



On the Grammatical Complexity of Finite Languages

Markus Holzer¹(✉) and Simon Wolfsteiner²

¹ Institut für Informatik, Universität Giessen, Arndtstr. 2, 35392 Giessen, Germany
holzer@informatik.uni-giessen.de

² Institut für Diskrete Mathematik und Geometrie, TU Wien,
Wiedner Hauptstr. 8–10, 1040 Wien, Austria
simon.wolfsteiner@tuwien.ac.at

Abstract. We study the grammatical production complexity of *finite* languages w.r.t. (i) different interpretations of approximation, i.e., equivalence, cover, and scattered cover, and (ii) whether the underlying grammar generates a finite or infinite language. In case the generated language is infinite, the intersection with all words up to a certain length has to be considered in order to obtain the finite language under consideration. In this way, we obtain six different measures for regular, linear context-free, and context-free grammars. We compare these measures according to the taxonomy introduced in [J. DASSOW, GH. PÄUN: *Regulated Rewriting in Formal Language Theory*, 1989] with each other by fixing the grammar type and varying the complexity measure and the other way around, that is, by fixing the complexity measure and varying the grammar type. In both of these cases, we develop an almost complete picture, which gives new and interesting insights into the old topic of grammatical production complexity.

1 Introduction

Measures of descriptional complexity or cost functions have a long and fruitful history. Most approaches to defining descriptional complexity measures are based on quantifying the ability of a device—automaton or grammar—to approximate languages (by finite subsets). The interesting quantities in the case of automata and grammars are, e.g., the number of states or transitions and the number of nonterminals or productions, respectively. For instance, finite languages can be represented by ordinary deterministic finite automata (DFAs) or by cover finite automata (CFAs)—roughly speaking a CFA is a DFA A and a natural number ℓ whose accepted language is defined as $L(A) \cap \Sigma^{\leq \ell}$. For a precise definition of

This research was completed while the author was on leave at the Institut für Informatik, Universität Giessen, Arndtstr. 2, 35392 Giessen, Germany, in 2017 and is supported by the Vienna Science Fund (WWTF) project VRG12-004.

CFA, we refer the reader to [5]. This gives rise to two different complexity measures: for a finite language $L \subseteq \Sigma^{\leq \ell}$, one defines the DFA state and the CFA state complexity as¹

$$\text{DFAc}(L) = \min\{|A| \mid A \text{ is a DFA and } L = L(A)\} \quad (1)$$

and

$$\text{DFAc}_\infty(L) = \min\{|A| \mid A \text{ is a DFA and } L = L(A) \cap \Sigma^{\leq \ell}\},$$

respectively, where $|A|$ refers to the number of states of the automaton A . Obviously, $\text{DFAc}_\infty(L) \leq \text{DFAc}(L)$, for every finite language L . Note that equality is possible for certain finite languages. It is worth mentioning that although these measures look very similar in their definitions, they can differ quite immensely when applied to the same language. For instance, there is a finite language L such that $\text{DFAc}_\infty(L) = 1$, but $\text{DFAc}(L) = n$, for every $n \geq 1$. Hence, the gap between both measures can be arbitrarily large. In general, complexity gaps can be classified according to the different growth rates of the complexity measures. To this end, a notion was introduced in [6], which defines three categories \leq^1 , \leq^2 , and \leq^3 of increasing complexity gaps; a precise definition is given later. The aforementioned complexity gap of arbitrary size between DFAs and CFAs is a gap of highest type, which is simply written as $c_\infty \leq_{\text{DFA}}^3 c$. This is one example of a complexity measure, but one can find legions of other automata-based descriptive complexity measures in the literature.

In this paper, we study the *grammatical* production complexity of finite languages. This topic is not new, and already in [2–4], measures similar in definition to (1), for regular (REG), linear context-free (LIN), and context-free grammars (CFG) have been investigated; they are named Xc , for $X \in \{\text{REG}, \text{LIN}, \text{CF}\}$. For instance, there it has been shown that there are incompressible finite languages for each grammar type mentioned above, i.e., languages that need at least as many productions of a certain type as there are words in that language. To the best of our knowledge, a classification of these grammatical measures in the sense of [6] has not been done yet. We close this gap and, moreover, also consider natural variants of Xc by varying the equivalence condition $L = L(A)$ in (1) to $L \subseteq L(A)$ and $L(A)$ finite (cover) or even $L \leq L(A)$ and $L(A)$ finite (scattered cover), where \leq refers to the scattered subword relation—similarly this can be done for $L = L(A) \cap \Sigma^{\leq \ell}$ in Xc_∞ , too. This leads to the additional grammatical measures (i) Xcc and Xcc_∞ (cover) and (ii) Xsc and Xsc_∞ (scattered cover), for $X \in \{\text{REG}, \text{LIN}, \text{CF}\}$. The variation Xcc is inspired by recent results on proof complexity in first-order logic, a research topic, which, from a first glance, seems completely unrelated to grammatical complexity, that connects the number of certain inference rules used in a specific logical calculus with the number of productions needed to *cover* a certain finite language [7]. For further results on the cover complexity of finite languages, see also [8].

¹ Observe that it is common in the literature to refer to the DFA and the CFA state complexity as sc and csc , respectively. We adapted the notation in order to be consistent with the notation used throughout this paper.

We compare these measures according to the taxonomy introduced in [6] with each other by (i) fixing the grammar type and varying the complexity measure and (ii) by fixing the complexity measure and varying the grammar type. In both of these cases, we develop an almost complete picture. As a byproduct, we also show that there are finite languages with large complexity. More precisely, the language of even length palindromes $P_n = \{ w\$w^R \mid w \in \{a, b\}^{\leq n} \}$ requires at least $\Omega(2^n)$ productions to be generated by a regular grammar. Moreover, the triple language $T_n = \{ w\$w\#w \mid w \in \{a, b\}^n \}$ can only be generated by grammars of type X , for $X \in \{\text{REG}, \text{LIN}, \text{CF}\}$, by simply listing all words in T_n , i.e., $\text{Xc}(T_n) = \Omega(2^n)$.

2 Preliminaries

We assume that the reader is familiar with the basic notions of formal language theory as contained in [10]. Nevertheless, to fix notation and terminology, we introduce the basic notions and results relevant to this paper in this section.

Let Σ be a finite alphabet. Then Σ^* denotes the set of all words over the finite alphabet Σ including the *empty word* ε and we write Σ^+ for $\Sigma^* \setminus \{\varepsilon\}$. The *length* of a word w in Σ^* is denoted by $|w|$. In particular, the length of the empty word ε is 0, i.e., $|\varepsilon| = 0$. The *reversal* of a word is defined as follows: $\varepsilon^R = \varepsilon$ and $(wa)^R = aw^R$, for $w \in \Sigma^*$ and $a \in \Sigma \cup \{\varepsilon\}$. Let $\ell \geq 0$. Then Σ^ℓ and $\Sigma^{\leq \ell}$ refers to the set of all words over Σ of length exactly ℓ and at most ℓ , respectively. A subset L of Σ^* is called a *language*. Any language $L \subseteq \Sigma^{\leq \ell}$, for $\ell \geq 0$, is called *finite* and, unless stated otherwise, we always assume $\ell = \max\{|w| \mid w \in L\}$. If L is a subset of Σ^ℓ , for $\ell \geq 0$, then L is called a *uniform language*. This means that in a uniform language all words have the same length.

A *context-free grammar* (CFG) is a quadruple $G = (N, \Sigma, P, S)$, where N and Σ are disjoint alphabets of *nonterminals* and *terminals*, respectively, $S \in N$ is the *start symbol*, and P is a finite set of *productions* of the form $A \rightarrow \alpha$, where $A \in N$ and $\alpha \in (N \cup \Sigma)^*$. As usual, the derivation relation of G is denoted by \Rightarrow_G and the reflexive and transitive closure of \Rightarrow_G is written as \Rightarrow_G^* . The *language generated* by G is defined as $L(G) = \{ w \in \Sigma^* \mid S \Rightarrow_G^* w \}$. We also consider the following restrictions of context-free grammars: (i) a context-free grammar is said to be *linear context-free* (LIN) if the productions are of the form $A \rightarrow \alpha$, where $A \in N$ and $\alpha \in \Sigma^*(N \cup \{\varepsilon\})\Sigma^*$, and (ii) a context-free grammar is said to be *right-linear* or *regular* (REG) if the productions are of the form $A \rightarrow \alpha$, where $A \in N$ and $\alpha \in \Sigma^*(N \cup \{\varepsilon\})$. Furthermore, Γ will denote the set of grammar types in the sequel, that is, $\Gamma = \{\text{REG}, \text{LIN}, \text{CF}\}$.

Let $G = (N, \Sigma, P, S)$ be a context-free grammar. By $|G|$, we denote the number of productions of G , i.e., the cardinality of P . Then the (*exact*) *X-complexity* (or *X-complexity* for short) of a finite language L w.r.t. X -grammars, for $X \in \Gamma$, is defined as

$$\text{Xc}(L) = \min\{|G| \mid G \text{ is an } X\text{-grammar with } L = L(G) \text{ and } L(G) \text{ finite}\}.$$

The additional condition that $L(G)$ is finite is redundant, but becomes important whenever we replace $L = L(G)$ by $L \subseteq L(G)$ or some other language-relating

property. Similarly, the *infinite X-complexity* of a finite language $L \subseteq \Sigma^{\leq \ell}$ is defined as

$$Xc_{\infty}(L) = \min\{|G| \mid G \text{ is an } X\text{-grammar with } L = L(G) \cap \Sigma^{\leq \ell}\}.$$

Note that in the definition of Xc_{∞} , the grammar G is allowed to generate an infinite language. If we replace $L = L(G)$ and $L = L(G) \cap \Sigma^{\leq \ell}$ in the definitions of $Xc(L)$ and $Xc_{\infty}(L)$, respectively, by $L \subseteq L(G)$ and $L \subseteq L(G) \cap \Sigma^{\leq \ell}$, respectively, then we get the definitions for the *X-cover-complexity* $Xcc(L)$ and the *infinite X-cover-complexity* $Xcc_{\infty}(L)$, respectively. The *scattered subword relation* \leq is defined as follows: let $w = w_1u_1w_2u_2 \dots u_{n-1}w_n$ be a word with $w_i, u_j \in \Sigma^*$, for $1 \leq i \leq n$ and $1 \leq j \leq n-1$. Then the word $w' = w_1w_2 \dots w_n$ is called a *scattered subword* of w and we write $w' \leq w$ in this case. We extend the relation \leq from words to languages L_1 and L_2 as follows: $L_1 \leq L_2$ if for all words $w_1 \in L_1$, there is a word $w_2 \in L_2$ such that $w_1 \leq w_2$. If $L_1 \leq L_2$ holds, we say that L_1 is a *scattered sublanguage* of L_2 . We, obtain the definitions for $Xsc(L)$ and $Xsc_{\infty}(L)$ if we replace $L = L(G)$ and $L = L(G) \cap \Sigma^{\leq \ell}$ by $L \leq L(G)$ and $L \leq L(G) \cap \Sigma^{\leq \ell}$, respectively, in the definitions of $Xc(L)$ and $Xc_{\infty}(L)$, respectively. This results in the *X-scattered-complexity* and the *infinite X-scattered-complexity*, respectively.

Note that the definitions of Xc , Xcc , and Xsc also have the additional requirement that $L(G)$ is a finite language. In the following, \mathcal{M} will denote the set of measure types, i.e., $\mathcal{M} = \{c, cc, sc, c_{\infty}, cc_{\infty}, sc_{\infty}\}$. By definition, for $\tau \in \mathcal{M}$, the following relations hold:

$$CFG \leq_{\tau} LIN \leq_{\tau} REG, \tag{2}$$

where, for $X, Y \in \Gamma$, we define $X \leq_{\tau} Y$ if and only if $X_{\tau}(L) \leq Y_{\tau}(L)$, for all finite languages L . In case that $X \leq_{\tau} Y$, we say that X is *more succinct than* Y (w.r.t. the complexity measure τ).

We say that G is a *minimal X-grammar*, for $X \in \Gamma$, w.r.t. the measure X_{τ} with $\tau \in \mathcal{M}$, if $|G| = X_{\tau}(L)$. In the case of the measure Xc , we speak of a *minimal X-grammar* generating a finite language.

Finally, we show that grammars that generate non-trivial uniform languages do not contain ε -productions. To this end, we first need the following result on the length of words generated by some nonterminal from a grammar that describes a finite uniform language.

Lemma 1. *Let $X \in \Gamma$ and $G = (N, \Sigma, P, S)$ be a minimal X-grammar generating a finite and uniform language. Then, for all $A \in N$, all words occurring in the set $\{w \in \Sigma^* \mid A \Rightarrow_{\mathcal{C}}^* w\}$ have the same length.*

An easy consequence of the previous lemma is the next theorem.

Theorem 2. *Let $X \in \Gamma$ and $G = (N, \Sigma, P, S)$ be a minimal X-grammar generating a finite and uniform language satisfying $L(G) \neq \{\varepsilon\}$. Then G is ε -free, i.e., P does not contain any rule of the form $A \rightarrow \varepsilon$, for all $A \in N$.*

3 Some Bounds on the Various X-Complexities

First, we prove some upper bounds for the finite variants of the introduced grammatical measures. We obtain the following results:

Theorem 3. *Let $L \subseteq \Sigma^{\leq \ell}$ be a finite language over the alphabet Σ . Then, for $X \in \Gamma$ and $Y \in \{\text{REG}, \text{LIN}\}$, we have*

1. $\text{Xc}(L) \leq \ell + 1$ if $|\Sigma| = 1$, and $\text{Xc}(L) \leq (|\Sigma|^{\ell+1} - 1)/(|\Sigma| - 1)$, otherwise,
2. $\text{CFcc}(L) \leq |\Sigma| + 2$ and $\text{Ycc}(L) \leq \ell + 1$ if $|\Sigma| = 1$, and $\text{Ycc}(L) \leq (|\Sigma| + 1) \cdot \ell$, otherwise, and
3. $\text{Xsc}(L) = 1$ if L is non-empty, and $\text{Xsc}(L) = 0$, otherwise.

Proof. We argue as follows:

1. Every finite language L can be generated by a grammar of type $X \in \Gamma$ by simply listing all words from L . Since there are at most $\sum_{i=0}^{\ell} |\Sigma|^i$ words of length at most ℓ in L , the upper bounds of $\ell + 1$ and $(|\Sigma|^{\ell+1} - 1)/(|\Sigma| - 1)$ follow for the cases $|\Sigma| = 1$ and $|\Sigma| \geq 2$, respectively.
2. Consider the context-free grammar $G' = (\{A, S'\}, \Sigma, P', S')$, where P' consists of the productions $S' \rightarrow A^\ell$ and $A \rightarrow a$, for $a \in \Sigma \cup \{\varepsilon\}$. Clearly, we have $L(G') = \Sigma^{\leq \ell}$ and $|G'| = |\Sigma| + 2$, i.e., $\text{CFcc}(\Sigma^{\leq \ell}) \leq |\Sigma| + 2$. By assumption, $L \subseteq \Sigma^{\leq \ell}$. Thus, every grammar generating $\Sigma^{\leq \ell}$ automatically covers the language L . For $|\Sigma| \geq 2$, consider the regular grammar $G = (N, \Sigma, P, S)$, where $N = \{A_1, A_2, \dots, A_\ell\}$ with $S = A_1$ and

$$P = \{A_i \rightarrow aA_{i+1} \mid a \in \Sigma \text{ and } 1 \leq i \leq \ell - 1\} \cup \{A_\ell \rightarrow a \mid a \in \Sigma\} \\ \cup \{A_i \rightarrow \varepsilon \mid 1 \leq i \leq \ell\}.$$

Obviously, $L(G) = \Sigma^{\leq \ell}$ and $|G| = (|\Sigma| + 1) \cdot \ell$. In the case that $|\Sigma| = 1$, we simply list all $\ell + 1$ words occurring in $\Sigma^{\leq \ell}$.

3. Assume that $\Sigma = \{a_1, a_2, \dots, a_n\}$ and consider the language $\{(a_1 a_2 \dots a_n)^\ell\}$, which is generated by $G = (\{S\}, \Sigma, \{S \rightarrow (a_1 a_2 \dots a_n)^\ell\}, S)$, a regular grammar with a single production rule. Clearly, we have $L \subseteq \{(a_1 a_2 \dots a_n)^\ell\}$, for all nonempty languages $L \subseteq \Sigma^{\leq \ell}$. Thus, $\text{Xsc}(L) \leq 1$. Since any grammar with empty production set can only generate the empty language, we also have $\text{Xsc}(L) \geq 1$. In case $L = \emptyset$, we obviously have $\text{Xsc}(L) = 0$. \square

What about lower bounds for these measures? Observe that $\text{Xsc}(L) = 1$ is already a lower bound result. In the seminal paper [4] on concise description of finite languages by different types of grammars, certain languages have been identified that require at least a *polynomial* number of productions. The proofs of these results are based on [4, Lemma 2.1] which states some easy facts about *minimal* context-free grammars:

Lemma 4. *Let $G = (N, \Sigma, P, S)$ be a minimal context-free grammar for the finite language L . Then, for every nonterminal $A \in N \setminus \{S\}$, there are strings α_1 and α_2 with $\alpha_1, \alpha_2 \in (N \cup \Sigma)^*$ and $\alpha_1 \neq \alpha_2$ such that $A \rightarrow \alpha_1$ and $A \rightarrow \alpha_2$ are*

in P . Moreover, for every $A \in N \setminus \{S\}$, the set $L_A(G) = \{w \in \Sigma^* \mid A \Rightarrow_G^* w\}$ contains at least two words and there is no derivation of the form $A \Rightarrow_G^+ \alpha A \beta$ with $\alpha, \beta \in (N \cup \Sigma)^*$. Finally, for every $A \in N \setminus \{S\}$, there are $u_1, u_2, v_1, v_2 \in \Sigma^*$ such that $u_1 A u_2 \neq v_1 A v_2$ as well as $S \Rightarrow_G^* u_1 A u_2$ and $S \Rightarrow_G^* v_1 A v_2$.

Our first candidate with a large X -complexity is the language of all even palindromes (with middle marker) $P_n = \{w\$w^R \mid w \in \{a, b\}^{\leq n}\}$. We show that any regular grammar generating this language needs at least an exponential number of productions.

Theorem 5. *Let $n \geq 1$. Then $\text{REGc}(P_n) \geq 2^n$.*

Proof. In the proof, we will use the following result from [4, Lemma 2.2]: let $G = (N, \Sigma, P, S)$ be a context-free grammar generating a finite language. Then there is a context-free grammar $G_{\max} = (N_{\max}, \Sigma, P_{\max}, S)$ such that $N_{\max} \subseteq N$, $P_{\max} \subseteq P$, and $L(G_{\max}) = L_{\max}$, where L_{\max} is the subset of $L(G)$ consisting of the words of maximal length. In light of this result, it suffices to show that $\text{REGc}(P'_n) \geq 2^n$, for the language $P'_n = \{w\$w^R \mid w \in \{a, b\}^n\}$. To this end, assume that $\Sigma = \{a, b, \$\}$ and $G = (N, \Sigma, P, S)$ is a minimal regular grammar generating P'_n that contains a nonterminal $A \in N \setminus \{S\}$. By Lemma 4, there are derivations $S \Rightarrow_G^* u_1 A$ and $S \Rightarrow_G^* u_2 A$ with $u_1, u_2 \in \Sigma^*$ and $u_1 \neq u_2$ as well as $v_1, v_2 \in \Sigma^*$ with $A \Rightarrow_G^* v_1$, $A \Rightarrow_G^* v_2$, and $v_1 \neq v_2$. Note that we must have both $|u_1| = |u_2|$ and $|v_1| = |v_2|$, for otherwise we would be able to derive words w_1 and w_2 with $|w_1| \neq |w_2|$, but P'_n only contains words of the same length. Since $v_1 \neq v_2$ and $|v_1| = |v_2|$, it follows that both $v_1 \neq \varepsilon$ and $v_2 \neq \varepsilon$. Let $w \in \{a, b\}^n$ be arbitrary. We distinguish the following cases:

1. Suppose $u_1 \in \{w\$\}\{a, b\}^*$. Then we must have $u_1 = w_1 w_2 \$w_2^R$, where $w = w_1 w_2$ and both $v_1 \in \{a, b\}^*$ and $v_2 \in \{a, b\}^*$ holds. Assume, w.l.o.g., that $v_1 = w_1^R$. Since $v_1 \neq v_2$, it follows that $v_2^R \neq w_1$. Thus, $u_1 v_2 \notin P'_n$. Contradiction.
2. Suppose $u_1 \in \{a, b\}^*$. Then we must have $v_1, v_2 \in \{a, b\}^* \{ \$w^R \}$. Assume $w = u_1 w_2$, for some $w_2 \in \{a, b\}^*$ and, w.l.o.g., $v_1 = w_2 \$w_2^R u_1^R$. Since $v_1 \neq v_2$, it follows that $v_2 = w'_2 \$w_2^R u_1^R$ with $w'_2 \neq w_2$. Thus, $u_1 v_2 \notin P'_n$. Contradiction.

Consequently, we have $N = \{S\}$ and so the only way to generate the language P'_n minimally with a regular grammar is to list all of its words using S . \square

For linear context-free and context-free grammars, one observes that both measures $\text{LINc}(P_n)$ and $\text{Cfc}(P_n)$ are at most linear, as witnessed by the linear context-free grammar $G = (N, \Sigma, P, S)$ with $N = \{S_0, S_1, \dots, S_n\}$, $\Sigma = \{a, b, \$\}$, start symbol $S = S_0$, and the productions

$$P = \{S_i \rightarrow aS_{i+1}a, S_i \rightarrow bS_{i+1}b, S_i \rightarrow S_{i+1} \mid 0 \leq i \leq n - 1\} \cup \{S_n \rightarrow \$\},$$

satisfying $L(G) = P_n$, for $n \geq 1$.

Using similar arguments as in the proof of Theorem 5, one can show that the triple language $T_n = \{w\$w\#w \mid w \in \{a, b\}^n\}$ has large X -complexity, for all grammar types X with $X \in \Gamma$. A detailed proof of this fact can be found in [9].

Theorem 6. *Let $X \in \Gamma$ and $n \geq 1$. Then $\text{Xc}(T_n) = 2^n$.*

4 Relating Finite and Infinite Complexity Measures

In this section, we will consider several different complexity measures on finite languages and relate them according to a group of relations that vary in strength. By the very nature of these relations, one can distinguish two main categories: the first category fixes the measure type $\tau \in \mathcal{M}$ and then compares the different grammar types in Γ with each other w.r.t. the measure type τ —see, e.g., (2). The second category swaps the roles of measure and grammar type, i.e., some grammar type $X \in \Gamma$ is fixed and then the different measure types in \mathcal{M} are compared with each other w.r.t. the grammar type X .

4.1 Relating Grammar Types

Now, we define several different relations on the grammar types in Γ w.r.t. some fixed measure type from \mathcal{M} . In this way, we classify the difference between different grammar types w.r.t. the same measure type. This is similar to the notion introduced in [6] for the nonterminal complexity and thus leads to a certain kind of production complexity hierarchy.

Let $X, Y \in \Gamma$ and $\tau \in \mathcal{M}$. Then we write

- $X \leq_\tau Y$ if and only if $X_\tau(L) \leq Y_\tau(L)$, for all finite languages L ;
- $X \leq_\tau^1 Y$ if and only if there is a constant c such that $X_\tau(L) \leq Y_\tau(L) + c$, for all finite languages L , and there is a sequence of finite languages $(L_i)_{i \geq 0}$ such that $Y_\tau(L_i) - X_\tau(L_i) \geq i$;
- $X \leq_\tau^2 Y$ if and only if there is a constant c such that $X_\tau(L) \leq Y_\tau(L) + c$, for all finite languages L , and there is a sequence of finite languages $(L_i)_{i \geq 0}$ such that

$$\lim_{i \rightarrow \infty} \frac{X_\tau(L_i)}{Y_\tau(L_i)} = 0;$$

- $X \leq_\tau^3 Y$ if and only if there is a constant c such that $X_\tau(L) \leq Y_\tau(L) + c$, for all finite languages L , and there is no function $f: \mathbb{N} \rightarrow \mathbb{N}$ such that $Y_\tau(L) \leq f(X_\tau(L))$, for all finite languages L .

Clearly, $X \leq_\tau^3 Y$ implies $X \leq_\tau^2 Y$, which, in turn, implies $X \leq_\tau^1 Y$. Moreover, $X \leq_\tau^3 Y$ holds if the first condition of its definition is satisfied and there is a sequence $(L_i)_{i \geq 0}$ of finite languages such that $X_\tau(L_i) \leq k$, for some constant k and $Y_\tau(L_i) \geq i$. We write $X =_\tau Y$ if both $X \leq_\tau Y$ and $Y \leq_\tau X$ hold.

Now, we are going to relate the different grammar types in Γ w.r.t. the finite complexity measure types under investigation. As already mentioned earlier, the relation $\text{CF} \leq_\tau \text{LIN} \leq_\tau \text{REG}$ holds by definition, for all $\tau \in \mathcal{M}$. The following result was shown in [4]:

Theorem 7. *It holds that $\text{CF} \leq_c^2 \text{LIN} \leq_c^2 \text{REG}$.*

Next, we show that we can only obtain $\text{CF} \leq_\tau^3 \text{LIN}$ and $\text{CF} \leq_\tau^3 \text{REG}$, for $\tau \in \{\text{c}, \text{cc}\}$. As a prerequisite, we need the following result from [4].

Lemma 8. *Let G be a linear grammar generating a finite language with $|G| \geq 1$. Then $|L(G)| \leq 2^{|G|-1}$ and $|G| \geq \log |L(G)| + 1$.*

Now, we are ready for the following theorem.

Theorem 9. *Let $\tau \in \{\text{c}, \text{cc}\}$. Then*

1. $\text{CF} \leq_{\tau}^3 X$, for $X \in \{\text{REG}, \text{LIN}\}$, and
2. $\text{LIN} \not\leq_{\tau}^3 \text{REG}$ and $\text{REG} \not\leq_{\tau}^3 \text{LIN}$.

Proof. 1. Let $L_n = \{a, b\}^{\leq n}$, for $n \geq 0$. It was shown in [4], that $\text{CFc}(L_n) \leq 4$ and $\text{LINc}(L_n) \geq n + 1$. From Lemma 8, it follows that also $\text{LINcc}(L_n) \geq n + 1$. Moreover, it holds that $\text{CFcc}(L_n) \leq \text{CFc}(L_n) \leq 4$. Consequently, $\text{CF} \leq_{\tau}^3 \text{LIN}$. Since $\text{REG}\tau(L_n) \geq \text{LIN}\tau(L_n) \geq n + 1$, we immediately get $\text{CF} \leq_{\tau}^3 \text{REG}$.

2. Let L be an arbitrary finite language and G be a minimal linear grammar with $L(G) = L$. Clearly, $\text{REG}\tau(L) \leq |L| = |L(G)|$. From Lemma 8, it follows that $|L(G)| \leq 2^{\text{LIN}\tau(L)-1}$. Thus, $\text{REG}\tau(L) \leq |L(G)| \leq 2^{\text{LIN}\tau(L)-1}$. The function $f: \mathbb{N} \rightarrow \mathbb{N}$ with $x \mapsto 2^{x-1}$ fulfills $\text{REG}\tau(L') \leq f(\text{LIN}\tau(L'))$, for all finite languages L' , i.e., $\text{LIN} \not\leq_{\tau}^3 \text{REG}$. Since $\text{LIN}\tau(L') \leq \text{REG}\tau(L')$, setting $f = \text{id}_{\mathbb{N}}$ yields $\text{REG} \not\leq_{\tau}^3 \text{LIN}$. \square

For the X -scattered-complexity, we have the following situation:

Theorem 10. *Let $X, Y \in \Gamma$. Then $\text{REG} =_{\text{sc}} \text{LIN} =_{\text{sc}} \text{CF}$, but $X \not\leq_{\text{sc}}^i Y$, for $i \in \{1, 2, 3\}$.*

It remains to relate the different grammar types in Γ w.r.t. the infinite complexity measure types $\text{c}_{\infty}, \text{cc}_{\infty}, \text{sc}_{\infty}$. For the infinite X -complexity we have:

Theorem 11. *Let $n \geq 1$. Then $\text{REGc}_{\infty}(P_n) = \Omega(2^n)$.*

Proof. The idea of the proof is as follows: let $G = (N, \Sigma, P, S)$ be a regular grammar that is a witness for $\text{REGc}_{\infty}(P_n)$. Then we construct a regular grammar generating the finite set $L(G) \cap \Sigma^{\leq 2n+1}$. Since this language is equal to P_n , we can apply Theorem 5 in order to obtain a lower bound on $|G|$. Although G may contain ε -productions, we can safely assume that G does not contain productions with a right hand-side longer than $2n + 2$, since none of these productions generates a word of length less than or equal to $2n + 1$ even with erasing rules.

To keep the presentation simple, assume for a moment that the grammar G is in 2-normal form.² Then we apply the triple construction—see, e.g., [10]—to the grammar G and the *nondeterministic* finite automaton $A = (Q, \Sigma, \delta, q_0, F)$ with $Q = \{0, 1, \dots, 2n + 1\}$, initial state $q_0 = 0$, final state set $F = \{2n + 1\}$, and the transition function $\delta(i, a) = \{i + 1, 2n + 1\}$, for $0 \leq i < 2n + 1$ and $a \in \Sigma$. Then the regular grammar $G' = (N', \Sigma, P', S')$ with $N' = Q \times N \times Q$, start symbol $S' = [0, S, 2n + 1]$, and the productions

$$P' = \{ [i, A, j] \rightarrow a \mid A \rightarrow a \in P, a \in \Sigma \cup \{\varepsilon\}, \text{ and } j \in \delta(i, a) \} \\ \cup \{ [i, A, k] \rightarrow a[j, B, k] \mid A \rightarrow aB \in P, a \in \Sigma \cup \{\varepsilon\}, \text{ and } j \in \delta(i, a) \}$$

² A regular grammar $G = (N, \Sigma, P, S)$ is in 2-normal form if all productions in P are of the form $A \rightarrow a$ and $A \rightarrow aB$, where $A, B \in N$ and $a \in \Sigma \cup \{\varepsilon\}$.

generates the finite language $L(G) \cap \Sigma^{\leq 2n+1}$. Observe that $|G'| \leq |G| \cdot (2n+2)^3$, because the productions of the form $A \rightarrow aB$ from P increase $|G'|$ the most.

Now, we are ready to prove the stated claim. Assume to the contrary that we have $|G| = \text{REG}_{c_\infty}(P_n) = o(2^n)$, for a regular grammar G . Then we transform G into an equivalent regular grammar in 2-normal form and apply the above described intersection construction. Let G' refer to the result of these construction steps. The transformation into 2-normal form is done by simply splitting the right hand-side of each production into a sequence of productions of the appropriate form. This increases $|G|$ by at most a factor of $2n+2$. Together with the increase of productions by the triple construction by a factor of at most $(2n+2)^3$, we conclude that $|G'| = o(2^n)$, because $o(2^n) \cdot (2n+2)^4 = o(2^n)$. Since the regular grammar G' with $|G'| = o(2^n)$ generates P_n , we get a contradiction to Theorem 5. Thus, we obtain $\text{REG}_{c_\infty}(P_n) = \Omega(2^n)$. \square

The taxonomy of the next theorem applies to the infinite exact X -complexity.

Theorem 12. *It holds that $\text{CF} \leq_{c_\infty}^i \text{REG}$ and $\text{LIN} \leq_{c_\infty}^i \text{REG}$, for all $i \in \{1, 2, 3\}$.*

Finally, both the infinite X -cover and the infinite X -scattered-complexity do not classify according to the used taxonomy.

Theorem 13. *Let $X, Y \in \Gamma$ and $\tau \in \{\text{cc}_\infty, \text{sc}_\infty\}$. Then we have $X \not\leq_\tau^i Y$, for all $i \in \{1, 2, 3\}$.*

We have to leave open the question whether $\text{LIN} \leq_{\text{cc}}^i \text{REG}$ and $\text{CF} \leq_{c_\infty}^j \text{LIN}$, for $i \in \{1, 2\}$ and $j \in \{1, 2, 3\}$, holds.

4.2 Relating Complexity Measure Types

Now, we introduce relations for measuring the difference between different measure types w.r.t. some fixed grammar type that we consider in this paper. Therefore, for $\tau, \sigma \in \mathcal{M}$ and $X \in \Gamma$, we similarly define the relations $\tau \leq_X \sigma$, $\tau \leq_X^1 \sigma$, $\tau \leq_X^2 \sigma$, and $\tau \leq_X^3 \sigma$ as in the beginning of the previous subsection. For instance, $\tau \leq_X^1 \sigma$ if and only if there is a constant c such that $X\tau(L) \leq X\sigma(L) + c$, for all finite languages L , and there is a sequence of finite languages $(L_i)_{i \geq 0}$ such that $X\sigma(L_i) - X\tau(L_i) \geq i$.

Clearly, the following chain of implications holds: $\tau \leq_X^3 \sigma$ implies $\tau \leq_X^2 \sigma$, which, in turn, implies $\tau \leq_X^1 \sigma$. Moreover, $\tau \leq_X^3 \sigma$ holds if the first condition of its definition is satisfied and there is a sequence $(L_i)_{i \geq 0}$ of finite languages such that $X\tau(L_i) \leq c$, for some constant c , and $X\sigma(L_i) \geq i$.

We start with comparing the finite X - with the infinite X -measures. Except for the X -scattered-complexity, the infinite versions are more succinct than their finite counterparts.

Lemma 14. *Let $X \in \Gamma$. Then (i) $c_\infty \leq_X c$ and (ii) $\text{cc}_\infty \leq \text{cc}$, but we have (iii) $\text{sc} \leq_X \text{sc}_\infty$.*

Proof. The first relation follows by definition. Let $L \subseteq \Sigma^{\leq \ell}$. Assume that G is a witness for $\text{Xc}(L)$, i.e., G generates a finite language, $L = L(G)$, and $\text{Xc}(L) = |G|$. But then also $L = L(G) \cap \Sigma^{\leq \ell}$, which implies $\text{Xc}_\infty(L) \leq \text{Xc}(L)$. A similar argumentation applies to the second relation. For the third relation, we argue as follows: in Theorem 3, it was shown that $\text{Xsc}(L) = 1$ if L is non-empty and $\text{Xsc}(L) = 0$ if $L = \emptyset$. In the latter case, we also have $\text{Xsc}_\infty(L) = 0$. Thus, we conclude that $\text{Xsc}_\infty(L) \leq \text{Xsc}(L)$, for every finite language L . \square

It is worth mentioning that in the previous lemma, the argumentation used in the proof of the first two relations does not apply to the third one. This is seen as follows: consider the uniform finite language $L = \{a, b\}$. Then the regular grammar $G = (\{S\}, \{a, b\}, \{S \rightarrow ab\}, S)$ witnesses $\text{Xsc}(L) = 1$, because L is a scattered sublanguage of $\{ab\}$, i.e., $L \leq \{ab\} = L(G)$. But $L(G) \cap \{a, b\}^{\leq 1} = \emptyset$. Thus, $|G|$ cannot be used as an upper bound for $\text{Xsc}_\infty(L)$ and hence $\text{sc}_\infty \leq_X \text{sc}$ does not hold in general. Thus, we conclude:

Corollary 15. *Let $X \in \Gamma$. Then $\text{sc}_\infty \not\leq_X \text{sc}$.*

Next, we compare the remaining X -complexities. Observe that $L_1 = L_2$ implies $L_1 \subseteq L_2$, which, in turn, implies $L_1 \leq L_2$. As an easy consequence, we deduce that the (infinite) X -scattered-complexity is more succinct than the (infinite) X -cover-complexity and it is also easy to see that the (infinite) X -cover-complexity is more succinct than the (infinite) X -complexity. We summarize:

Lemma 16. *Let $X \in \Gamma$. Then (i) $\text{sc} \leq_X \text{cc} \leq_X \text{c}$ and (ii) $\text{sc}_\infty \leq_X \text{cc}_\infty \leq_X \text{c}_\infty$.*

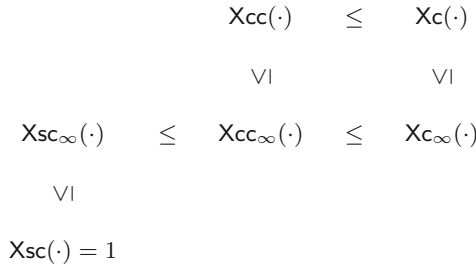


Fig. 1. Relations between the different grammatical complexity measures on finite languages. In this figure, a \leq -relation is of type (i, j, k) , for $i, j, k \in \{1, 2, 3\}$, if (i) \leq_{CFG}^i , (ii) \leq_{LIN}^j , and (iii) \leq_{REG}^k hold for the appropriate X -measures. All \leq -relations are of type $(3, 3, 3)$, except $\text{Xcc} \leq \text{Xc}$, $\text{Xcc}_\infty \leq \text{Xcc}$, $\text{Xsc} \leq \text{Xsc}_\infty$, and $\text{Xsc}_\infty \leq \text{Xcc}_\infty$, which are of types $(3, 2, 2)$, $(-, 3, 3)$, $(-, -, -)$, and $(-, -, -)$, respectively; the $-$ sign means that it cannot be classified by the taxonomy.

The obtained \leq -relations are visualized in Fig. 1. For the measures cc and c_∞ we show incomparability w.r.t. the \leq -relation. Before we can prove this, we need two prerequisites. The first one is a lower bound on T_n w.r.t. the infinite

X -complexity, which reads as follows—the proof is similar to the proof of Theorem 11 with the slight modification that we use 2-Greibach normal form³ instead of 2-normal form, since we cannot apply the reasoning regarding the absence of ε -productions to context-free grammars in case of the c_∞ -measure. The change to 2-Greibach normal form induces the fourth root in the lower bound:

Theorem 17. *Let $X \in \Gamma$. Then $Xc_\infty(T_n) = \Omega(2^{n/4})$.*

The second prerequisite is an exact complexity bound for the language $\Sigma^{\leq \ell}$ w.r.t. the finite X -cover-complexity.

Lemma 18. *Let $X \in \Gamma$ and $\ell \geq 2$. Then $CFcc(\Sigma^{\leq \ell}) = |\Sigma| + 2$.*

Now, we are ready for the incomparability results:

Theorem 19. *Let $X \in \Gamma$. Then (i) $cc \not\leq_X c_\infty$ and (ii) $c_\infty \not\leq_X cc$.*

Proof. For the proof of (i), observe that the grammar $G = (N, \Sigma, P, S)$ with the productions $P = \{S \rightarrow aS, S \rightarrow bS, S \rightarrow \varepsilon\}$ shows that $Xc_\infty(\{a, b\}^{\leq n}) \leq 3$. By Lemma 8, we have $n + 1 \leq Ycc(\{a, b\}^{\leq n})$, for $n \geq 3$ and $Y \in \{\text{REG}, \text{LIN}\}$. As a consequence, $Yc_\infty(\{a, b\}^{\leq n}) < Ycc(\{a, b\}^{\leq n})$. From Lemma 18, it follows that $CFcc(\{a, b\}^{\leq n}) \geq 4$, i.e., $CFc_\infty(\{a, b\}^{\leq n}) < CFcc(\{a, b\}^{\leq n})$. Next, we prove (ii). By Theorem 17, $Xc_\infty(T_n) = \Omega(2^{n/4})$ and, by Theorem 3, we have $XccT_n \leq 15n + 10$. Thus, $XccT_n < Xc_\infty(T_n)$, for large enough n . \square

In the remainder of this subsection, we classify the relations between the X -measures under consideration according to the taxonomy from [6].

Theorem 20. *Let $X \in \Gamma$ and $Y \in \{\text{REG}, \text{LIN}\}$. Then*

1. $cc \leq_{CFG}^3 c$ and $cc \leq_Y^2 c$, but $\tau \not\leq_Y^3 \sigma$, for all $\tau, \sigma \in \{c, cc\}$ with $\tau \neq \sigma$,
2. $sc \leq_X^3 c$,
3. $sc \leq_Y^3 cc$, but $\tau \not\leq_{CFG}^i \sigma$, for all $\tau, \sigma \in \{sc, cc\}$ with $\tau \neq \sigma$ and all $i \in \{1, 2, 3\}$,
4. $c_\infty \leq_X^3 c$,
5. $\tau \leq_X^3 c$, for all $\tau \in \{cc_\infty, sc_\infty\}$,
6. $cc_\infty \leq_Y^3 cc$, but $\tau \not\leq_{CFG}^i \sigma$, for $\tau, \sigma \in \{cc, cc_\infty\}$ with $\tau \neq \sigma$ and $i \in \{1, 2, 3\}$,
7. $cc \leq_{CFG}^3 c_\infty$, and
8. $\tau \leq_X^3 c_\infty$, for all $\tau \in \{sc, sc_\infty, cc_\infty\}$.

Proof. We only prove the first statement. The remaining results can be shown with similar arguments. Let $L_n = \{a^j b^j c^j \mid 1 \leq j \leq n\}$. In [2], it was shown that $CFc(L_n) = n$. On the other hand, by Theorem 3, we have $CFcc(L_n) \leq 5$.

Let $Y \in \{\text{REG}, \text{LIN}\}$ and $T_n = \{w\$w\#w \mid w \in \{a, b\}^n\}$. By Theorem 6, we have $Yc(T_n) = 2^n$. On the other hand, we have $Ycc(T_n) \leq 15n + 10$ by Theorem 3.

³ A context-free grammar $G = (N, \Sigma, P, S)$ is in 2-Greibach normal form if all productions in P are of the form $A \rightarrow a$, $A \rightarrow aB$, $A \rightarrow aBC$, or $S \rightarrow \varepsilon$, where $A \in N$, $a \in \Sigma$, and $B, C \in N \setminus \{S\}$. The transformation increases the number of productions by at most a polynomial of fourth degree [1].

Hence, $\text{cc} \leq_Y^2 c$. Finally, let L be an arbitrary finite language and $f: \mathbb{N} \rightarrow \mathbb{N}$ a function defined by $x \mapsto 2^{x-1}$. Moreover, let G be a minimal Y -grammar with $L(G) \supseteq L$. By Lemma 8, we know that $|L| \leq |L(G)| \leq 2^{\text{Ycc}(L)-1}$. It holds that $\text{Yc}(L) \leq |L|$. Thus,

$$\text{Yc}(L) \leq |L| \leq |L(G)| \leq 2^{\text{Ycc}(L)-1} = f(\text{Ycc}(L)).$$

Hence, $\text{cc} \not\leq_Y^3 c$. If we set $f = \text{id}_{\mathbb{N}}$, it follows that $c \not\leq_Y^3 \text{cc}$. \square

Finally, we list some incomparability results.

Theorem 21. *Let $X \in \Gamma$ and $Y \in \{\text{REG}, \text{LIN}\}$. Then, for $i \in \{1, 2, 3\}$, we have*

1. $\tau \not\leq_X^i \sigma$, for every $\tau, \sigma \in \{\text{cc}_\infty, \text{sc}_\infty\}$ with $\tau \neq \sigma$,
2. $\tau \not\leq_X^i \sigma$, for every $\tau, \sigma \in \{\text{cc}, \text{sc}_\infty\}$ with $\tau \neq \sigma$,
3. $\tau \not\leq_X^i \sigma$, for every $\tau, \sigma \in \{\text{sc}, \text{sc}_\infty\}$ with $\tau \neq \sigma$,
4. $\tau \not\leq_X^i \sigma$, for every $\tau, \sigma \in \{\text{sc}, \text{cc}_\infty\}$ with $\tau \neq \sigma$, and
5. $\tau \not\leq_Y^i \sigma$, for every $\tau, \sigma \in \{\text{cc}, \text{c}_\infty\}$ with $\tau \neq \sigma$.

Proof. We only prove the first statement. The remaining results can be shown with similar arguments. First note that $L \subseteq \Sigma^* \cap \Sigma^{\leq \ell}$ as well as $L \leq \Sigma^* \cap \Sigma^{\leq \ell}$, for all finite languages L over Σ with $\ell = \max\{|w| \mid w \in L\}$. The universal language Σ^* , for $\Sigma = \{a_1, a_2, \dots, a_n\}$, can be produced with the following regular productions: $S \rightarrow a_1 S \mid a_2 S \mid \dots \mid a_n S \mid \varepsilon$. Thus, $\text{Xcc}_\infty(L) \leq |\Sigma| + 1$ and $\text{Xsc}_\infty(L) \leq |\Sigma| + 1$, for all finite languages L over Σ . This, however, means that $\tau \leq_X^i \sigma$ does *not* hold for $\tau, \sigma \in \{\text{cc}_\infty, \text{sc}_\infty\}$ and all $i \in \{1, 2, 3\}$. \square

References

1. Blum, N., Koch, R.: Greibach normal form transformation revisited. *Inform. Comput.* **150**(1), 112–118 (1999)
2. Bucher, W.: A note on a problem in the theory of grammatical complexity. *Theoret. Comput. Sci.* **14**(3), 337–344 (1981)
3. Bucher, W., Maurer, H.A., Culik II, K.: Context-free complexity of finite languages. *Theoret. Comput. Sci.* **28**(3), 277–285 (1983)
4. Bucher, W., Maurer, H.A., Culik II, K., Wotschke, D.: Concise description of finite languages. *Theoret. Comput. Sci.* **14**(3), 227–246 (1981)
5. Câmpeanu, C., Sântean, N., Yu, S.: Minimal cover-automata for finite languages. *Theoret. Comput. Sci.* **267**(1–2), 3–16 (2001)
6. Dassow, J., Păun, Gh.: *Regulated Rewriting in Formal Language Theory*. EATCS Monographs in Theoretical Computer Science, vol. 18. Springer, Heidelberg (1989)
7. Eberhard, S., Hetzl, S.: On the compressibility of finite languages and formal proofs. *Inform. Comput.* **259**(2), 191–213 (2018)
8. Hetzl, S., Wolfsteiner, S.: Cover complexity of finite languages. In: Konstantinidis, S., Pighizzini, G. (eds.) *DCFS 2018*. LNCS, vol. 10952, pp. 139–150. Springer, Cham (2018)
9. Gruber, H., Holzer, M., Wolfsteiner, S.: *On Minimal Grammar Problems for Finite Languages* (2018, Submitted for publication)
10. Wood, D.: *Theory of Computation*. Wiley, New York (1987)