = uvxyzwith $|uvxyz| \ge k_2$, for some $k_2 \ge 0$, then uvxyz satisfies $uv^ixy^iz \in L(G')$, for $i \ge 1$. Hence, the derivations starting from the symbols of v and y can be repeated in G. Since (G, G') is $\prod_{n} PC$ grammar, it follows that the derivations starting from the symbols of v in G are common for all $s \in C$ and the derivations starting from the symbols of y in G are unique for each $s \in C$.

Consider the derivations starting from *v* in *G*. This leads to pumping of two substrings v_1 , v_{2n+2} of *z*—one in the left-hand side, one in the right-hand side of common part of all $s \in C$.

Consider the derivations starting from y in G. This leads to pumping of two substrings of z—one in the left-hand side, one in the right-hand side of each $s \in C$. For each $s_{i+1} \in C$, denote this two substrings v_{2i+2} , v_{2i+3} , i = 0, 1, ..., n-1. Since (G, G') is ${}_{n}^{III}PC$ grammar, we obtain 2n pumped substrings of z.

By pumping lemma for context-free languages, the substrings $v_1, v_2, \ldots, v_{2n+2}$ are bounded in length. So, the total length of the 2n+2 pumped substrings of *z* is bounded by a constant *q*. \Box

Corollary 3. III-n-PC cannot count to 2n+3, but can count to 2n+2, e.g. $L = \{a^i b^i c^i d^i e^i f^i g^i | i \ge 1\} \notin III-2-PC$, but $L \in III-3-PC$.

Corollary 4. There is an infinite hierarchy of $\bigcup_{i=0}^{n}$ III-i-PC languages; i.e. $\bigcup_{i=0}^{n}$ III-i-PC $\subset \bigcup_{i=0}^{n+1}$ III-i-PC, for $n \ge 0$, is proper.

Corollary 5. III-n-PC is not closed under concatenation, e.g. $L = \{a^i a^i a^i a^i a^i a^i a^i a^i | i \ge 1\} \in$ **III-2-PC**, but $LL \notin$ **III-2-PC**.

Example 1. Let us have ${}_{n}^{III}PC$ grammar (G, G'), $n \ge 0$, where

$$G_{1} = (\{S\} \cup \{A_{i}, B_{i} | i = 1, ..., n\} \cup \{a_{i} | i = 1, ..., 2n + 2\}, \{a_{i} | i = 1, ..., 2n + 2\}, P, S),$$

$$P = \{S \rightarrow a_{1}Sa_{2n+2}, S \rightarrow a_{1}A_{1} ... A_{n}a_{2n+n}\} \cup \{A_{i+1} \rightarrow a_{2i+2}A_{i+1}a_{2i+3}, A_{i+1} \rightarrow B_{i+1}, B_{i+1} \rightarrow a_{2i+2}a_{2i+3} | i = 0, ..., n-1\},$$

$$L(G') = \bigcup_{i=1}^{n} \{S^{k}A_{i}^{k}B_{i}a_{2i} | k \ge 1\} - \text{clearly}, L(G') \in \textbf{LIN}.$$

Consider a derivation in (G, G'):

$$\begin{split} S &\Rightarrow^{k} a_{1}^{k} S a_{2n+2}^{k} \\ &\Rightarrow a_{1}^{k} a_{1} A_{1} \dots A_{n} a_{2n+2} a_{2n+2}^{k} \\ &\Rightarrow^{n \times k} a^{k+1} a_{2}^{k} B_{1} a_{3}^{k} \dots a_{2n}^{k} B_{n} a_{2n+1}^{k} a_{2n+2}^{k+1} \\ &\Rightarrow^{n} a^{k+1} a_{2}^{k+1} a_{3}^{k+1} \dots a_{2n}^{k+1} a_{2n+1}^{k+1} a_{2n+2}^{k+1} \end{split}$$

Clearly, *n* different paths are described by L(G'). This way, by ${}_n^{III}PC$ grammar (G,G'), $n \ge 0$, we can generate the language $L(G,G') = \{a_1^k, \ldots, a_{2n+2}^k | k \ge 1\}$.

To be more concrete, consider $_{2}^{III}PC$ grammar (G, G'), where

$$G = (\{S, X, Y, U, V, a, b, c, d, e, f\}, \{a, b, c, d, e, f\}, P, S),$$

$$P = \{S \rightarrow aSf, S \rightarrow aXYf, X \rightarrow bXc, X \rightarrow U, U \rightarrow bc, Y \rightarrow dYe, Y \rightarrow V, V \rightarrow de\},$$

$$L(G') = \{S^n X^n U b \cup S^n Y^n V d | n \ge 1\},$$

$$L(G, G') = \{a^n b^n c^n d^n e^n f^n | n \ge 1\}.$$

Example 2. Let $m \ge 0$ with $m \mod 2 = 0$. Let us have $\prod_{n=1}^{HI} PC$ grammar $(G, G'), n \ge 0$, where

$$\begin{split} G &= (\{A_j, B_j, a_j | j = 1, \dots, m\} \cup \{C\}, \{a_j | j = 1, \dots, m\}, P, A_1), \\ P &= \{A_1 \to a_1 A_1, A_1 \to a_1 A_2, B_1 \to B_1 a_1, B_1 \to C, C \to a_1, A_m \to A_m a_m, A_m \to \{B_m\}^n\} \cup \\ \{A_i \to A_i a_i, A_i \to A_{i+1} | i = 2, \dots, m-1 \text{ with } i \mod 2 = 0\} \cup \\ \{A_i \to a_i A_i, A_i \to A_{i+1} | i = 3, \dots, m-1 \text{ with } i \mod 2 = 1\} \cup \\ \{B_i \to a_i B_i, B_i \to B_{i-1} | i = 2, \dots, m \text{ with } i \mod 2 = 0\} \cup \\ \{B_i \to B_i a_i, B_i \to B_{i-1} | i = 3, \dots, m \text{ with } i \mod 2 = 1\}, \\ L(G') &= \{A_1^{k_1} A_2^{k_2} \dots A_m^{k_m} B_m^{k_m} B_{m-1}^{k_{m-1}} \dots B_2^{k_2} B_1^{k_1} Ca_1 | k_i \ge 0, i = 1, \dots, m\} \text{--clearly, } L(G') \in \text{LIN} \end{split}$$

Consider a derivation in (G, G'), for some $m \ge 0$ with $m \mod 2 = 0$:

$$\begin{split} A_{1} &\Rightarrow^{k_{1}} a_{1}^{k_{1}} A_{2} \\ &\Rightarrow^{k_{2}} a_{1}^{k_{1}+1} A_{2} a_{2}^{k_{2}} \\ &\Rightarrow^{k_{2}} a_{1}^{k_{1}+1} A_{3} a_{2}^{k_{2}} \\ &\Rightarrow^{k_{1}} a_{1}^{k_{1}+1} A_{3} a_{2}^{k_{2}} \\ &\Rightarrow^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} A_{m} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{B_{m}\}^{n} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n \times k_{m}} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} B_{m}\}^{n} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} B_{m-1}\}^{n} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n \times k_{m-1}} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} B_{m-1} a_{m-1}^{k_{m-1}}\}^{n} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n \times k_{m-1}} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} B_{m-2} \dots a_{2}^{k_{2}} B_{1} a_{1}^{k_{1}} \dots a_{m-3}^{k_{m-3}} a_{m-1}^{k_{m-1}}\}^{n} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} a_{m-2}^{k_{2}} \dots a_{2}^{k_{2}} Ca_{1}^{k_{1}} \dots a_{m-3}^{k_{m-1}} n_{m-1}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} a_{m-2}^{k_{m-2}} \dots a_{2}^{k_{2}} Ca_{1}^{k_{1}} \dots a_{m-3}^{k_{m-1}} n_{m-1}^{k_{m}} a_{m}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} a_{m-2}^{k_{2}} \dots a_{2}^{k_{2}} a_{1}^{k_{1}} \dots a_{m-3}^{k_{m-1}} n_{m-1}^{k_{m}} \dots a_{6}^{k_{6}} a_{4}^{k_{4}} a_{2}^{k_{2}} \\ &\Rightarrow^{n} a_{1}^{k_{1}+1} a_{3}^{k_{3}} a_{5}^{k_{5}} \dots a_{m-1}^{k_{m-1}} \{a_{m}^{k_{m}} a_{m-2}^{k_{2}} \dots a_{2}^{k_{2}} a_{1}^{k_{1}} \dots a_{m-3}^{k_{m-1}} a_{m-1}^{k_{m}} n_{m} \dots a_{6}^$$

Thus, *n* different paths are described by L(G'). This way, by ${}_{n}^{III}PC$ grammar (G,G'), $n \ge 0$, $m \ge 0$ with $m \mod 2 = 0$, we can generate the language

$$L(G,G') = \{ (a_1^{k_1+1}a_3^{k_3}\dots a_{m-1}^{k_{m-1}}a_m^m a_{m-2}^{k_{m-2}}a_{m-4}^{k_{m-4}}\dots a_2^{k_2})^{n+1} | k_i \ge 0, i = 1,\dots,m \}.$$

To be more concrete, consider m = 4 and ${}_{3}^{III}PC$ grammar (G, G'), where

$$\begin{split} G &= (\{A, B, C, D, E, F, G, H, I, a, b, c, d\}, \{a, b, c, d\}, P, A), \\ P &= \{A \rightarrow aA, A \rightarrow aB, B \rightarrow Bb, B \rightarrow C, C \rightarrow cC, C \rightarrow D, D \rightarrow Dd, D \rightarrow HHH, E \rightarrow Ea, \\ E \rightarrow I, F \rightarrow bF, F \rightarrow E, G \rightarrow Gc, G \rightarrow F, H \rightarrow dH, H \rightarrow G, I \rightarrow a\}, \\ L(G') &= \{A^r B^s C^t D^u H^u G^t F^s E^r Ia | r, s, t, u \ge 0\}, \\ L(G, G') &= \{a^r c^t d^u b^s a^r c^t d^u b^s a^r c^t d^u b^s | r > 0, s, t, u \ge 0\}. \end{split}$$

4 CONCLUSION

We have considered a new type of restriction in derivation: ${}_{n}PC$ grammars as a generalization of *PC* grammars introduced in [5]. In relation to pumping lemma for linear languages, we have demonstrated that there are several types of such ${}_{n}PC$ grammars. We have found several properties of such model, especially pumping property for ${}_{n}^{III}PC$ grammars and some consequences that follow this property.

It seems that the most useful is the family of **III-n-PC** languages and for this family, there are still many questions to by answered—for instance, generative power, closure properties, decidability properties, parsing properties, and descriptional complexity. For **I-n-PC** and **V-n-PC**, however, also many properties can be found. For **II-n-PC** and **IV-n-PC**, there are currently no answers, only questions. As the requirements for a set *C* of all controlled paths looks fairly restrictive, weaker control could be of interest, i.e. the request for *C* that all controlled paths have to be divided in just one node can be weakened. The formal study of such variants remains to be carried out; we hope to return to this topic in a forthcoming paper.

ACKNOWLEDGEMENT

This work was partially supported by the FRVŠ MŠMT grant FR2581/2010/G1, the BUT FIT grant FIT-10-S-2, and the research plan MSM0021630528.

REFERENCES

- [1] K. Čulik and H. A. Maurer. Tree controlled grammars. *Computing*, 19:129–139, 1977.
- [2] J. Dassow and Gh. Păun. *Regulated Rewriting in Formal Language Theory*. Springer, Berlin, 1989.
- [3] J. Dassow and B. Truthe. Subregularly tree controlled grammars and languages. In Automata and Fromal Languages - 12th International Conference AFL 2008, Balatonfured, pages 158–169. Computer and Automation Research Institute of the Hungarian Academy of Sciences, 2008.
- [4] J. Koutný. Regular paths in derivation trees of context-free grammars. In Proceedings of the 15th Conference and Competition STUDENT EEICT 2009 Volume 4, pages 410–414. Faculty of Information Technology BUT, 2009.
- [5] S. Marcus, C. Martín-Vide, V. Mitrana, and Gh. Păun. A new-old class of linguistically motivated regulated grammars. In *CLIN*, pages 111–125, 2000.
- [6] C. Martín-Vide and V. Mitrana. Further properties of path-controlled grammars. In *Formal Grammar / Mathematics of Language 2005*, pages 219–230. Edimburgh, 2005.
- [7] A. Meduna. Automata and Languages: Theory and Applications. Springer Verlag, 2005.
- [8] Gh. Păun. On the generative capacity of tree controlled grammars. *Computing*, 21(3):213–220, 1979.