

# Dissecting power of intersection of two context-free languages

Josef Rukavicka

Faculty of Nuclear Sciences and Physical Engineering, CTU in Prague, Czech Republic

revisions 8<sup>th</sup> Feb. 2022, 26<sup>th</sup> Oct. 2022, 26<sup>th</sup> Jan. 2023, 6<sup>th</sup> June 2023; accepted 6<sup>th</sup> July 2023.

---

We say that a language  $L$  is *constantly growing* if there is a constant  $c$  such that for every word  $u \in L$  there is a word  $v \in L$  with  $|u| < |v| \leq c + |u|$ . We say that a language  $L$  is *geometrically growing* if there is a constant  $c$  such that for every word  $u \in L$  there is a word  $v \in L$  with  $|u| < |v| \leq c|u|$ . Given two infinite languages  $L_1, L_2$ , we say that  $L_1$  *dissects*  $L_2$  if  $|L_2 \setminus L_1| = \infty$  and  $|L_1 \cap L_2| = \infty$ . In 2013, it was shown that for every constantly growing language  $L$  there is a regular language  $R$  such that  $R$  dissects  $L$ .

In the current article we show how to dissect a geometrically growing language by a homomorphic image of intersection of two context-free languages.

Consider three alphabets  $\Gamma, \Sigma$ , and  $\Theta$  such that  $|\Sigma| = 1$  and  $|\Theta| = 4$ . We prove that there are context-free languages  $M_1, M_2 \subseteq \Theta^*$ , an erasing alphabetical homomorphism  $\pi : \Theta^* \rightarrow \Sigma^*$ , and a nonerasing alphabetical homomorphism  $\varphi : \Gamma^* \rightarrow \Sigma^*$  such that: If  $L \subseteq \Gamma^*$  is a geometrically growing language then there is a regular language  $R \subseteq \Theta^*$  such that  $\varphi^{-1}(\pi(R \cap M_1 \cap M_2))$  dissects the language  $L$ .

**Keywords:** Dissecting of infinite languages, Context-free languages, Intersection of context-free languages

---

## 1 Introduction

In the theory of formal languages, the regular and the context-free languages constitute a fundamental concept that attracted a lot of attention in the past several decades.

In contrast to regular languages, the context-free languages are closed neither under intersection nor under complement. The intersection of context-free languages have been systematically studied; see for instance [4, 6, 9]. Let  $\text{CFL}_k$  denote the family of all languages such that for each  $L \in \text{CFL}_k$  there are  $k$  context-free languages  $L_1, L_2, \dots, L_k$  with  $L = \bigcap_{i=1}^k L_i$ . For each  $k$ , it has been shown that there is a language  $L \in \text{CFL}_{k+1}$  such that  $L \notin \text{CFL}_k$ . Thus the  $k$ -intersections of context-free languages form an infinite hierarchy in the family of all formal languages lying between context-free and context sensitive languages [6].

Dissection of infinite languages belongs to the topics of the theory of formal languages that have been studied in recent years. Let  $L_1$  and  $L_2$  be infinite languages. We say that  $L_1$  *dissects*  $L_2$  if  $|L_2 \setminus L_1| = \infty$  and  $|L_1 \cap L_2| = \infty$ . Let  $\mathcal{C}$  be a family of languages. We say that a language  $L_2$  is  $\mathcal{C}$ -dissectible if there is  $L_1 \in \mathcal{C}$  such that  $L_1$  dissects  $L_2$ . Let REG denote the family of regular languages. In [10]

the REG-dissectibility has been investigated. Several families of REG-dissectible languages have been presented. Moreover, it has been shown that there are infinite languages that cannot be dissected with a regular language. Also some open questions for REG-dissectibility can be found in [10]. For example, it is not known if the complement of a context-free language is REG-dissectible.

Given two countable sets  $A$  and  $B$ , we write  $A \subseteq_{ae} B$  if  $|A - B| < \infty$  and we write  $A =_{ae} B$  if both  $A \subseteq_{ae} B$  and  $B \subseteq_{ae} A$  hold. The subscript “ae” stands for “almost everywhere”. We say that  $A$  covers  $B$  with an infinite margin (or  $A$   $i$ -covers  $B$ , in short) if both  $B \subseteq A$  and  $A \neq_{ae} B$  hold. We represent a pair of languages  $A$  and  $B$  such that  $A$   $i$ -covers  $B$  by  $i(A, B)$ . A language  $C$  is said to separate  $i(B, A)$  with infinite margins (or  $i$ -separates  $i(B, A)$ , in short) if  $B \subseteq C \subseteq A$ ,  $A \neq_{ae} C$ , and  $B \neq_{ae} C$ . In addition, given two language families  $\mathcal{A}, \mathcal{B}$  we define

$$i(\mathcal{B}, \mathcal{A}) = \{i(B, A) \mid A \in \mathcal{A} \text{ and } B \in \mathcal{B} \text{ and } A \text{ } i\text{-covers } B\}.$$

We say that a language family  $\mathcal{C}$   $i$ -separates  $i(\mathcal{B}, \mathcal{A})$  if for every  $i(B, A) \in i(\mathcal{B}, \mathcal{A})$  there is a language  $C \in \mathcal{C}$  such that  $C$   $i$ -separates  $i(B, A)$ .

Given  $\mathcal{A}$  and  $\mathcal{B}$  be any two language families, let

$$\mathcal{A} - \mathcal{B} = \{A - B \mid A \in \mathcal{A} \text{ and } B \in \mathcal{B}\}.$$

In [10], a connection between REG-dissectibility and  $i$ -separation has been shown:

**Lemma 1.1** (see [10, Lemma 5.1]) *Let  $\mathcal{A}$  and  $\mathcal{B}$  be any two language families and assume that  $\mathcal{A} - \mathcal{B}$  is REG-dissectible. It then holds that, for any  $A \in \mathcal{A}$  and any  $B \in \mathcal{B}$ , if  $A$   $i$ -covers  $B$ , then there exists a language in  $\mathcal{E}$  that  $i$ -separates  $i(B, A)$ , where  $\mathcal{E}$  expresses the set  $\{B \cup (C \cap A) \mid A \in \mathcal{A}, B \in \mathcal{B}, C \in \text{REG}\}$ . In other words,  $\mathcal{E}$   $i$ -separates  $i(\mathcal{B}, \mathcal{A})$ .*

Although Lemma 1.1 is stated explicitly for REG-dissectibility, from the proof in [10] it is clear that any family of languages could be applied. For convenience, in Section 5 we present such generalization with a proof; see Lemma 5.1. This generalized connection between dissectibility and  $i$ -separation adds another argument for the study of a dissection of infinite languages.

There is a longstanding open question in [1]: Given two context-free languages  $L_1, L_2$  such that  $L_1 \subset L_2$  and  $L_2 \setminus L_1$  is an infinite language, is there a context-free language  $L_3$  such that  $L_3 \subset L_2$ ,  $L_1 \subset L_3$ , and both the languages  $L_3 \setminus L_1$  and  $L_2 \setminus L_3$  are infinite? This question was mentioned also in [10] using the  $i$ -separation: Let CFL denote the family of all context free languages. Does CFL  $i$ -separate  $i(\text{CFL}, \text{CFL})$ ? Understanding the dissectibility could help to solve this open question or at least it could help to identify “minimal” language families  $\mathcal{C}$  such that  $\mathcal{C}$   $i$ -separates  $i(\text{CFL}, \text{CFL})$ .

Some other results concerning the dissection of infinite languages may be found in [5]. Two related topics are the construction of minimal covers of languages, [2], and the immunity of languages, [3, 7, 10]. Recall that a language  $L_1$  is called  $\mathcal{C}$ -immune if there is no infinite language  $L_2 \subseteq L_1$  such that  $L_2 \in \mathcal{C}$ .

Let  $\mathbb{N}^+$  denote the set of all positive integers. An infinite language  $L$  is called *constantly growing* if there is a constant  $c$  such that for every word  $u \in L$  there is a word  $v \in L$  with  $|u| < |v| \leq c + |u|$ . In [10], it has been proved that every constantly growing language  $L$  is REG-dissectible.

We introduce a “natural” generalization of constantly growing languages as follows. Let  $\mathbb{R}^+$  denote the set of all positive real numbers. We define that a language  $L$  is a *geometrically growing* language if there is a constant  $c \in \mathbb{R}^+$  such that for every  $u \in L$  there exists  $v \in L$  with  $|u| < |v| \leq c|u|$ . We say also that  $L$  is  $c$ -geometrically growing. In the current article we show how to dissect geometrically growing language by a homomorphic image of intersection of two context-free languages.

Consider two alphabets  $\Sigma$  and  $\Theta$  such that  $|\Sigma| = 1$  and  $|\Theta| = 4$ ; that is  $\Sigma$  and  $\Theta$  denote alphabets with one letter and four letters, respectively. The main results of the current article are the following theorem and its corollary below.

**Theorem 1.2** *There are context-free languages  $M_1, M_2 \subseteq \Theta^*$  and an erasing alphabetical homomorphism  $\pi : \Theta^* \rightarrow \Sigma^*$  such that: If  $L \subseteq \Sigma^*$  is a geometrically growing language then there is a regular language  $R \subseteq \Theta^*$  such that  $\pi(R \cap M_1 \cap M_2)$  dissects the language  $L$ .*

To emphasize the essential result of our article, we consider in Theorem 1.2 that  $L$  is a language over an alphabet with one letter. Let  $\Gamma$  denote a finite alphabet. The next corollary shows a generalization for a geometrically growing language over the alphabet  $\Gamma$ .

**Corollary 1.3** *There are context-free languages  $M_1, M_2 \subseteq \Theta^*$ , an erasing alphabetical homomorphism  $\pi : \Theta^* \rightarrow \Sigma^*$ , and a nonerasing alphabetical homomorphism  $\varphi : \Gamma^* \rightarrow \Sigma^*$  such that: If  $L \subseteq \Gamma^*$  is a geometrically growing language then there is a regular language  $R \subseteq \Theta^*$  such that  $\varphi^{-1}(\pi(R \cap M_1 \cap M_2))$  dissects the language  $L$ .*

**Proof:** Let  $L_1 \subseteq \Gamma^*$  be a language and let  $\varphi : \Gamma^* \rightarrow \Sigma^*$  be an alphabetical homomorphism defined as follows:  $\varphi(a) = z$  for every  $a \in \Gamma$ , where  $z$  is the only letter of the alphabet  $\Sigma$ . If  $L_2 \subseteq \Sigma^*$  and  $L_2$  dissects  $\varphi(L_1)$  then clearly the language  $L_3 = \varphi^{-1}(L_2) = \{w \in L_1 \mid \varphi(w) \in L_2\}$  dissects  $L_1$ . This completes the proof.  $\square$

**Remark 1.4** *Since the intersection of a regular language and a context-free language is a context-free language we have  $R \cap M_1 \cap M_2$  is also intersection of two context-free languages. This explains why we do not mention the regular language in the title of the article.*

We sketch the basic ideas of our proof. Note that a non-associative word on the letter  $a$  is a “well parenthesized” word containing a given number of occurrences of  $a$ . It is known that the number of non-associative words containing  $n + 1$  occurrences of  $a$  is equal to the  $n$ -th Catalan number [8]. For example for  $n = 3$  we have five distinct non-associative words:  $((aa)a)a$ ,  $((aa)(aa))$ ,  $(a(a(aa)))$ ,  $(a((aa)a))$ , and  $((a(aa))a)$ . Every non-associative word contains the prefix  $(^k a$  for some  $k \in \mathbb{N}^+$ , where  $(^k$  denotes the  $k$ -th power of the opening bracket. We show that there are non-associative words such that  $k$  equals “approximately”  $\log_2 n$ . We construct two context-free languages whose intersection accepts such words and we call these words *balanced extended non-associative words*. By counting the number of opening brackets of a balanced extended non-associative word with  $n$  occurrences of  $a$  we can compute the logarithm of the number of occurrences of  $a$ . If  $L$  is a geometrically growing language then the language

$$\widehat{L} = \{a^j \mid j = \lceil \log_2 |w| \rceil \text{ and } w \in L\}$$

is obviously constantly growing. Hence, by means of intersection of two context-free languages we transform the challenge of dissecting a geometrically growing language to the challenge of dissecting a constantly growing language. This approach allows us to prove our result.

## 2 Preliminaries

Let  $\epsilon$  denote the empty word. Given a finite alphabet  $A$ , let  $A^+$  denote the set of all finite nonempty words over the alphabet  $A$  and let  $A^* = A^+ \cup \{\epsilon\}$ .

Let  $\text{Fac}(w)$  denote the set of all factors of the word  $w \in A^*$ . We have  $\epsilon, w \in \text{Fac}(w)$ .

Let  $\text{Pref}(w), \text{Suf}(w) \subseteq \text{Fac}(w)$  denote the set of all prefixes and suffixes of  $w \in A^*$ , respectively. We have  $\epsilon, w \in \text{Pref}(w) \cap \text{Suf}(w)$ . Let  $\text{occur}(w, t)$  denote the number of occurrences of the factor  $t \in A^+$  in the word  $w \in A^*$ ; formally

$$\text{occur}(w, t) = |\{v \in \text{Suf}(w) \mid t \in \text{Pref}(v)\}|.$$

Given two finite alphabets  $A_1, A_2$ , a *homomorphism* from  $A_1^*$  to  $A_2^*$  is a function  $\tau : A_1^* \rightarrow A_2^*$  such  $\tau(uv) = \tau(u)\tau(v)$ , where  $u, v \in A_1^*$ . It follows that in order to define a homomorphism  $\tau$ , it suffices to define  $\tau(a)$  for every  $a \in A_1$ ; such definition “naturally” extends to every word  $u \in A_1^*$ . We say that  $\tau$  is an *alphabetical homomorphism* if  $\tau(a) \in A_2$  for every  $a \in A_1$ . We say that  $\tau$  is an *erasing alphabetical homomorphism* if  $\tau(a) \in A_2 \cup \{\epsilon\}$  for every  $a \in A_1$  and there is at least one  $a \in A_1$  such that  $\tau(a) = \epsilon$ .

### 3 Balanced non-associative words

Let  $\Theta = \{x, y, z, p\}$ . We reserve the symbols  $x, y, z, p$  for the letters of the alphabet  $\Theta$ . It means that wherever in our article we use the symbols  $x, y, z, p$ , we refer to the letters of  $\Theta$ .

Let  $\text{ENW} \subseteq \Theta^*$  be the language generated by the following context-free grammar, where  $S$  is a start non-terminal symbol,  $P$  is a non-terminal symbol, and  $x, y, z, p \in \Theta$  are terminal symbols:

- $S \rightarrow xPPy$ ,
- $P \rightarrow S \mid pzp \mid pzzp$ .

We call the words from  $\text{ENW}$  *extended non-associative words*.

**Remark 3.1** *Let the letter  $x$  represent an opening bracket and the letter  $y$  a closing bracket. It is easy to see that if  $v_1, v_2 \in \{pzp, pzzp\} \cup \text{ENW}$  then  $xv_1v_2y \in \text{ENW}$ . Note that if  $w \in \text{ENW}$  then  $w$  is “well parenthesized” with brackets  $x$  and  $y$ . Also note that if  $w \in \text{ENW}$ ,  $xvy \in \text{Fac}(w)$ ,  $\text{occur}(v, x) = 0$ , and  $\text{occur}(v, y) = 0$ , then  $v \in \{pzppzzp, pzppzp, pzzppzpzp, pzzppzpzp\}$ .*

**Remark 3.2** *Recall from [8] that a “standard” non-associative word on the letter  $a$ , mentioned in the introduction, can be represented as a full binary rooted tree, where every inner node represents a corresponding pair of brackets and every leaf represents the letter  $a$ . It is known that the number of inner nodes plus one is equal to the number of leaves in a full binary rooted tree.*

*Obviously we can also represent the extended non-associative words from  $\text{ENW}$  as full binary rooted trees, where the factors  $pzp$  and  $pzzp$  represent the leaves. It follows that if  $w \in \text{ENW}$  then*

$$\text{occur}(w, x) + 1 = \text{occur}(w, pzp) + \text{occur}(w, pzzp).$$

*If  $w$  is a non-associative word on the symbol  $a$  with brackets  $x, y$  having  $n + 1$  occurrences of  $a$  then we get  $2^{n+1}$  extended non-associative words by replacing  $a$  with  $pzp$  or  $pzzp$ ; for example if  $K = \{xxa_1a_2ya_3y \mid a_1, a_2, a_3 \in \{pzp, pzzp\}\}$  then  $|K| = 2^3 = 8$  and  $K \subseteq \text{ENW}$ . Since the number of non-associative words containing  $n + 1$  occurrences of  $a$  is equal to the  $n$ -th Catalan number  $C_n$  [8], it is clear that*

$$|\{w \in \text{ENW} \mid \text{occur}(w, pzp) + \text{occur}(w, pzzp) = n + 1\}| = 2^{n+1}C_n,$$

where  $n \in \mathbb{N}^+$ .

Let  $BAL \subseteq \Theta^*$  be the language generated by the following context-free grammar, where  $S$  is a start non-terminal symbol,  $T, V, Z$  are non-terminal symbols, and  $x, y, z, p \in \Theta$  are terminal symbols:

- $S \rightarrow xSy \mid pZVZp$ ,
- $V \rightarrow VZV \mid pTp$ ,
- $T \rightarrow yTx \mid \epsilon$ ,
- $Z \rightarrow z \mid zz$ .

We call the words from  $BAL$  *balanced words*. The reason for the name “balanced” comes from the following lemma.

**Lemma 3.3** *If  $u \in BAL$ ,  $w \in \text{Fac}(u)$ , and  $pwp \in \text{Fac}(u)$  then  $\text{occur}(w, x) = \text{occur}(w, y)$ .*

**Proof:** The proof is by induction on  $j = \text{occur}(w, p)$ . From the definition of the language  $BAL$ , it is clear that if  $\text{occur}(w, p) = 0$  then  $w = y^i x^i$  for some  $i \in \{0\} \cup \mathbb{N}^+$ . Hence we have the base case for  $j = 0$ . Suppose  $j > 0$ . Then it follows that  $w = w_1 p w_2$  for some  $w_1, w_2 \in \text{Fac}(w)$ . Since  $\text{occur}(w_1, p) + 1 + \text{occur}(w_2, p) = j$ , we have  $\text{occur}(w_1, p), \text{occur}(w_2, p) < j$ . Hence the lemma holds for both  $p w_1 p$  and  $p w_2 p$  and in consequence lemma holds also for  $pwp$ . This completes the proof.  $\square$

Let  $\Omega = \text{ENW} \cap BAL \subseteq \Theta^*$ . We call the words from  $\Omega$  *balanced extended non-associative words*. Let  $\Omega(n) = \{w \in \Omega \mid \text{occur}(w, z) = n\}$ , where  $n \in \mathbb{N}^+$ .

**Remark 3.4** *To understand the idea of balanced extended non-associative words, suppose  $w \in \Omega$  and let  $G$  be the full binary rooted tree that represents  $w$  (as explained in Remark 3.2). Then in  $G$ , the length of the path from the root to a leaf does not depend on the leaf; it means the number of inner nodes lying on the path from a leaf to the root is a constant for  $G$ .*

**Example 3.5** *Let  $v_1 = xpzpxzppzzyy$  and  $v_2 = xpzppzpyxpzppzzyy$ . We have  $v_1, v_2 \in \text{ENW}$ ,  $v_1 \notin \Omega$ , and  $v_2 \in \Omega$ .*

Given a word  $w \in \Theta^*$ , let  $\text{height}(w) = \max\{j \mid x^j \in \text{Fac}(w)\}$ . We call  $\text{height}(w)$  the *height* of  $w$ . We show that if  $w \in \Omega$  and  $h$  is the height of  $w$  then  $x^h$  is a prefix of  $w$  and  $y^h$  is a suffix of  $w$ .

**Lemma 3.6** *If  $w \in \Omega$  and  $h = \text{height}(w)$  then  $x^h \in \text{Pref}(w)$  and  $y^h \in \text{Suf}(w)$ .*

**Proof:** Since  $\Omega \subseteq \text{ENW}$ , there is  $\hat{h} \in \mathbb{N}^+$  such that  $x^{\hat{h}} p \in \text{Pref}(w)$ . To get a contradiction suppose that  $\hat{h} < h$ . Because  $\Omega \subseteq BAL$  it follows that  $w = x^{\hat{h}} p w_1 p y^{\hat{h}} x^{\hat{h}} p w_2$  for some  $w_1 \in \text{Fac}(w)$ ,  $w_2 \in \text{Suf}(w)$ , and

$$\text{occur}(x^{\hat{h}} p w_1 p y^{\hat{h}} x^{\hat{h}}, x^{\hat{h}}) = 1.$$

Lemma 3.3 implies that  $\text{occur}(p w_1 p, x) = \text{occur}(p w_1 p, y)$ . Let  $r = x^{\hat{h}} p w_1 p y^{\hat{h}}$ . It follows that

$$\text{occur}(r, x) < \text{occur}(r, y).$$

This is a contradiction, since for every prefix  $v \in \text{Pref}(w)$  of an extended non-associative word  $w \in \text{ENW}$  (a well parenthesized word) we have  $\text{occur}(v, x) \geq \text{occur}(v, y)$ . We conclude that  $\hat{h} = h$  and  $x^h \in \text{Pref}(w)$ . In an analogous way we can show that  $y^h \in \text{Suf}(w)$ . This completes the proof.  $\square$

For a word  $w \in \Omega$ , we show the relation between the height of  $w$  and the number of occurrences of  $z$  in  $w$ .

**Proposition 3.7** *If  $w \in \Omega$  and  $h = \text{height}(w)$  then*

$$2^h \leq \text{occur}(w, z) \leq 2^{h+1}.$$

**Proof:** From the definition of ENW it follows that  $h \geq 1$ . We prove the proposition by induction. Obviously if  $h = 1$  then

$$w \in \{xpzppzpy, xpzzppzpy, xpzppzzpy, xpzzppzzpy\}.$$

Thus the proposition holds for  $h = 1$ . Suppose that the proposition holds for all  $\hat{h} < h$  and let  $h \geq 2$ . Since  $\Omega \subseteq \text{ENW}$ , it follows that  $h \geq 2$  implies that  $w = xw_1w_2y$  for some  $w_1, w_2 \in \text{ENW} \cup \{pzp, pzzp\}$  with  $\{w_1, w_2\} \cap \text{ENW} \neq \emptyset$ . Without loss of generality suppose that  $w_1 \in \text{ENW}$ .

Let  $h_1 = \text{height}(w_1)$ . Lemma 3.6 implies that  $x^{h_1} \in \text{Pref}(w_1)$  and  $y^{h_1} \in \text{Suf}(w_1)$ . Lemma 3.3 implies that  $x^{h_1} \in \text{Pref}(w_2)$  and in consequence  $w_2 \in \Omega$ ,  $y^{h_1} \in \text{Suf}(w_2)$ , and  $\text{height}(w_2) = h_1$ .

Because  $x^{h_1} \in \text{Pref}(w_1)$  it follows that  $x^{h_1+1} \in \text{Pref}(w)$ . Thus  $h_1 + 1 = h$ . As we assumed that the proposition holds for all  $\hat{h} < h$ , we can derive that

$$\text{occur}(w, z) = \text{occur}(w_1, z) + \text{occur}(w_2, z) \leq 2^{h_1+1} + 2^{h_1+1} = 2^{h_1+2} = 2^{h+1}$$

and

$$\text{occur}(w, z) = \text{occur}(w_1, z) + \text{occur}(w_2, z) \geq 2^{h_1} + 2^{h_1} = 2^{h_1+1} = 2^h.$$

This completes the proof. □

**Remark 3.8** *Proposition 3.7 could be also proven using tree graphs as follows: There are exactly  $2^h$  leaves in a complete tree of height  $h$ , since there is a bijection (as mentioned in Remark 3.2), there are  $2^h$  occurrences of  $pzp$  and  $pzzp$ , and hence the number of occurrences of  $z$  in  $w$  is between  $2^h$  and  $2^{h+1}$ .*

Proposition 3.7 has the following obvious corollary.

**Corollary 3.9** *If  $n \in \mathbb{N}^+$ ,  $w \in \Omega(n)$ , and  $h = \text{height}(w)$  then*

$$\log_2 n - 1 \leq h \leq \log_2 n.$$

**Remark 3.10** *Note that the number of occurrence of  $z$  does not uniquely determine the height. For example if  $w_1 = xpzppzpyxpzppzpyy$  and  $w_2 = xpzzppzzpy$ , then  $w_1, w_2 \in \Omega(4)$  and  $1 = \text{height}(w_2) < \text{height}(w_1) = 2$ .*

Given  $w, u, v \in \Theta^+$ , let  $\text{replace}(w, v, u)$  denote the word built from  $w$  by replacing the first occurrence of  $v$  in  $w$  by  $u$ . Formally, if  $\text{occur}(w, v) = 0$  then  $\text{replace}(w, v, u) = w$ . If  $\text{occur}(w, v) = j > 0$  and  $w = w_1vw_2$ , where  $\text{occur}(vw_2, v) = j$  then  $\text{replace}(w, v, u) = w_1uw_2$ .

We prove that the set of balanced extended non-associative words  $\Omega(n)$  having  $n$  occurrences of  $z$  is nonempty for each  $n \geq 2$ .

**Proposition 3.11** *If  $n \in \mathbb{N}^+$  and  $n \geq 2$  then  $\Omega(n) \neq \emptyset$ .*

**Proof:** Let  $j \in \mathbb{N}^+$  be such that  $2^{j-1} < n \leq 2^j$ . Obviously such  $j$  exists and is uniquely determined. Let  $w_1 = xpzzppzppy$ . Let  $w_{i+1} = xw_iw_iy$  for every  $i \in \mathbb{N}^+$ . Clearly  $\text{occur}(w_j, z) = 2^{j+1}$  and  $w_j \in \Omega(2^j)$ . Note that  $\text{occur}(w_j, pzzp) = 2^j$ .

Let  $w_{j,0} = w_j$  and  $w_{j,i+1} = \text{replace}(w_{j,i}, pzzp, pzp)$ , where  $i \in \mathbb{N}^+ \cup \{0\}$  and  $i < 2^j$ . Let  $\alpha = 2^j - n$ . Then one can easily verify that  $\text{occur}(w_{j,\alpha}, z) = n$  and  $w_{j,\alpha} \in \Omega(n)$ .

In principle, we construct a balanced extended non-associative word  $w_j$  having  $2^j$  occurrences of  $pzzp$  and then we replace a certain number of occurrences of  $pzzp$  with the factor  $pszp$  to achieve the required number of occurrences of  $z$ . This completes the proof.  $\square$

## 4 Dissection of infinite languages

In [10] it was shown that every constantly growing language can be dissected by some regular language.

**Lemma 4.1** (see [10, Lemma 3.3]) *Every infinite constantly growing language is REG-dissectible.*

In the next proposition we show under which condition we can dissect a language  $L \subseteq \Omega$  by a regular language. Informally, the proposition says that a geometrically growing subset of balanced extended non-associative words is *REG*-dissectible.

**Proposition 4.2** *If  $\beta \in \mathbb{N}^+$ ,  $\beta \geq 2$ , and  $L \subseteq \Omega$  is an infinite language such that for each  $w_1 \in L$  there is  $w_2 \in L$  with*

$$\text{occur}(w_1, z) < \text{occur}(w_2, z) \leq \beta \text{occur}(w_1, z)$$

*then there is a regular language  $R \subseteq \Theta^*$  such that  $R$  dissects  $L$ .*

**Proof:** Let  $\alpha \in \mathbb{N}^+$  be such that  $2^{\alpha-1} < \beta \leq 2^\alpha$ . Obviously such  $\alpha$  exists and is uniquely determined. Given  $w_1 \in L$ , let  $w_2 \in L$  be such that

$$n_1 < n_2 \leq \beta n_1 \leq 2^\alpha n_1, \tag{1}$$

where  $n_1 = \text{occur}(w_1, z)$  and  $n_2 = \text{occur}(w_2, z)$ . From the conditions of the proposition such  $w_2$  exists.

Without loss of generality suppose that  $\log_2 n_1 \geq 2$ . Note that there are only finitely many words  $v \in \Omega$  with  $\log_2(\text{occur}(v, z)) < 2$ .

Let  $h_1 = \text{height}(w_1)$  and  $h_2 = \text{height}(w_2)$ . Corollary 3.9 implies that

$$\log_2 n_1 - 1 \leq h_1 \text{ and } h_2 \leq \log_2 n_2 \tag{2}$$

From (1) and (2) it follows that

$$h_2 \leq \log_2 n_2 \leq \log_2 (2^\alpha n_1) = \alpha + \log_2 n_1 \leq \alpha + 1 + h_1. \tag{3}$$

Since we selected  $w_1$  arbitrarily, it follows from (3) that

$$H = \{x^h \mid h = \text{height}(v) \text{ and } v \in L\}$$

is a constantly growing language.

Lemma 4.1 implies that  $H \subseteq \{x\}^*$  is *REG*-dissectible. Let  $\widehat{R} \subseteq \{x\}^*$  be a regular language that dissects  $H$ . Let  $R = \{rvv \mid r \in \widehat{R} \text{ and } v \in \Theta^*\} \subseteq \Theta^*$ . Obviously  $R$  is a regular language that dissects  $L$ ; to see this, recall that if  $w \in \Omega$ ,  $i \in \mathbb{N}^+$ ,  $a \in \Theta \setminus \{x\}$ , and  $x^i a \in \text{Pref}(w)$  then  $a = p$ .

This completes the proof.  $\square$

We step to the proof of the main theorem of the current article.

**Proof of Theorem 1.2:** Without loss of generality let  $\Sigma = \{z\}$  be the alphabet with the letter  $z \in \Theta$ . Let  $\pi : \Theta^* \rightarrow \Sigma^*$  be an erasing alphabetical homomorphism defined as follows:

$$\pi(a) = \begin{cases} z & \text{If } a = z. \\ \epsilon & \text{If } a \in \{x, y, p\}. \end{cases}$$

Thus  $\pi$  erases all letters except for the letter  $z$ . From the definition of ENW, BAL, and  $\Omega$ , it follows that the language  $\Omega$  is an intersection of two context-free languages ENW and BAL. Let  $M_1 = \text{ENW}$  and let  $M_2 = \text{BAL}$ ; recall that  $M_1$  and  $M_2$  are used in the statement of Theorem 1.2.

Let  $\widehat{L} = \{w \in \Omega \mid \pi(w) \in L\}$ . Note that  $\widehat{L}$  contains  $w \in \Omega$  if and only if there is a word  $v \in L$  such that the number of occurrences of  $z$  in  $w$  is equal to the length of  $v$ ; formally  $\text{occ}(w, z) = |v|$ . Proposition 3.11 implies that  $\widehat{L}$  is an infinite language.

Let  $c \in \mathbb{R}^+$  be such that for every  $u \in L$  there exists  $v \in L$  with  $|u| < |v| \leq c|u|$ . Since  $L$  is a geometrically growing language, we know that such  $c$  exists. Let  $\beta \in \mathbb{N}^+$  be such that  $\beta \geq 2$  and  $\beta \geq c$ . Hence  $L$  is  $\beta$ -geometrically growing language. It follows that if  $w_1 \in \widehat{L}$  then there is a word  $w_2 \in \widehat{L}$  with

$$\text{occ}(w_1, z) < \text{occ}(w_2, z) \leq \beta \text{occ}(w_1, z).$$

Then Proposition 4.2 implies that there is a regular language  $R \subseteq \Theta^*$  that dissects  $\widehat{L}$ . This implies that the homomorphic image  $\pi(R \cap \text{ENW} \cap \text{BAL})$  dissects the language  $L$ . This completes the proof.  $\square$

## 5 Dissection and i-separation

As mentioned in the introduction, for convenience we present a generalization of Lemma 1.1 ([10, Lemma 5.1]), which demonstrates the connection between dissectibility and i-separation. The presented proof is just a copy of the proof in [10] by changing REG to  $\mathcal{C}$ .

**Lemma 5.1** *Let  $\mathcal{A}$ ,  $\mathcal{B}$ , and  $\mathcal{C}$  be any three language families and assume that  $\mathcal{A} - \mathcal{B}$  is  $\mathcal{C}$ -dissectible. It then holds that, for any  $A \in \mathcal{A}$  and any  $B \in \mathcal{B}$ , if  $A$  i-covers  $B$ , then there exists a language in  $\mathcal{E}$  that i-separates  $i(B, A)$ , where  $\mathcal{E}$  expresses the set  $\{B \cup (C \cap A) \mid A \in \mathcal{A}, B \in \mathcal{B}, C \in \mathcal{C}\}$ . In other words,  $\mathcal{E}$  i-separates  $i(\mathcal{B}, \mathcal{A})$ .*

**Proof:** Let  $A \in \mathcal{A}$  and  $B \in \mathcal{B}$  be two infinite languages. Let  $D = A - B$  and assume that  $D$  is infinite. Our assumption guarantees the existence of a language  $C \in \mathcal{C}$  for which  $C$  dissects  $D$ . We set  $E = B \cup (A \cap C)$ . Since  $C$  dissects  $D$ , it follows that  $|(A \cap C) - B| = \infty$  and  $|(A \cap \overline{C}) - B| = \infty$ . It follows that  $B \subseteq E \subseteq A$  and  $|A - E| = |E - B| = \infty$ . Thus,  $E$  i-separates  $i(B, A)$ . Since  $C \in \mathcal{C}$ ,  $E$  belongs to the language family  $\mathcal{E}$ . This completes the proof.  $\square$

## 6 Open questions

In the current article we applied a new idea of dissecting a language  $L$  by a homomorphic image of a language  $\widehat{L}$  from the family of languages  $\text{CFL}_2$ . The idea can be generalized for every  $\text{CFL}_k$ , where  $k \in \mathbb{N}^+$ . Let us introduce a notation for this technique. Given an alphabet  $A$  and a positive integer  $k$ , let

$$\begin{aligned} \text{hiCFL}_{k,A} = \{L \subseteq A^* \mid & \text{there are an alphabet } \widehat{A} \text{ and} \\ & \text{context-free languages } L_i \subseteq \widehat{A}^* \text{ with } i \in \{1, 2, \dots, k\} \\ & \text{and a homomorphism } \phi : \widehat{A}^* \rightarrow A^* \\ & \text{such that } L = \phi\left(\bigcap_{i=1}^k L_i\right)\}. \end{aligned}$$

The prefix ‘‘hi’’ stands for ‘‘homomorphic image’’. Using the set  $\text{hiCFL}_{k,A}$  we can restate Theorem 1.2 as follows:

**Corollary 6.1** *Every geometrically growing language over an alphabet  $A$  with  $|A| = 1$  is  $\text{hiCFL}_{2,A}$ -dissectible.*

Moreover we introduced in the current article the notion of a geometrically growing language. We generalize this concept as follows. Let

$$\Pi = \{\sigma : \mathbb{N}^+ \rightarrow \mathbb{N}^+ \mid \sigma(n) > n \text{ for all } n \in \mathbb{N}^+\}.$$

Given  $\sigma \in \Pi$ , we say that a language  $L$  is  $\sigma$ -growing if for every word  $u \in L$  there is a word  $v \in L$  such that  $|u| < |v| \leq \sigma(|u|)$ .

**Remark 6.2** *Let  $\tilde{\sigma}(n) = cn$  for some  $c \in \mathbb{N}^+$  with  $c > 1$ . Obviously  $\tilde{\sigma} \in \Pi$ . A language  $L$  is  $\tilde{\sigma}$ -growing language if and only if  $L$  is  $c$ -geometrically growing.*

Let  $\mathcal{C}$  be a family of languages. Using the notion of  $\sigma$ -growing languages, we present the following open questions and problems:

- Find  $\sigma \in \Pi$  such that there exists a  $\sigma$ -growing language  $L$  that is not  $\mathcal{C}$ -dissectible or show that such  $\sigma$  does not exist.
- Find  $\sigma \in \Pi$  such that:
  - Every  $\sigma$ -growing language  $L$  is  $\mathcal{C}$ -dissectible.
  - If  $\hat{\sigma} \in \Pi$  and  $\hat{\sigma}(n) > \sigma(n)$  for all  $n \in \mathbb{N}^+$  then there is a  $\hat{\sigma}$ -growing language  $L$  that is not  $\mathcal{C}$ -dissectible.

Concerning the family of languages  $\mathcal{C}$ , we are particularly interested in  $\text{REG}$ ,  $\text{CFL}_k$ , and  $\text{hiCFL}_{k,A}$  for all  $k \in \mathbb{N}^+$ . However the questions may be of interest also for other families.

We list some more open questions and problems in spite of the fact that some of them are already mentioned (directly or indirectly) above.

- Is the family of geometrically growing languages  $\text{REG}$ -dissectible?

- Does CFL  $i$ -separate  $i(\text{CFL}, \text{CFL})$  (mentioned in [1] and [10])?
- Describe the “minimal” families of languages  $\mathcal{C}$  such that  $\mathcal{C}$   $i$ -separates  $i(\text{CFL}, \text{CFL})$ .
- Let  $\text{CFLGL} \subseteq \text{CFL}$  be the family of geometrically growing context-free languages. Describe the “minimal” families of languages  $\mathcal{C}$  such that  $\mathcal{C}$   $i$ -separates  $i(\text{CFLGL}, \text{CFLGL})$ . In particular, are all geometrically growing context-free languages REG-dissectible?
- Is there  $\sigma \in \Pi$  such that if  $i(L_1, L_2) \in i(\text{CFL}, \text{CFL})$  then the language  $L_2 - L_1$  is  $\sigma$ -growing?
- Describe languages that are CFL-dissectible.
- Find an example of a language  $L$  such that  $L$  is CFL-dissectible and not REG-dissectible.

For more open questions about dissectibility, we recommend the readers to review the section “6. Future challenges” in [10].

## Acknowledgements

This work was supported by the Grant Agency of the Czech Technical University in Prague, grant No. SGS20/183/OHK4/3T/14.

## References

- [1] W. BUCHER, *A density problem for context-free languages*, Bull. Eur. Assoc. Theor. Comput. Sci. EATCS 10, (1980).
- [2] M. DOMARATZKI, J. SHALLIT, AND S. YU, *Minimal covers of formal languages*, in *Developments in Language Theory*, 2001.
- [3] P. FLAJOLET AND J. M. STEYAERT, *On sets having only hard subsets*, in *Automata, Languages and Programming*, J. Loeckx, ed., Berlin, Heidelberg, 1974, Springer Berlin Heidelberg, pp. 446–457.
- [4] S. GINSBURG AND S. GREIBACH, *Deterministic context free languages*, *Information and Control*, 9 (1966), pp. 620 – 648.
- [5] J. JULIE, J. BASKAR BABUJEE, AND V. MASILAMANI, *Dissecting power of certain matrix languages*, in *Theoretical Computer Science and Discrete Mathematics*, S. Arumugam, J. Bagga, L. W. Beineke, and B. Panda, eds., Cham, 2017, Springer International Publishing, pp. 98–105.
- [6] L. LIU AND P. WEINER, *An infinite hierarchy of intersections of context-free languages*, *Math. Systems Theory* 7, 185–192., (1973).
- [7] E. L. POST, *Recursively enumerable sets of positive integers and their decision problems*, *Bull. Amer. Math. Soc.*, 50 (1944), pp. 284–316.
- [8] R. P. STANLEY AND S. FOMIN, *Enumerative Combinatorics*, vol. 2 of *Cambridge Studies in Advanced Mathematics*, Cambridge University Press, 1999.

- [9] T. YAMAKAMI, *Intersection and union hierarchies of deterministic context-free languages and pumping lemmas*, in Language and Automata Theory and Applications, A. Leporati, C. Martín-Vide, D. Shapira, and C. Zandron, eds., Cham, 2020, Springer International Publishing, pp. 341–353.
- [10] T. YAMAKAMI AND Y. KATO, *The dissecting power of regular languages*, Information Processing Letters, 113 (2013), pp. 116 – 122.