

Generating All Permutations by Context-Free Grammars in Chomsky Normal Form

Peter R.J. Asveld

*Department of Computer Science, Twente University of Technology
P.O. Box 217, 7500 AE Enschede, the Netherlands*

Abstract

Let L_n be the finite language of all $n!$ strings that are permutations of n different symbols ($n \geq 1$). We consider context-free grammars G_n in Chomsky normal form that generate L_n . In particular we study a few families $\{G_n\}_{n \geq 1}$, satisfying $L(G_n) = L_n$ for $n \geq 1$, with respect to their descriptonal complexity, i.e. we determine the number of nonterminal symbols and the number of production rules of G_n as functions of n .

Key words: context-free grammar, Chomsky normal form, permutation, descriptonal complexity, unambiguous grammar.

1 Introduction

The set L_n of all permutations of n different symbols consists of $n!$ elements [9,14]. So being a finite language, L_n can be trivially generated by a context-free grammar with a single nonterminal symbol and $n!$ productions. However, this is no longer true when we require that L_n is generated by a context-free grammar G_n in Chomsky normal form.

In this paper we investigate a few families $\{G_n\}_{n \geq 1}$ of context-free grammars in Chomsky normal form that generate $\{L_n\}_{n \geq 1}$. In particular we are interested in the grammatical or descriptonal complexity of these families. As complexity measures we use the number of nonterminal symbols and the number of production rules of G_n , both considered as functions of n . These measures have been used frequently in investigating context-free grammars; cf. e.g. [10,12,13,6,4,1,5].

This paper is organized as follows. After preliminaries on notation and terminology (Section 2) we consider some elementary properties of grammars G_n in

Email address: `infprja@cs.utwente.nl` (Peter R.J. Asveld).

Chomsky normal form that generate L_n (Section 3). In Section 4 we consider a straightforward approach based on the power set of the terminal alphabet Σ_n of G_n . Looking at regular (i.e. right-linear) grammars to generate $\{L_n\}_{n \geq 1}$ gives rise to a family of context-free grammars in Chomsky normal form in Section 5 with less productions than the ones in Section 4 (provided $n \geq 3$).

The families $\{G_n\}_{n \geq 1}$ studied in Sections 6 and 7 are obtained in a different way; viz. we exhibit G_1 and G_2 explicitly and then we proceed inductively by means of a grammatical transformation to obtain G_{n+1} from G_n ($n \geq 2$). Section 8 is devoted to a divide-and-conquer approach; although it leads to “concise” grammars, determining their descriptive complexity is less straightforward. Finally, Section 9 consists of some concluding remarks.

The present paper has been inspired by G. Satta who conjectured in 2002 [16] that “any context-free grammar G_n in Chomsky normal form that generates L_n must have a number of nonterminal symbols that is not bounded by any polynomial in n ”. Recently, this statement has been proved by K. Ellul, B. Krawetz, J. Shallit and M.-w. Wang in [7]. However, in [7] it is not shown how to generate the languages $\{L_n\}_{n \geq 1}$ by context-free grammars $\{G_n\}_{n \geq 1}$ in Chomsky normal form. The present paper provides some straightforward approaches to obtain a few such families $\{G_n^i\}_{n \geq 1}$ ($1 \leq i \leq 7$). None of these approaches is surprising but their relative descriptive complexity (expressed in terms of the number of nonterminal symbols and of the number of productions) is by no means obvious; cf. Section 9. In this way the paper is a taxonomy of basic grammar families for $\{L_n\}_{n \geq 1}$ and it might serve as a starting point for more involved approaches as well as for the quest for optimal grammars, i.e. grammars that are minimal with respect to these or other descriptive complexity measures.

2 Preliminaries

For each set X , let $\mathcal{P}(X)$ denote the power set of X , and $\mathcal{P}_+(X)$ the set of nonempty subsets of X , i.e. $\mathcal{P}_+(X) = \mathcal{P}(X) - \{\emptyset\}$. For each finite set X , $\#X$ denotes the cardinality (i.e. the number of elements) of X .

For background and elementary results on discrete mathematics, particularly on combinatorics (counting, recurrence relations or difference equations), we refer to texts like [9,14,15]. In order to save space we often use $C(n, k)$ to denote the binomial coefficient $C(n, k) = n!/(k!(n - k)!)$; in displayed formulas we apply the usual notation.

We assume familiarity with basic concepts, terminology and notation from formal language theory; cf. e.g. [11]. We will denote the empty word by λ . Recall that a λ -free context-free grammar $G = (V, \Sigma, P, S)$ is in *Chomsky normal form* if $P \subseteq N \times (N - \{S\})^2 \cup N \times \Sigma$ where $N = V - \Sigma$. For each

context-free grammar $G = (V, \Sigma, P, S)$, let $L(G, A)$ be the language over Σ defined by $L(G, A) = \{w \in \Sigma^* \mid A \Rightarrow^* w\}$. Then for the language $L(G)$ generated by G , we have $L(G) = L(G, S)$. Note that, if G is in Chomsky normal form, then $L(G, A)$ is a nonempty language for each A in N .

Henceforth, we use $\Sigma_n = \{a_1, a_2, \dots, a_n\}$ to denote an alphabet of n different symbols ($n \geq 1$). As mentioned earlier, L_n is the finite language over Σ_n that consists of the $n!$ permutations of a_1, a_2, \dots, a_n . Since L_n is finite, we have that each context-free grammar G_n in Chomsky normal form that generates L_n , possesses the property that each nonterminal symbol of G_n is not recursive.

The length of word w will be denoted by $|w|$, as usual. For each word w over Σ_n , $\mathcal{A}(w)$ is the set of all symbols from Σ_n that really do occur in w . Formally, $\mathcal{A}(\lambda) = \emptyset$, and $\mathcal{A}(ax) = \{a\} \cup \mathcal{A}(x)$ for each $a \in \Sigma_n$ and $x \in \Sigma_n^*$. This mapping is extended to languages L over Σ_n by $\mathcal{A}(L) = \bigcup \{\mathcal{A}(w) \mid w \in L\}$.

In the sequel we often restrict ourselves to context-free grammars $G_n = (V_n, \Sigma_n, P_n, S_n)$ in Chomsky normal form with the following property: if $A \rightarrow BC$ is a production in P_n , then so is $A \rightarrow CB$, and we abbreviate $A \rightarrow BC \mid CB$ by $A \twoheadrightarrow BC$. The underlying rationale is, of course, that we want to keep the number of nonterminal symbols as low as possible. However, the reader should always realize that $A \twoheadrightarrow BC$ counts for two productions.

3 Elementary Properties

In this section we discuss some straightforward properties of context-free grammars in Chomsky normal form that generate L_n . Examples of these properties will be given at appropriate places in subsequent sections. Throughout this section $G_n = (V_n, \Sigma_n, P_n, S_n)$ is a context-free grammar in Chomsky normal form that generates L_n and N_n is defined by $N_n = V_n - \Sigma_n$.

For each word w over Σ_n in $L(G_n, A)$, $D(A, w)$ denotes a derivation tree for w from A according to the rules of G_n .

Proposition 3.1. (1) *For each nonterminal A in N_n , the language $L(G_n, A)$ is a nonempty subset of an isomorphic copy M_k of the language L_k for some k ($1 \leq k \leq n$). Consequently, each string z in $L(G_n, A)$ has length k , z consists of k different symbols, and $\mathcal{A}(z) = \mathcal{A}(L(G_n, A)) = \mathcal{A}(M_k)$.*

(2) *Let A and B be nonterminal symbols in N_n . If $L(G_n, A) \cap L(G_n, B) \neq \emptyset$, then $\mathcal{A}(L(G_n, A)) = \mathcal{A}(L(G_n, B))$.*

Proof. (1) Let w be a word in $L(G_n)$ with derivation tree $D(S_n, w)$ in which the nonterminal symbol A occurs. Consider the subtree $D(A, x)$ of $D(S_n, w)$, rooted by the nonterminal A , the leaves of which constitute a substring x of w ; so there exist words u and v with $w = uxv$. If $|x| = k$ for some k ($1 \leq k \leq n$), then $\mathcal{A}(x)$ has precisely k elements, since w is a permutation in L_n .

Suppose that $L(G_n, A)$ contains a string y with $|y| \neq k$: thus there is a derivation tree $D(A, y)$ according to G_n for y . Replacing $D(A, x)$ by $D(A, y)$ in $D(S_n, w)$ yields a derivation of uyv with $|uyv| \neq n$ and $uyv \notin L_n$. Hence each word in $L(G_n, A)$ has length k .

By a similar argument we can conclude that $L(G_n, A)$ is a language over the alphabet $\mathcal{A}(x)$ with the property that for each word z in $L(G_n, A)$, we have $\mathcal{A}(z) = \mathcal{A}(x)$. Consequently, $L(G_n, A)$ is a subset of an isomorphic copy M_k of L_k , i.e. $L(G_n, A) \subseteq M_k$.

(2) Suppose $L(G_n, A) \cap L(G_n, B) \neq \emptyset$: so it contains a word of length k for some $k \geq 1$. Then by Proposition 3.1(1), we have that both $L(G_n, A)$ and $L(G_n, B)$ are subsets of the same isomorphic copy M_k of L_k . Consequently, $\mathcal{A}(L(G_n, A)) = \mathcal{A}(L(G_n, B)) = \mathcal{A}(M_k)$. \square

This result gives rise to an equivalence relation on N_n ; viz.

Definition 3.2. Two nonterminal symbols A and B from N_n are called *equivalent* if $|x| = |y|$ for some $x \in L(G_n, A)$ and some $y \in L(G_n, B)$. The corresponding equivalence classes are $\{E_{n,k}\}_{k=1}^n$. The number of elements $\#E_{n,k}$ of the equivalence class $E_{n,k}$ will be denoted by $D(n, k)$ ($1 \leq k \leq n$). \square

Next we consider the effect of a single rewriting step with respect to the equivalence classes $\{E_{n,k}\}_{k=1}^n$.

Proposition 3.3. (1) *If $A \rightarrow BC$ is a rule in G_n , then $\mathcal{A}(L(G_n, B)) \cap \mathcal{A}(L(G_n, C)) = \emptyset$ and $\mathcal{A}(L(G_n, B)) \cup \mathcal{A}(L(G_n, C)) = \mathcal{A}(L(G_n, A))$.*

(2) *If $A \rightarrow BC$ is a rule in G_n with $A \in E_{n,k}$, $B \in E_{n,i}$ and $C \in E_{n,j}$, then $i + j = k$. Consequently, $1 \leq i < k$ and $1 \leq j < k$.*

Proof. (1) Suppose that the intersection is nonempty: if it contains a symbol a , then we have a subderivation $A \Rightarrow BC \Rightarrow^* x_1 a x_2 a x_3$ which cannot be a subderivation of a derivation that yields a permutation.

The inclusion $\mathcal{A}(L(G_n, B)) \cup \mathcal{A}(L(G_n, C)) \subseteq \mathcal{A}(L(G_n, A))$ is obvious. Suppose that this inclusion is proper; so there exists a symbol a with $a \in \mathcal{A}(L(G_n, A)) - (\mathcal{A}(L(G_n, B)) \cup \mathcal{A}(L(G_n, C)))$. Clearly, there is a rule $A \rightarrow DE$ with $a \in \mathcal{A}(L(G_n, D)) \cup \mathcal{A}(L(G_n, E))$. Consider the derivation $S_n \Rightarrow^* uAv \Rightarrow uBCv \Rightarrow^* u xv$ with $a \in \mathcal{A}(uv)$ and $a \notin \mathcal{A}(x)$, yielding the permutation $u xv$. Using this alternative rule $A \rightarrow DE$ for A we obtain the derivation $S_n \Rightarrow^* uAv \Rightarrow uDEv \Rightarrow^* u y v$ with $a \in \mathcal{A}(y)$; hence $u y v$ is not a permutation. Consequently, the inclusion cannot be proper; hence we have equality.

(2) follows from Propositions 3.1 and 3.3(1). \square

By Propositions 3.1 and 3.3 the set N_n inherits a partial order from the power

set $\mathcal{P}(\Sigma_n)$ of the alphabet Σ_n . This partial order, induced by the inclusion relation on $\mathcal{P}(\Sigma_n)$, is a more general notion than the linear order present in the concept of *sequential grammar*; cf. [8,3].

We will now define this partial order relation formally as follows.

Definition 3.4. Let A and B be nonterminal symbols from N_n . Then the partial order \sqsubseteq on N_n and the corresponding strict order \sqsubset are defined by:

$$A \sqsubseteq B \text{ if and only if } \mathcal{A}(L(G_n, A)) \subseteq \mathcal{A}(L(G_n, B)),$$

$$A \sqsubset B \text{ if and only if } \mathcal{A}(L(G_n, A)) \subset \mathcal{A}(L(G_n, B)). \quad \square$$

As complexity measures of a context-free grammar G_n we use the number $\nu(n)$ of nonterminal symbols and the number $\pi(n)$ of productions of G_n ; so $\nu(n) = \#N_n$ and $\pi(n) = \#P_n$. As the notation suggests, we will view both ν and π as functions of n . For a more general and thorough treatment of descriptonal complexity issues in relation to context-free grammars and their languages we refer to [10,12,13,6,4,1,5].

4 A Simple Approach

In view of Section 3 a straightforward way to generate L_n is to define G_n in terms of subsets of Σ_n : to each X of $\mathcal{P}_+(\Sigma_n)$ we associate a nonterminal A_X that generates all permutations over X , i.e. if $\#X = k$ ($1 \leq k \leq n$), then $L(G_n, A_X) \subset X^k$ and $L(G_n, A_X)$ is an isomorphic copy of L_k .

Definition 4.1. The family $\{G_n^1\}_{n \geq 1}$ is given by $\{(V_n, \Sigma_n, P_n, S_n)\}_{n \geq 1}$ with

- $N_n = V_n - \Sigma_n = \{A_X \mid X \in \mathcal{P}_+(\Sigma_n)\}$,
- $P_n = \{A_{\{a\}} \rightarrow a \mid a \in \Sigma_n\} \cup \{A_{X \cup Y} \rightarrow A_X A_Y \mid X, Y \in \mathcal{P}_+(\Sigma_n), X \cap Y = \emptyset\}$,
- $S_n = A_{\Sigma_n}$. □

Clearly, $A_X \sqsubset A_Y$ [$A_X \sqsubseteq A_Y$, respectively] holds if and only if $X \subset Y$ [$X \subseteq Y$] for all X and Y in $\mathcal{P}_+(\Sigma_n)$.

Example 4.2. We consider the case $n = 3$ in detail; instead of subsets of Σ_3 , we use subsets of $\{1, 2, 3\}$ as indices of nonterminals. Then we have $G_3^1 = (V_3, \Sigma_3, P_3, S_3)$ with $S_3 = A_{123}$, $N_3 = \{A_{123}, A_{12}, A_{13}, A_{23}, A_1, A_2, A_3\}$ and $P_3 = \{A_{123} \rightarrow A_{12}A_3 \mid A_{13}A_2 \mid A_{23}A_1, A_{12} \rightarrow A_1A_2, A_{13} \rightarrow A_1A_3, A_{23} \rightarrow A_2A_3, A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3\}$.

Now $E_{3,3} = \{A_{123}\}$, $E_{3,2} = \{A_{12}, A_{13}, A_{23}\}$, $E_{3,1} = \{A_1, A_2, A_3\}$, $D(3,3) = 1$, $D(3,2) = D(3,1) = 3$, $\nu_1(3) = 7$ and $\pi_1(3) = 15$. □

Proposition 4.3. For the family $\{G_n^1\}_{n \geq 1}$ of Definition 4.1 we have

- (1) $D(n, k) = C(n, k)$ with $1 \leq k \leq n$,

$$(2) \nu_1(n) = 2^n - 1,$$

$$(3) \pi_1(n) = 3^n - 2^{n+1} + n + 1.$$

Proof. (1) and (2) follow from Definition 4.1 and $\nu_1(n) = \sum_{k=1}^n D(n, k) = \sum_{k=1}^n C(n, k) = 2^n - 1$ [9]. By the definition of N_n and P_n , we have $\pi_1(n) = n + h(n)$ where $h(n) = \#\{A_{X \cup Y} \rightarrow A_X A_Y \mid X, Y \in \mathcal{P}_+(\Sigma_n), X \cap Y = \emptyset\}$.

If the set $X \cup Y$ possesses k elements ($k \geq 2$), then the set $\{A_{X \cup Y} \rightarrow A_X A_Y \mid X, Y \in \mathcal{P}_+(\Sigma_n), X \cap Y = \emptyset\}$ contains $2^k - 2$ elements, because both X and Y are nonempty. Then

$$\begin{aligned} h(n) &= \sum_{k=2}^n \binom{n}{k} (2^k - 2) = \sum_{k=1}^n \binom{n}{k} (2^k - 2) = \\ &= \sum_{k=1}^n \binom{n}{k} 2^k - 2 \cdot \sum_{k=1}^n \binom{n}{k} = \sum_{k=1}^n \binom{n}{k} 2^k 1^{n-k} - 2 \cdot \sum_{k=1}^n \binom{n}{k} = \\ &= \sum_{k=0}^n \binom{n}{k} 2^k 1^{n-k} - 2^0 1^n - 2 \cdot \left(\sum_{k=0}^n \binom{n}{k} - \binom{n}{0} \right) = \\ &= (2 + 1)^n - 1 - 2 \cdot (2^n - 1) = 3^n - 2^{n+1} + 1. \end{aligned}$$

Consequently, we have $\pi_1(n) = n + h(n) = 3^n - 2^{n+1} + n + 1$. \square

5 An Improvement

As a kind of intermezzo we briefly discuss a way to generate $\{L_n\}_{n \geq 1}$ by regular grammars $\{G_n^R\}_{n \geq 1}$. Although regular grammars are by no means context-free grammars in Chomsky normal form, Proposition 3.3 and Definition 4.1 suggest the following family $\{G_n^R\}_{n \geq 1}$.

Definition 5.1. The family $\{G_n^R\}_{n \geq 1}$ is given by $\{(V_n, \Sigma_n, P_n, S_n)\}_{n \geq 1}$ with

- $N_n = V_n - \Sigma_n = \{A_X \mid X \in \mathcal{P}_+(\Sigma_n)\}$,
- $P_n = \{A_{\{a\}} \rightarrow a \mid a \in \Sigma_n\} \cup \{A_X \rightarrow a A_{X - \{a\}} \mid X \subseteq \Sigma_n, a \in X, \#X \geq 2\}$,
- $S_n = A_{\Sigma_n}$. \square

Notice that in each rule of the form $A \rightarrow BC$ from G_n^1 (Definition 4.1) we first restricted B by some symbol A_i from $E_{n,1}$ and then we replaced A_i by the right-hand side of the unique rule $A_i \rightarrow a_i$.

Example 5.2. Again we show the case $n = 3$: $G_3^R = (V_3, \Sigma_3, P_3, S_3)$ with $S_3 = A_{123}$, $N_3 = \{A_{123}, A_{12}, A_{13}, A_{23}, A_1, A_2, A_3\}$ and $P_3 = \{A_{123} \rightarrow a_1 A_{23} \mid a_2 A_{13} \mid a_3 A_{12}, A_{12} \rightarrow a_1 A_2 \mid a_2 A_1, A_{13} \rightarrow a_1 A_3 \mid a_3 A_1, A_{23} \rightarrow a_2 A_3 \mid a_3 A_2\}$.

$a_3A_2, A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3\}$. The entities $E_{n,k}$ and $D(n,k)$ are as in Example 4.2; $\nu_R(3) = 7$ but now we have $\pi_R(3) = 12$. \square

Proposition 5.3. *For the family $\{G_n^R\}_{n \geq 1}$ of Definition 5.1 we have*

- (1) $D(n, k) = C(n, k)$ with $1 \leq k \leq n$,
- (2) $\nu_R(n) = 2^n - 1$,
- (3) $\pi_R(n) = n \cdot 2^{n-1}$.

Proof. We proceed as in the proof of Proposition 4.3(3): so let $\pi_R(n) = n + h(n)$ where $h(n) = \#\{A_X \rightarrow aA_{X-\{a\}} \mid X \subseteq \Sigma_n, a \in X, \#X \geq 2\}$.

If X has k ($k \geq 2$) elements, then $\{A_X \rightarrow aA_{X-\{a\}} \mid X \subseteq \Sigma_n, a \in X, \#X \geq 2\}$ contains k elements. Thus

$$\begin{aligned} h(n) &= \sum_{k=2}^n \binom{n}{k} k = \sum_{k=1}^n \frac{n! \cdot k}{k! (n-k)!} - \binom{n}{1} = \\ &= n \cdot \sum_{k=1}^n \frac{(n-1)!}{(k-1)! (n-k)!} - n = n \cdot \sum_{k=0}^{n-1} \frac{(n-1)!}{k! (n-1-k)!} - n = \\ &= n \cdot \sum_{k=0}^{n-1} \binom{n-1}{k} - n = n \cdot 2^{n-1} - n. \end{aligned}$$

Hence, we have $\pi_R(n) = n + h(n) = n + n \cdot 2^{n-1} - n = n \cdot 2^{n-1}$. \square

From Definition 5.1 we can obtain a family of context-free grammars in Chomsky normal form that generates $\{L_n\}_{n \geq 1}$ as follows.

Definition 5.4. The family $\{G_n^2\}_{n \geq 1}$ is given by $\{(V_n, \Sigma_n, P_n, S_n)\}_{n \geq 1}$ with

- $N_n = V_n - \Sigma_n = \{A_X \mid X \in \mathcal{P}_+(\Sigma_n)\}$,
- $P_n = \{A_{\{a\}} \rightarrow a \mid a \in \Sigma_n\} \cup \{A_X \rightarrow A_{\{a\}}A_{X-\{a\}} \mid X \subseteq \Sigma_n, a \in X, \#X \geq 2\}$,
- $S_n = A_{\Sigma_n}$. \square

Clearly, we have substituted A_i for a_i in all right-hand sides of rules from Definition 5.1 with left-hand side in $E_{n,2} \cup E_{n,3} \cup \dots \cup E_{n,n}$.

Example 5.5. For the case $n = 3$ we obtain: $G_3^2 = (V_3, \Sigma_3, P_3, S_3)$ with $S_3 = A_{123}$, $N_3 = \{A_{123}, A_{12}, A_{13}, A_{23}, A_1, A_2, A_3\}$ and $P_3 = \{A_{123} \rightarrow A_1A_{23} \mid A_2A_{13} \mid A_3A_{12}, A_{12} \rightarrow A_1A_2 \mid A_2A_1, A_{13} \rightarrow A_1A_3 \mid A_3A_1, A_{23} \rightarrow A_2A_3 \mid A_3A_2, A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3\}$ with $\nu_2(3) = 7$ and $\pi_2(3) = 12$. \square

Proposition 5.6. *For the family $\{G_n^2\}_{n \geq 1}$ of Definition 5.4 we have*

- (1) $D(n, k) = C(n, k)$ with $1 \leq k \leq n$,

- (2) $\nu_2(n) = 2^n - 1$,
(3) $\pi_2(n) = n \cdot 2^{n-1}$.

Proof. The one-to-one correspondence between G_n^R and G_n^2 for each $n \geq 1$ also implies that Proposition 5.6 follows from the proof of Proposition 5.3. \square

Proposition 5.7. *For each $n \geq 1$, G_n^2 is an unambiguous context-free grammar in Chomsky normal form.*

Proof. By Definition 5.4, G_n^2 is a context-free grammar in Chomsky normal form. The rules in P_n imply that for each word $w = a_{i_1}a_{i_2}\cdots a_{i_n}$ in L_n , the linearly ordered sequence (i_1, i_2, \dots, i_n) uniquely determines the order in which the productions have to be applied in a leftmost way in order to obtain w ; viz.

$$\begin{aligned} S &= A_{\{i_1, i_2, \dots, i_n\}} \Rightarrow A_{i_1}A_{\{i_2, i_3, \dots, i_n\}} \Rightarrow a_{i_1}A_{\{i_2, i_3, \dots, i_n\}} \Rightarrow a_{i_1}A_{i_2}A_{\{i_3, \dots, i_n\}} \Rightarrow \\ &\Rightarrow a_{i_1}a_{i_2}A_{\{i_3, \dots, i_n\}} \Rightarrow a_{i_1}a_{i_2}A_{i_3}A_{\{i_4, \dots, i_n\}} \Rightarrow \cdots \Rightarrow a_{i_1}a_{i_2}a_{i_3}\cdots a_{i_n} = w. \end{aligned}$$

So there is exactly one leftmost derivation for each w in L_n ; hence G_n^2 is unambiguous. \square

With respect to the number of productions the grammars G_n^2 are superior to the ones of Definition 4.1 since for $n \geq 3$, we have $\pi_2(n) = n \cdot 2^{n-1} < 3^n - 2^{n+1} + n + 1 = \pi_1(n)$.

6 Inserting an Additional Terminal Symbol — 1

In this section we provide a family $\{G_n^3\}_{n \geq 1}$ that —apart from the first two elements which are given explicitly— is defined inductively by means of a grammatical transformation. First, we have a look at the three most simple grammars ($n = 1, 2, 3$).

Example 6.1. (1) ($n = 1$). Consider G_1^3 with $P_1 = \{S_1 \rightarrow a_1\}$. Then $L(G_1^3) = \{a_1\} = L_1$, $\nu_3(1) = 1$ and $\pi_3(1) = 1$.

(2) ($n = 2$). Let G_2^3 be defined by $P_2 = \{S_2 \rightarrow A_1A_2, A_1 \rightarrow a_1, A_2 \rightarrow a_2\}$. Now we have $L(G_2^3) = \{a_1a_2, a_2a_1\} = L_2$, $\nu_3(2) = 3$ and $\pi_3(2) = 4$.

(3) ($n = 3$). For G_3^3 we define $P_3 = \{S_3 \rightarrow A_1A_{23} \mid A_{13}A_2, A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3, A_{13} \rightarrow A_1A_3, A_{23} \rightarrow A_2A_3\}$. Then $L(G_3^3) = \{a_1a_2a_3, a_1a_3a_2, a_2a_1a_3, a_2a_3a_1, a_3a_1a_2, a_3a_2a_1\} = L_3$, $\nu_3(3) = 6$ and $\pi_3(3) = 11$.

(4) Adding another nonterminal A_{12} together with rules $S_3 \rightarrow A_3A_{12}$ and $A_{12} \rightarrow A_1A_2$ to G_3^3 does not affect the language $L(G_3^3)$; the resulting grammar has 7 nonterminals and 15 productions. \square

Note that in both grammars G_n^3 ($n = 2, 3$) of Example 6.1(2–3) all nonterminals are not recursive and that $P_n \subseteq N_n \times (N_n - \{S_n\})^2 \cup (N_n - \{S_n\}) \times \Sigma_n$.

Definition 6.2. The family $\{G_n^3\}_{n \geq 1}$ is given by $\{(V_n, \Sigma_n, P_n, S_n)\}_{n \geq 1}$ with

- (1) G_1^3 is as in Example 6.1(1).
- (2) G_2^3 is as in Example 6.1(2).
- (3) G_{n+1}^3 is obtained from G_n^3 ($n \geq 2$) by the steps (a), (b), (c) and (d).

N.B. First, note that L_n with $L_n = L(G_n^3)$ is a language over Σ_n , whereas L_{n+1} is a language over Σ_{n+1} . More precisely, we obtain the elements of L_{n+1} by inserting the new symbol a_{n+1} at each available spot in the strings of L_n . This observation is the crux of our grammatical transformation. We obtain the new grammar G_{n+1}^3 from G_n^3 as follows.

- (a) Each initial rule, e.g. $S_n \rightarrow AB$, is replaced by two rules: $S_{n+1} \rightarrow A'B \mid AB'$. A primed symbol indicates that in the subtree rooted by that primed symbol still an occurrence of the new symbol a_{n+1} should be inserted.
- (b) To each noninitial rule in G_n^3 of the form $A \rightarrow BC$, there correspond in G_{n+1}^3 three rules: $A \rightarrow BC$ and $A' \rightarrow B'C \mid BC'$. The latter two rules are added “to propagate the primes”.
- (c) For each (noninitial) rule in G_n^3 of the form $A \rightarrow a$, there are the following associated rules in G_{n+1}^3 : $A \rightarrow a$ and $A' \rightarrow AA_{n+1}$, where A_{n+1} is new nonterminal symbol not yet present in N_n . The last rule will place a_{n+1} to the left or the right, respectively, of the a generated by A ; cf. also the next, final step in the construction.
- (d) Finally, we add the new rule $A_{n+1} \rightarrow a_{n+1}$ to P_{n+1} . □

It is now a routine matter to verify that (i) $L(G_{n+1}^3) = L_{n+1}$, (ii) each nonterminal symbol is not recursive in G_{n+1}^3 , and (iii) P_{n+1} does not contain a rule of the form $S_{n+1} \rightarrow a$ ($a \in \Sigma_{n+1}$).

Example 6.3. (1) By the grammatical transformation of Definition 6.2(3) we can obtain G_3^3 of Example 6.1(3) from G_2^3 from Example 6.1(2): $A'_1 = A_{13}$ and $A'_2 = A_{23}$.

(2) Next we apply this grammatical transformation to obtain G_4^3 from G_3^3 ; cf. Example 6.1(3) for the definition of G_3^3 .

The first step (a) yields: $S_4 \rightarrow A'_1 A_{23} \mid A_1 A'_{23} \mid A'_{13} A_2 \mid A_{13} A'_2$.

From the second step (b) we get: $A_{13} \rightarrow A_1 A_3$ and $A'_{13} \rightarrow A'_1 A_3 \mid A_1 A'_3$ as well as $A_{23} \rightarrow A_2 A_3$ and $A'_{23} \rightarrow A'_2 A_3 \mid A_2 A'_3$.

The last two steps (c) and (d) produce: $A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3$ together with $A'_1 \rightarrow A_1 A_4, A'_2 \rightarrow A_2 A_4, A'_3 \rightarrow A_3 A_4$ and $A_4 \rightarrow a_4$.

It is now easy to show that $L(G_4^3) = L_4, \nu_3(4) = 12$ and $\pi_3(4) = 30$. Of course, we may rename the nonterminal symbols: e.g., A'_{ij} by A_{ij4} and A'_i by A_{i4} ; cf. Section 4. □

Example 6.4. (1) Consider G_3^3 of Example 6.1(3). Then $E_{3,1} = \{A_1, A_2, A_3\}$, $E_{3,2} = \{A_{13}, A_{23}\}$ and $E_{3,3} = \{S_3\}$. The strict order of N_3 is: $A_1 \sqsubset A_{13} \sqsubset S_3$,

$A_3 \sqsubset A_{13}$, $A_2 \sqsubset A_{23} \sqsubset S_3$ and $A_3 \sqsubset A_{23}$.

(2) For the grammar G_4^3 of Example 6.3(2), we have $E_{4,1} = \{A_1, A_2, A_3, A_4\}$, $E_{4,2} = \{A_{13}, A_{23}, A'_1, A'_2, A'_3\}$, $E_{4,3} = \{A'_{13}, A'_{23}\}$ and $E_{4,4} = \{S_4\}$. The strict order of N_4 is given by $A_1 \sqsubset A'_1 \sqsubset A'_{13} \sqsubset S_4$, $A_1 \sqsubset A_{13} \sqsubset A'_{13}$, $A_2 \sqsubset A'_2 \sqsubset A'_{23} \sqsubset S_4$, $A_2 \sqsubset A_{23} \sqsubset A'_{23}$, $A_3 \sqsubset A'_3 \sqsubset A'_{23}$, $A_3 \sqsubset A_{13}$, $A_3 \sqsubset A_{23}$, $A_4 \sqsubset A'_1$, $A_4 \sqsubset A'_2$ and $A_4 \sqsubset A'_3 \sqsubset A'_{13}$. \square

Proposition 6.5. *For the family $\{G_n^3\}_{n \geq 1}$ of Definition 6.2 we have*

(1) $D(n, 1) = n$, $D(n, n-1) = 2$ ($n \geq 2$), $D(n, n) = 1$ and for each k with $2 \leq k \leq n-2$,

$$D(n, k) = D(n-1, k) + D(n-1, k-1),$$

(2) $\nu_3(1) = 1$ and for $n \geq 2$, $\nu_3(n) = 3 \cdot 2^{n-2}$,

(3) $\pi_3(1) = 1$ and for $n \geq 2$, $\pi_3(n) = \frac{5}{2} \cdot 3^{n-2} + 2^{n-1} - \frac{1}{2}$.

Proof. (1) Clearly, $D(n, n) = 1$ and $D(n, 1) = n$ as $E_{n,n} = \{S_n\}$ and $E_{n,1} = \{A_1, \dots, A_n\}$ because $A_i \rightarrow a_i$ are the only rules in P_n with terminal right-hand sides.

The other two equalities are easily established by induction over n using the properties of G_2^3 —particularly, the fact that $E_{2,1} = \{A_1, A_2\}$ —and the effect of the transformation given in Definition 6.2(3).

(2) From Definition 6.2(3) it follows that for the new set of nonterminal symbols N_{n+1} of G_{n+1}^3 we have

$$N_{n+1} = (N_n - \{S_n\}) \cup \{A' \mid A \in N_n - \{S_n\}\} \cup \{S_{n+1}, A_{n+1}\}.$$

This implies that $\nu_3(n+1) = 2 \cdot \nu_3(n)$. Solving this difference equation with initial condition $\nu_3(2) = 3$ (Definition 6.2(2) and Example 6.1(2)) yields $\nu_3(n) = 3 \cdot 2^{n-2}$ for $n \geq 2$.

(3) We write $\pi_3(n) = f(n) + g(n)$ for $n \geq 2$, where $f(n)$ is the number of initial productions and $g(n)$ is the number of noninitial productions in G_n^3 . By the transformation of Definition 6.2(3) we obtain the following recurrence relations: $f(n+1) = 2 \cdot f(n)$ with $f(2) = 2$, and $g(n+1) = 3 \cdot g(n) + 1$ with $g(2) = 2$. Solving these equations yields $f(n) = 2^{n-1}$ and $g(n) = \frac{5}{2} \cdot 3^{n-2} - \frac{1}{2}$ ($n \geq 2$); hence the result. \square

Proposition 6.5(2)–(3) may be rewritten as $\nu_3(n) = \lfloor 3 \cdot 2^{n-2} \rfloor$ and $\pi_3(n) = \lfloor \frac{5}{2} \cdot 3^{n-2} + 2^{n-1} - \frac{1}{2} \rfloor$, respectively ($n \geq 1$).

Note that the recurrence relation in Proposition 6.5(1) is identical to the one for the binomial coefficients $C(n, k)$, although the boundary conditions are different. It results in the Pascal-like triangle of Table 1.

n	$D(n, k)$									
	$k = 1$	2	3	4	5	6	7	8	9	10
1	1									
2	2	1								
3	3	2	1							
4	4	5	2	1						
5	5	9	7	2	1					
6	6	14	16	9	2	1				
7	7	20	30	25	11	2	1			
8	8	27	50	55	36	13	2	1		
9	9	35	77	105	91	49	15	2	1	
10	10	44	112	182	196	140	64	17	2	1

Table 1

$D(n, k)$ for G_n^3 ($1 \leq n \leq 10$).

Finally, we remark that the grammatical transformation of Definition 6.2(3) is of general interest in the following sense: given *any* context-free grammar G_n in Chomsky normal form that generates L_n (thus not just G_n^3), then it produces a context-free grammar G_{n+1} in Chomsky normal form for L_{n+1} . We will apply this observation in Section 9.

7 Inserting an Additional Terminal Symbol — 2

The family $\{G_n^3\}_{n \geq 1}$ is rather efficient with respect to the number of nonterminals as compared to the family $\{G_n^2\}_{n \geq 1}$: $\nu_3(n) = 3 \cdot 2^{n-2} < 2^n - 1 = \nu_2(n)$ for $n \geq 3$. The price we have to pay is an increase of the number of productions, since $\pi_3(n) = \frac{5}{2} \cdot 3^{n-2} + 2^{n-1} - \frac{1}{2} > n \cdot 2^{n-1} = \pi_2(n)$ for $n \geq 5$. In addition the degree of ambiguity of G_n^3 is rather high as can be seen from the following sample subderivations. Let $A \Rightarrow BC \Rightarrow^* w_B w_C$ with $B \Rightarrow^* w_B$ and $C \Rightarrow^* w_C$ be a subderivation according to G_n^3 . From the new grammar G_{n+1}^3 the substring $w_B a_{n+1} w_C$ can be obtained by $A' \Rightarrow B'C \Rightarrow^* w_B a_{n+1} w_C$ or by $A' \Rightarrow BC' \Rightarrow^* w_B a_{n+1} w_C$.

In this section we will modify the grammatical transformation of Definition 6.2 in such a way that the second subderivation is not possible, because the occurrence of a_{n+1} will always be introduced to the right of the terminal symbols a_1, a_2, \dots, a_n . This results in a family of grammars $\{G_n^4\}_{n \geq 1}$ with $A' \rightarrow AA_{n+1}$ rather than $A' \rightarrow AA_{n+1}$ in G_{n+1}^4 . In order to derive permutations from $\{a_{n+1}\}L_n$ we need the rule $S_{n+1} \rightarrow A_{n+1}S_n$ and to preserve S_n as well as all rules from G_n^4 .

Definition 7.1. The family $\{G_n^4\}_{n \geq 1}$ is given by $\{(V_n, \Sigma_n, P_n, S_n)\}_{n \geq 1}$ with

- (1) $G_1^4 = G_1^3$ (as in Example 6.1(1)).

(2) $G_2^4 = G_2^3$ (as in Example 6.1(2)).

(3) G_{n+1}^4 is obtained from G_n^4 ($n \geq 2$) by the steps (a), (b) and (c):

- (a) To each rule in G_n^4 of the form $A \rightarrow BC$, there corresponds in G_{n+1}^4 three rules: $A \rightarrow BC$ and $A' \rightarrow B'C \mid BC'$. The latter two rules are added “to propagate the primes”. The primed version S'_n of S_n becomes the initial symbol S_{n+1} of G_{n+1}^4 ; so $S_{n+1} = S'_n$.
- (b) We add the rules $S_{n+1} \rightarrow A_{n+1}S_n$ and $A_{n+1} \rightarrow a_{n+1}$ to P_{n+1} , where A_{n+1} is new nonterminal symbol not yet present in N_n .
- (c) For each (noninitial) rule in G_n^4 of the form $A \rightarrow a$, there are the following associated rules in G_{n+1}^4 : $A \rightarrow a$ and $A' \rightarrow AA_{n+1}$. The last rule will place a_{n+1} to the right of the a generated by A . \square

Example 7.2. (1) We construct G_3^4 from G_2^4 (i.e. G_2^3 , Example 6.1(1)). Definition 7.1(a)–(c) yields: $S_3 \rightarrow A'_1A_2 \mid A_1A'_2 \mid A'_2A_1 \mid A_2A'_1 \mid A_3S_2$, $S_2 \rightarrow A_1A_2 \mid A_2A_1$, $A'_1 \rightarrow A_1A_3$, $A'_2 \rightarrow A_2A_3$, $A_1 \rightarrow a_1$, $A_2 \rightarrow a_2$, $A_3 \rightarrow a_3$. Then $E_{3,3} = \{S_3\}$, $E_{3,2} = \{S_2, A'_1, A'_2\}$, $E_{3,1} = \{A_1, A_2, A_3\}$, and hence $D(3,3) = 1$, $D(3,2) = D(3,1) = 3$, $\nu_4(3) = 7$ and $\pi_4(3) = 12$.

(2) Next we derive G_4^4 from G_3^4 ; but first we rename A'_i by B_i ($i = 1, 2$) in G_3^4 of Example 7.2.(1) in order to avoid two types of primes with different meanings. Then we obtain: $S_4 \rightarrow B'_1A_2 \mid B_1A'_2 \mid A'_1B_2 \mid A_1B'_2 \mid B'_2A_1 \mid B_2A'_1 \mid A'_2B_1 \mid A_2B'_1 \mid A'_3S_2 \mid A_3S'_2 \mid A_4S_3$, $S_3 \rightarrow B_1A_1 \mid A_1B_2 \mid B_2A_1 \mid A_2B_1 \mid A_3S_2$, $S'_2 \rightarrow A'_1A_2 \mid A_1A'_2 \mid A'_2A_1 \mid A_2A'_1$, $S_2 \rightarrow A_1A_2 \mid A_2A_1$, $B_1 \rightarrow A_1A_3$, $B'_1 \rightarrow A'_1A_3 \mid A_1A'_3$, $B_2 \rightarrow A_2A_3$, $B'_2 \rightarrow A'_2A_3 \mid A_2A'_3$, $A'_1 \rightarrow A_1A_4$, $A'_2 \rightarrow A_2A_4$, $A'_3 \rightarrow A_3A_4$, $A_1 \rightarrow a_1$, $A_2 \rightarrow a_2$, $A_3 \rightarrow a_3$ and $A_4 \rightarrow a_4$. Hence $E_{4,4} = \{S_4\}$, $E_{4,3} = \{S_3, S'_2, B'_1, B'_2\}$, $E_{4,2} = \{S_2, B_1, B_2, A'_1, A'_2, A'_3\}$, $E_{4,1} = \{A_1, A_2, A_3, A_4\}$, $\nu_4(4) = 15$ and $\pi_4(4) = 35$. \square

There is a one-to-one correspondence between the nonterminals of G_n^4 and the elements of $\mathcal{P}_+(\Sigma_n)$. E.g. in Example 7.2(2) we have $S'_2 \leftrightarrow \{a_1, a_2, a_4\}$, $B_1 \leftrightarrow \{a_1, a_3\}$ and $B'_2 \leftrightarrow \{a_2, a_3, a_4\}$; cf. also Proposition 7.4(1)–(2) below.

Proposition 7.3. *For each $n \geq 1$, G_n^4 is an unambiguous context-free grammar in Chomsky normal form.*

Proof. Clearly, each G_n^4 is in Chomsky normal form. So it remains to show that each G_n^4 is unambiguous; this will be done by induction on n .

Basis ($n = 1, 2$): Obviously, both G_1^4 and G_2^4 are unambiguous grammars.

Induction hypothesis: G_n^4 is an unambiguous grammar.

Induction step: Let w be a word from $L(G_{n+1}^4)$. Then we distinguish two cases:

(i) $w \in \{a_{n+1}\} \cdot L(G_n^4)$, i.e. $w = a_{n+1}v$ for some $v \in L(G_n^4)$. Since v does not possess an occurrence of a_{n+1} , a leftmost derivation of w has the form $S_{n+1} \Rightarrow A_{n+1}S_n \Rightarrow a_{n+1}S_n \Rightarrow^* a_{n+1}v$. By the induction hypothesis there is

only one leftmost derivation according to G_n for v from S_n . And notice that $P_n \subset P_{n+1}$, whereas rules from $P_{n+1} - P_n$ cannot interfere in the subderivation $S_n \Rightarrow^* v$. Consequently, $S_{n+1} \Rightarrow A_{n+1}S_n \Rightarrow a_{n+1}S_n \Rightarrow^* a_{n+1}v$ is the only leftmost derivation of w in G_{n+1}^4 .

(ii) $w \notin \{a_{n+1}\} \cdot L(G_n^4)$, i.e. $w = ua_i a_{n+1}v$ with $ua_i v \in L(G_n^4)$ and $a_i \in \Sigma_n$; note that $i \neq n+1$. As $ua_i \neq \lambda$, the occurrence of a_{n+1} in w cannot be introduced by the initial rule $S_{n+1} \rightarrow A_{n+1}S_n$, but it must be obtained by a leftmost subderivation $A'_i \Rightarrow A_i A_{n+1} \Rightarrow a_i A_{n+1} \Rightarrow a_i a_{n+1}$ using the unique rule $A_i \rightarrow a_i$ from P_n and the unique rule $A_{n+1} \rightarrow a_{n+1}$ from $P_{n+1} - P_n$. Consider, the following leftmost derivation of w :

$$S_{n+1} \Rightarrow^+ uA'_i \omega \Rightarrow uA_i A_{n+1} \omega \Rightarrow ua_i A_{n+1} \omega \Rightarrow ua_i a_{n+1} \omega \Rightarrow^+ ua_i a_{n+1} v = w.$$

Suppose there are two such derivations according to G_{n+1}^4 . Then we can obtain two different leftmost derivations for $ua_i v$ according to G_n^4 as follows: (1) replace the subderivation $uA'_i \omega \Rightarrow uA_i A_{n+1} \omega \Rightarrow ua_i A_{n+1} \omega \Rightarrow ua_i a_{n+1} \omega$ by $uA_i \omega \Rightarrow ua_i \omega$, (2) remove all primes from primed symbols, and (3) change all remaining occurrences of a_{n+1} into λ .

However, the existence of two different leftmost derivations for $ua_i v$ in G_n^4 contradicts the induction hypothesis, i.e. the unambiguity of G_n^4 . \square

Proposition 7.4. *For the family $\{G_n^4\}_{n \geq 1}$ of Definition 7.1 we have*

- (1) $D(n, k) = C(n, k)$ for $1 \leq k \leq n$,
- (2) $\nu_4(n) = 2^n - 1$,
- (3) $\pi_4(n) = \frac{5}{4} \cdot 3^{n-1} + \frac{1}{2}n - \frac{3}{4}$.

Proof. (1) From Definition 7.1(a)–(c) it follows that $D(n, k) = D(n-1, k) + D(n-1, k-1)$ with $D(n, n) = 1$ and $D(n, 1) = n$ ($1 \leq k \leq n$). Hence $D(n, k) = C(n, k)$; cf. [9,14].

(2) Obviously, $\nu_4(n) = \sum_{k=1}^n D(n, k) = \sum_{k=1}^n C(n, k) = 2^n - 1$ for $n \geq 2$ [9]. Alternatively, we have $N_{n+1} = N_n \cup \{A' \mid A \in N_n - \{S_n\}\} \cup \{S_{n+1}, A_{n+1}\}$ which yields the difference equation $\nu_4(n+1) = 2 \cdot \nu_4(n) + 1$ with $\nu_4(2) = 3$. Solving this equation gives the same result.

(3) We write π_4 as $\pi_4(n) = f(n) + g(n)$ where $g(n)$ is the number of terminal rules $A_i \rightarrow a_i$ and $f(n)$ the number of remaining rules. Then $g(n) = n$, whereas $f(n+1) = 3 \cdot f(n) + n + 1$ with $f(2) = 2$. Let f_h be the solution of the corresponding homogeneous equation $f_h(n+1) = 3 \cdot f_h(n)$, i.e. $f_h(n) = c \cdot 3^n$. For a particular solution we try $f_p(n) = an + b$ which yields $a = -\frac{1}{2}$ and $b = -\frac{3}{4}$; thus $f_p(n) = -\frac{1}{2}n - \frac{3}{4}$. Finally, we use the initial condition $f(2) = 2$ to determine the constant c from $f(n) = f_h(n) + f_p(n) = c \cdot 3^n - \frac{1}{2}n - \frac{3}{4}$. Then $c = \frac{5}{12}$ which implies $\pi_4(n) = f(n) + g(n) = \frac{5}{12} \cdot 3^n - \frac{1}{2}n - \frac{3}{4} + n = \frac{5}{4} \cdot 3^{n-1} + \frac{1}{2}n - \frac{3}{4}$

($n \geq 2$). Substituting $n = 1$ in this expression gives $\pi_4(1) = 1$ as well. \square

Although we obtained unambiguous grammars (Proposition 7.3), the price we have to pay for this is high (Proposition 7.4): viz. $\nu_4(n) = 2^n - 1 > 3 \cdot 2^{n-2} = \nu_3(n)$ and $\pi_4(n) = \frac{5}{4} \cdot 3^{n-1} + \frac{1}{2}n - \frac{3}{4} > \frac{5}{2} \cdot 3^{n-2} + 2^{n-1} - \frac{1}{2} = \pi_3(n)$ for $n \geq 3$.

The grammatical transformation of Definition 7.1(3) is as general as the one of Definition 6.2(3): it is applicable to any context-free grammar G_n in Chomsky normal form for L_n and it yields a context-free grammar G_{n+1} in Chomsky normal form with $L(G_{n+1}) = L_{n+1}$; cf. Section 9 for an application.

8 Divide and Conquer

The families of grammars considered in the previous sections all share the property that $E_{n,k} \neq \emptyset$ for all k ($1 \leq k \leq n$). In this section we consider a family of grammars $\{G_n^5\}_{n \geq 1}$ that is a divide-and-conquer modification of $\{G_n^1\}_{n \geq 1}$ of Section 4 in the sense that —instead of dividing $X \cup Y$ in all possible disjoint nonempty X and Y — we restrict the subdivisions of $X \cup Y$ to almost equally sized X and Y . As a consequence we have that for some k , the sets $E_{n,k}$ are empty, whenever $n \geq 4$.

```

 $E_{n,1} := \{A_{\{a\}} \mid a \in \Sigma_n\}; \quad N_n := E_{n,1}; \quad M := \{A_{\Sigma_n}\};$ 
 $P_n := \{A_{\{a\}} \rightarrow a \mid a \in \Sigma_n\};$ 
while  $M - E_{n,1} \neq \emptyset$  [i.e.  $\exists A_X \in M : X \subseteq \Sigma_n$  and  $\#X \geq 2$ ] do
  begin
     $S(X) := \{(Y, Z) \mid Y \subset X, \#Y = \lceil \frac{1}{2} \#X \rceil, Z = X - Y\};$ 
     $P_n := P_n \cup \{A_X \rightarrow A_Y A_Z \mid (Y, Z) \in S(X)\};$ 
     $N_n := N_n \cup \{A_X\};$ 
     $M := (M - \{A_X\}) \cup \{A_Y, A_Z \mid (Y, Z) \in S(X)\}$ 
  end

```

Fig. 1. Algorithm to determine N_n and P_n of G_n^5 .

Definition 8.1. The family $\{G_n^5\}_{n \geq 1}$ is given by $\{(V_n, \Sigma_n, P_n, S_n)\}_{n \geq 1}$ with

- $S_n = A_{\Sigma_n}$, and
- the sets N_n and P_n are determined by the algorithm in Figure 1. \square

Example 8.2. (1) For $n = 4$ Definition 8.1 yields G_4^5 with $S_4 = A_{1234}$, $N_4 = E_{4,1} \cup E_{4,2} \cup E_{4,3} \cup E_{4,4}$, $E_{4,1} = \{A_1, A_2, A_3, A_4\}$, $E_{4,2} = \{A_{12}, A_{13}, A_{14}, A_{23}, A_{24}, A_{34}\}$, $E_{4,3} = \emptyset$, $E_{4,4} = \{A_{1234}\}$, $P_4 = \{A_{1234} \rightarrow A_{12}A_{34} \mid A_{13}A_{24} \mid A_{14}A_{23}, A_{12} \rightarrow A_1A_2, A_{13} \rightarrow A_1A_3, A_{14} \rightarrow A_1A_4, A_{23} \rightarrow A_2A_3, A_{24} \rightarrow A_2A_4, A_{34} \rightarrow A_3A_4, A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3, A_4 \rightarrow a_4\}$, $\nu_5(4) = 11$ and $\pi_5(4) = 22$.

(2) Similarly, for $n = 5$ we obtain G_5^5 with $S_5 = A_{12345}$, $N_5 = E_{5,5} \cup E_{5,4} \cup E_{5,3} \cup E_{5,2} \cup E_{5,1}$, $E_{5,5} = \{A_{12345}\}$, $E_{5,4} = \emptyset$, $E_{5,3} = \{A_{123}, A_{124}, A_{125}, A_{134}, A_{135}, A_{145}, A_{234}, A_{235}, A_{245}, A_{345}\}$, $E_{5,2} = \{A_{12}, A_{13}, A_{14}, A_{15}, A_{23}, A_{24}, A_{25}, A_{34}, A_{35}, A_{45}\}$, $E_{5,1} = \{A_1, A_2, A_3, A_4, A_5\}$, $P_5 = \{A_{12345} \rightarrow A_{123}A_{45} \mid A_{124}A_{35} \mid A_{125}A_{34} \mid A_{134}A_{25} \mid A_{135}A_{24} \mid A_{145}A_{23} \mid A_{234}A_{15} \mid A_{235}A_{14} \mid A_{245}A_{13} \mid A_{345}A_{12}, A_{123} \rightarrow A_{12}A_3 \mid A_{13}A_2 \mid A_{23}A_1, A_{124} \rightarrow A_{12}A_4 \mid A_{14}A_2 \mid A_{24}A_1, A_{125} \rightarrow A_{12}A_5 \mid A_{15}A_2 \mid A_{25}A_1, A_{134} \rightarrow A_{13}A_4 \mid A_{14}A_3 \mid A_{34}A_1, A_{135} \rightarrow A_{13}A_5 \mid A_{15}A_3 \mid A_{35}A_1, A_{145} \rightarrow A_{14}A_5 \mid A_{15}A_4 \mid A_{45}A_1, A_{234} \rightarrow A_{23}A_4 \mid A_{24}A_3 \mid A_{34}A_2, A_{235} \rightarrow A_{23}A_5 \mid A_{25}A_3 \mid A_{35}A_2, A_{245} \rightarrow A_{24}A_5 \mid A_{25}A_4 \mid A_{45}A_2, A_{345} \rightarrow A_{34}A_5 \mid A_{35}A_4 \mid A_{45}A_3, A_{12} \rightarrow A_1A_2, A_{13} \rightarrow A_1A_3, A_{14} \rightarrow A_1A_4, A_{15} \rightarrow A_1A_5, A_{23} \rightarrow A_2A_3, A_{24} \rightarrow A_2A_4, A_{25} \rightarrow A_2A_5, A_{34} \rightarrow A_3A_4, A_{35} \rightarrow A_3A_5, A_{45} \rightarrow A_4A_5, A_1 \rightarrow a_1, A_2 \rightarrow a_2, A_3 \rightarrow a_3, A_4 \rightarrow a_4, A_5 \rightarrow a_5\}$, $\nu_5(5) = 26$ and $\pi_5(5) = 65$.

(3) For $n = 8$ the algorithm of Definition 8.1 produces a grammar G_8^5 with $E_{8,7} = E_{8,6} = E_{8,5} = E_{8,3} = \emptyset$. Similarly, the grammar G_{10}^5 satisfies $E_{10,9} = E_{10,8} = E_{10,7} = E_{10,6} = E_{10,4} = \emptyset$. \square

The next result follows from the structure of the algorithm in Definition 8.1; cf. Figure 1.

Proposition 8.3. *For the family $\{G_n^5\}_{n \geq 1}$ of Definition 8.1 we have*

- (1) $D(n, k) = \mathbf{if } k \in \{\lceil n/2^i \rceil, \lfloor n/2^i \rfloor \mid 0 \leq i \leq \lceil \log_2 n \rceil\} \mathbf{ then } C(n, k) \mathbf{ else } 0$,
- (2) $\nu_5(n) = \sum_{k=1}^n D(n, k)$,
- (3) $\pi_5(n) = \sum_{k=1}^n D(n, k) \cdot C(k, \lceil k/2 \rceil)$. \square

n	$D(n, k)$									
	$k = 1$	2	3	4	5	6	7	8	9	10
1	1									
2	2	1								
3	3	3	1							
4	4	6	0	1						
5	5	10	10	0	1					
6	6	15	20	0	0	1				
7	7	21	35	35	0	0	1			
8	8	28	0	70	0	0	0	1		
9	9	36	84	126	126	0	0	0	1	
10	10	45	120	0	252	0	0	0	0	1

Table 2
 $D(n, k)$ for G_n^5 ($1 \leq n \leq 10$).

The values of $D(n, k)$ for $1 \leq n \leq 10$ are given in Table 2. Unfortunately, a closed form for $\nu_5(n)$ and $\pi_5(n)$ is very hard or even impossible to obtain; a situation very common in analyzing these divide-and-conquer approaches; cf. e.g. pp. 62–78 in [17] or [20]. A numerical evaluation and a comparison

with $\nu_i(n)$ and $\pi_i(n)$ ($1 \leq i \leq 4$) can be found in Section 9. These numerical values suggest that both functions ν_5 and π_5 satisfy $f(n+2) > 2 \cdot f(n)$ and $f(n+1) > f(n)$, confirming the exponential growth of these complexity measures; cf. Section 1 [16,7].

9 Concluding Remarks

In the previous sections we discussed a few ways to generate the set of all permutations of an alphabet of n symbols by context-free grammars in Chomsky normal form. For the resulting families of grammars $\{G_n^i\}_{n \geq 1}$ ($1 \leq i \leq 5$) we considered the values of the descriptive complexity measures $\nu_i(n)$ (i.e. the number of nonterminal symbols) and $\pi_i(n)$ (i.e. the number of productions) of G_n^i . A comparison of actual values for $1 \leq n \leq 16$ of these measures is given in Tables 3 and 4.

n	$\nu_1(n)$	$\nu_2(n)$	$\nu_3(n)$	$\nu_4(n)$	$\nu_5(n)$	$\nu_6(n)$	$\nu_7(n)$
1	1	1	1	1	1	1	1
2	3	3	3	3	3	3	3
3	7	7	6	7	7	7	6
4	15	15	12	15	11	11	11
5	31	31	24	31	26	23	22
6	63	63	48	63	42	42	42
7	127	127	96	127	99	85	84
8	255	255	192	255	107	107	107
9	511	511	384	511	382	215	214
10	1023	1023	768	1023	428	428	428
11	2047	2047	1536	2047	1156	857	856
12	4095	4095	3072	4095	1223	1223	1223
13	8191	8191	6144	8191	4525	2447	2446
14	16383	16383	12288	16383	4903	4903	4903
15	32767	32767	24576	32767	14811	9807	9806
16	65535	65535	49152	65535	14827	14827	14827

Table 3
 $\nu_i(n)$ ($1 \leq i \leq 7$; $1 \leq n \leq 16$).

Note that, for instance, the grammars $\{G_n^1\}_{n \geq 1}$ and $\{G_n^3\}_{n \geq 1}$ from Sections 4 and 6 respectively, are ambiguous. Now let for each $G_n = (V_n, \Sigma_n, P_n, S_n)$ that generates L_n , $\delta(n)$ denote the total number of possible leftmost derivations according to P_n ; thus $\delta(n) \geq n!$. E.g. for G_3^3 we have $\delta_3(3) = 8 > 3!$; so G_3^3 is not minimal with respect to this complexity measure. And the family of trivial

n	$\pi_1(n)$	$\pi_2(n)$	$\pi_3(n)$	$\pi_4(n)$	$\pi_5(n)$	$\pi_6(n)$	$\pi_7(n)$
1	1	1	1	1	1	1	1
2	4	4	4	4	4	4	4
3	15	12	11	12	12	12	11
4	54	32	30	35	22	22	22
5	185	80	83	103	65	64	59
6	608	192	234	306	116	116	116
7	1939	448	671	914	399	344	317
8	6058	1024	1950	2737	554	554	554
9	18669	2304	5723	8205	2475	1556	1535
10	57012	5120	16914	24608	3232	3232	3232
11	173063	11264	50231	73816	14938	9688	9185
12	523262	24576	149670	221439	20208	20208	20208
13	1577953	53248	446963	664307	101413	60614	58577
14	4750216	114688	1336794	1992910	130846	130846	130846
15	14283387	245760	4002191	5978718	691890	392526	384347
16	42915666	524288	11990190	17936141	924946	924946	924946

Table 4

$\pi_i(n)$ ($1 \leq i \leq 7$; $1 \leq n \leq 16$).

i	1	2	3	4	5
$\nu_i(n)$	A000225	A000225	A003945	A000225	A012272*
$\pi_i(n)$	A090326*	A001787	A090327*	A090328*	A077277*

Table 5

Integer sequences.

grammars mentioned in Section 1 —viz. $\{G_n^0\}_{n \geq 1}$ with $G_n^0 = (V_n, \Sigma_n, P_n, S_n)$, $N_n = \{S_n\}$ and $P_n = \{S_n \rightarrow w \mid w \in L_n\}$, although not in Chomsky normal form— satisfies $\nu_0(n) = 1$, and $\pi_0(n) = \delta_0(n) = n!$. From Propositions 5.7 and 7.3 it follows that for the families $\{G_n^2\}_{n \geq 1}$ and $\{G_n^4\}_{n \geq 1}$, we have $\delta_2(n) = \delta_4(n) = n!$ as well. Quite generally, one may ask whether there exist trade-offs between the complexity measures ν , π and δ . And, of course, the question remains whether there exists a family of minimal grammars with respect to the descriptive complexity measures $\nu(n)$ and $\pi(n)$.

It is rather straightforward to show that the family of grammars $\{G_n^R\}_{n \geq 1}$ is minimal with respect to both $\nu_R(n)$ and $\pi_R(n)$ for the class of regular (or right-linear) grammars that generate $\{L_n\}_{n \geq 1}$. But for the class of context-free grammars in Chomsky normal form that generate $\{L_n\}_{n \geq 1}$ the situation is not that clear. For the families $\{G_n^i\}_{n \geq 1}$ ($1 \leq i \leq 5$) studied in Sections 4–8,

$\{G_n^5\}_{n \geq 1}$ happens to have the least number of nonterminals, whereas $\{G_n^2\}_{n \geq 1}$ has the least number of productions. Note that the family $\{G_n^5\}_{n \geq 1}$ is not minimal with respect to ν . We can slightly improve upon $\{G_n^5\}_{n \geq 1}$ in the following way:

- (i) for even values of n we take G_n equal to G_n^5 , and
- (ii) for odd values of n —i.e. in case $n = 2k + 1$ — we take G_{2k}^5 and we apply the grammatical transformation of Section 6 or 7 to obtain G_n ; cf. the remarks at the end of Sections 6 and 7.

Applying the grammatical transformation from Definition 7.1(3) in this way, together with the recurrence relations $\nu_4(n+1) = 2 \cdot \nu_4(n) + 1$ and $\pi_4(n+1) = 3 \cdot \pi_4(n) - n + 2$, yields the family $\{G_n^6\}_{n \geq 1}$. Similarly, the family $\{G_n^7\}_{n \geq 1}$ is obtained by using the grammatical transformation of Definition 6.2(3) and the recurrences $\nu_3(n+1) = 2 \cdot \nu_3(n)$ and $\pi_3(n+1) = 3 \cdot \pi_3(n) - 2^{n-1} + 1$. The resulting values of $\nu_6(n)$, $\pi_6(n)$, $\nu_7(n)$ and $\pi_7(n)$ for $1 \leq n \leq 16$ are in Tables 3 and 4. These modifications of $\{G_n^5\}_{n \geq 1}$ have a profitable effect on the $\pi(n)$ -values for odd n as well.

In Section 5 we defined a regular grammar G_n^R for L_n ($n \geq 1$). By standard methods G_n^R can be converted into a deterministic finite automaton for L_n . So Proposition 5.3 or 5.6 determines the state complexity [21] (and the “transition complexity”) of this automaton.

The construction of the grammar families in this paper has something in common with designing algorithms to generate permutations, although in our case we are somewhat limited: we are unable to apply transpositions (“swapping of symbols”) because a transposition—even in the simple case of swapping adjacent elements—is a context-dependent rewriting step inherently. For a classification of (functional) programs for generating permutations we refer to [19]. The family $\{G_n^3\}_{n \geq 1}$ corresponds to Algorithm A in [19], whereas the family $\{G_n^R\}_{n \geq 1}$ is more or less a “mirrored” instance of its Algorithm B.

In this paper we restricted ourselves to generating permutations. Of course, there are other algebraic or combinatorial objects that—restricted to size n or parameterized by n in an other way— can be represented as a finite formal language L_n for which one may proceed as in the previous sections. An example is in [2] where we restrict our attention to “circular shifts”; these special permutations give rise to functions $\nu(n)$ and $\pi(n)$ that are polynomially bounded in n rather than the exponential functions of the present paper; cf. Section 1 [16,7].

Finally, we mention that the result of evaluating functions like $\nu_i(n)$ and $\pi_i(n)$ for $n = 1, 2, 3, \dots$ ($1 \leq i \leq 7$) is a so-called integer sequence. Some of these are well known, other ones seem to be new. In Table 5 we give an overview: the codes in this table refer to N.J.A. Sloane’s “Database of Integer Sequences”

[18]: the starred items have been added recently as being new, whereas the sequences for $i = 6, 7$ have not been included because of their ad hoc character. Tables 1 and 2 are known in [18] as A029635 and A090349*, respectively.

Acknowledgements. I am indebted to Giorgio Satta for suggesting the subject and to Jeffrey Shallit for sending me a copy of [7].

References

- [1] B. Alspach, P. Eades & G. Rose, A lower-bound for the number of productions for a certain class of languages, *Discrete Appl. Math.* **6** (1983) 109-115.
- [2] P.R.J. Asveld, Generating all circular shifts by context-free grammars in Chomsky normal form, (in preparation).
- [3] P.R.J. Asveld & A. Nijholt, The inclusion problem for some subclasses of context-free languages, *Theor. Comp. Sci.* **230** (2000) 247-256.
- [4] W. Bucher, A note on a problem in the theory of grammatical complexity, *Theor. Comp. Sci.* **14** (1981) 337-344.
- [5] W. Bucher, H.A. Maurer & K. Culik II, Context-free complexity of finite languages, *Theor. Comp. Sci.* **28** (1984) 277-285.
- [6] W. Bucher, H.A. Maurer, K. Culik II & D. Wotschke, Concise description of finite languages, *Theor. Comp. Sci.* **14** (1981) 227-246.
- [7] K. Ellul, B. Krawetz, J. Shallit & M.-w. Wang, Regular expressions: New results and open problems (2003), manuscript, Dept. of Comp. Sci., University of Waterloo, Waterloo, Ontario, Canada,
- [8] S. Ginsburg & H.G. Rice, Two families of languages related to ALGOL, *J. Assoc. Comp. Mach.* **9** (1962) 350-371.
- [9] R.L. Graham, D.E. Knuth & O. Patashnik, *Concrete Mathematics* (1989), Addison-Wesley, Reading, MA.
- [10] J. Gruska, Some classifications of context-free languages, *Inform. Contr.* **14** (1969) 152-179.
- [11] M.A. Harrison, *Introduction to Formal Language Theory* (1978), Addison-Wesley, Reading, MA.
- [12] V.A. Iljuškin, The complexity of the grammatical description of context-free languages, *Dokl. Akad. Nauk SSSR* **203** (1972) 1244-1245 / *Soviet Math. Dokl.* **13** (1972) 533-535.
- [13] A. Kelemenová, Complexity of normal form grammars, *Theor. Comp. Sci.* **28** (1984) 299-314.

- [14] C.L. Liu, *Introduction to Combinatorial Mathematics* (1968), McGraw-Hill, New York, etc.
- [15] R.E. Mickens, *Difference Equations — Theory and Applications* (1987), Second Edition (1990), Chapman & Hall, New York, London.
- [16] G. Satta, personal communication (2002).
- [17] R. Sedgewick & Ph. Flajolet, *An Introduction to the Analysis of Algorithms* (1996), Addison-Wesley, Reading, Ma.
- [18] N.J.A. Sloane, *Database of Integer Sequences*,
<http://www.research.att.com/~njas/sequences/Seis.html>
An earlier, printed version appeared as: N.J.A. Sloane & S. Plouffe, *The Encyclopedia of Integer Sequences* (1995), Academic Press, San Diego CA, etc.
- [19] R.W. Topor, Functional programs for generating permutations, *Computer J.* **25** (1982) 257–263.
- [20] X. Wang & Q. Fu, A frame for general divide-and-conquer recurrences, *Inform. Process. Lett.* **59** (1996) 45–51.
- [21] S. Yu, Q. Zhuang & K. Salomaa, The state complexity of some basic operations on regular languages, *Theor. Comp. Sci.* **125** (1994) 315–328.