

# Characterizations of families of morphisms and words via binomial complexities\*

Michel Rigo<sup>†</sup>, Manon Stipulanti<sup>‡</sup> and Markus A. Whiteland<sup>§¶</sup>

Department of Mathematics, University of Liège,  
Allée de la Découverte 12, 4000 Liège, Belgium

{m.rigo,m.stipulanti,mwhiteland}@uliege.be

## Abstract

Two words are  $k$ -binomially equivalent if each subword of length at most  $k$  occurs the same number of times in both words. The  $k$ -binomial complexity of an infinite word is a counting function that maps  $n$  to the number of  $k$ -binomial equivalence classes represented by its factors of length  $n$ . Cassaigne et al. [Int. J. Found. Comput. S., 22(4) (2011)] characterized a family of morphisms, which we call Parikh-collinear, as those morphisms that map all words to words with bounded 1-binomial complexity. Firstly, we extend this characterization: they map words with bounded  $k$ -binomial complexity to words with bounded  $(k + 1)$ -binomial complexity. As a consequence, fixed points of Parikh-collinear morphisms are shown to have bounded  $k$ -binomial complexity for all  $k$ . Secondly, we give a new characterization of Sturmian words with respect to their  $k$ -binomial complexity. Then we characterize recurrent words having, for some  $k$ , the same  $j$ -binomial complexity as the Thue–Morse word for all  $j \leq k$ . Finally, inspired by questions raised by Lejeune, we study the relationships between the  $k$ - and  $(k + 1)$ -binomial complexities of infinite words; as well as the link with the usual factor complexity.

**Keywords:** Factor complexity, Abelian complexity, Binomial complexity, powers of the Thue–Morse morphism, Sturmian words.

**2020 Mathematics Subject Classification:** Primary: 68R15. Secondary: 05A05.

## 1 Introduction

The combinatorial structure of an infinite word  $\mathbf{x} \in A^{\mathbb{N}}$  over a finite alphabet  $A$  may reveal important aspects of  $\mathbf{x}$  itself. This structure is often studied through its language  $\mathcal{L}(\mathbf{x})$ , i.e., the set of its factors, and in particular to inspect the set  $\mathcal{L}_n(\mathbf{x}) := \mathcal{L}(\mathbf{x}) \cap A^n$  of factors of length  $n$ . Even plain counting the cardinality of this set turns out to be a useful concept: with  $p_{\mathbf{x}}$  the *factor complexity* function defined as  $p_{\mathbf{x}}: \mathbb{N} \rightarrow \mathbb{N}$ ,  $n \mapsto \#\mathcal{L}_n(\mathbf{x})$ , the celebrated Morse–Hedlund theorem asserts that an infinite word  $\mathbf{x}$  is aperiodic if and only if  $p_{\mathbf{x}}(n) \geq n + 1$  for all  $n \geq 1$ . For instance, the Thue–Morse word  $\mathbf{t} = 01101001 \dots$  (also known as the Prouhet–Thue–Morse word), the fixed point of the morphism  $\varphi: 0 \mapsto 01, 1 \mapsto 10$ , is aperiodic because its factor complexity is given by

$$p_{\mathbf{t}}(2^m + r) = \begin{cases} 3 \cdot 2^m + 4(r - 1), & \text{if } 1 \leq r \leq 2^{m-1}; \\ 4 \cdot 2^m + 2(r - 1), & \text{if } 2^{m-1} < r \leq 2^m, \end{cases} \quad (1)$$

---

\*Markus Whiteland dedicates this paper to the memory of his father Alan Whiteland (1940–2021).

<sup>†</sup>Supported by the FNRS Research grant T.196.23 (PDR)

<sup>‡</sup>Supported by the FNRS Research grant 1.C.104.24F.

<sup>§</sup>Supported by the FNRS Research grant 1.B.466.21F.

<sup>¶</sup>Corresponding author

for all  $m \geq 0$  [9, §4]. The factor complexity has proved its importance in a number of areas of mathematics. For example, in number theory, Adamczewski and Bugeaud [2] proved that the base- $b$  expansion of a real algebraic irrational number has a factor complexity function satisfying

$$\liminf_{n \rightarrow \infty} \frac{p(n)}{n} = \infty.$$

As a consequence, the number having  $t$  as base-2 expansion is transcendental.

One can conversely define families of words using the factor complexity function. For example, a word  $x$  is called *Sturmian* if  $p_x(n) = n + 1$  for all  $n$ . Such words, studied also in this note, turn out to have many interesting properties. For general references about combinatorics on (Sturmian) words, see, for instance, [5, 9, 32].

Many variations of the factor complexity have been introduced. Some of these counting functions only take into account factors with specific properties such as palindromes or privileged words [16, 37]. Other functions count subwords extracted along subsequences of prescribed forms like maximal pattern or arithmetical complexities [25, 7]. Closely related to the subject discussed in this paper, abelian,  $k$ -abelian or cyclic complexities are functions of the form  $n \mapsto \#(\mathcal{L}_n(x)/\sim)$  for the quotient by a relevant equivalence relation  $\sim$  [42, 26, 10]. For more on abelian combinatorics on words (and related notions), we refer the reader to the recent excellent survey [18]. For each of the above complexity functions usual questions naturally arise:

- What is the complexity of well-known families of words such as Sturmian, Arnoux-Rauzy, automatic, (pure) morphic, or Toeplitz words?
- This leads to the more interesting problem of classifying or characterizing infinite words with respect to their complexity. As an example, Coven and Hedlund in [15] show that an infinite word is purely periodic if and only if its abelian complexity function attains the value 1. Further, a binary aperiodic word has its abelian complexity function equal to the constant function 2 if and only if it is Sturmian. Words with linear factor complexity are characterized in [11].
- What are the possible growth rates of the complexity function?
- Which non-periodic words may achieved the lowest complexity?

These questions have intrinsic theoretical interests but also provide particular insight about the combinatorial structure of the studied words. Depending on the properties of interest, one focuses on the appropriate complexity function. Infinite words with specific combinatorial properties are, for instance, sought to construct particular symbolic dynamical systems or tilings of the line. For instance, the Thue–Morse minimal subshift is completely characterized by its factor complexity together with its abelian complexity [40]. The Thue–Morse word, which is central in our paper, plays an important role in many areas of mathematics, e.g., see [4, 3].

In this paper, the complexity function of interest is built from binomial coefficients of words.

**Definition 1.1.** Let  $u, w \in A^*$ . The *binomial coefficient* of  $u$  and  $w$  is the number of times  $w$  occurs as a subword of  $u$ , i.e. writing  $u = u_1 \cdots u_n$  with  $u_i \in A$ ,

$$\binom{u}{w} = \# \{i_1 < i_2 < \cdots < i_{|w|} : u_{i_1} u_{i_2} \cdots u_{i_{|w|}} = w\}.$$

These binomial coefficients have proven to be useful in a variety of domains: generalizations of Pascal-like triangles [31], algebra and topology [38, 8], formal languages or relationship with the extensively studied Parikh matrices and Simon’s congruence [6, 20, 22]. For more on these binomial coefficients, see, for instance, [32, §6].

We mention a well-known and actively researched problem related to binomial coefficients. A word  $u$  is  *$k$ -reconstructible* whenever the knowledge of the binomial coefficients  $\binom{u}{v}$ , for all subwords  $v$  of length  $k$ , uniquely determines  $u$ . Inspired by the problem of reconstructing graphs

from vertex-deleted subgraphs, the famous reconstruction problem is to determine the function  $f(n) = k$  where  $k$  is the least integer for which all words of length  $n$  (over a given alphabet) are  $k$ -reconstructible. For articles on this problem, we mention [19, 24, 34] and references therein.

Let us now introduce our main object of study. Let  $k \geq 1$  be an integer. The  $k$ -binomial complexity function introduced in [43] is the central theme of Lejeune's thesis [28]. It is built on the  $k$ -binomial equivalence where factors are distinguished with respect to the number of occurring subwords.

**Definition 1.2.** Two words  $u, v \in A^*$  are  $k$ -binomially equivalent, and we write  $u \sim_k v$ , if

$$\binom{u}{x} = \binom{v}{x}, \quad \forall x \in A^{\leq k}.$$

As an example the words  $u = 0110$  and  $v = 1001$  are 2-binomially equivalent because for  $z \in \{u, v\}$

$$\binom{z}{0} = 2, \binom{z}{1} = 2, \binom{z}{00} = 1, \binom{z}{01} = 2, \binom{z}{10} = 2, \binom{z}{11} = 1$$

and  $u \not\sim_3 v$  because 011 is a subword of  $u$  and not of  $v$ .

In [34, Lem. 1], it is observed that one may replace the condition  $\forall x \in A^{\leq k}$  with  $\forall x \in A^k$  as soon as  $|u|, |v| \geq k$ . Observe that the word  $u$  is obtained as a permutation of the letters in  $v$  if and only if  $u \sim_1 v$ . The latter relation is the *abelian equivalence* already introduced by Erdős [17]. This leads to the following definition, introduced in [43].

**Definition 1.3.** Let  $k \geq 1$  be an integer. The  $k$ -binomial complexity function of an infinite word  $x$  is defined as  $b_x^{(k)}: \mathbb{N} \rightarrow \mathbb{N}, n \mapsto \#(\mathcal{L}_n(x)/\sim_k)$ .

As an example, the first few values for the Thue–Morse word  $t$  are given in Table 1.

	0	1	2	3	4	5	6	7	8	9	10
$b_t^{(1)}$	1	2	3	2	3	2	3	2	3	2	3
$b_t^{(2)}$	1	2	4	6	9	8	8	8	9	8	8
$p_t$	1	2	4	6	10	12	16	20	22	24	28

Table 1: The first few values of  $b_t^{(1)}$ ,  $b_t^{(2)}$  and  $p_t$ .

It is clear that we have a series of refinements of the abelian equivalence: for all  $k \geq 1$ ,  $u \sim_{k+1} v$  implies  $u \sim_k v$ . Thus, for all  $n$ , we have the inequalities

$$b_x^{(1)}(n) \leq b_x^{(2)}(n) \leq \dots \leq b_x^{(k)}(n) \leq b_x^{(k+1)}(n) \leq \dots \leq p_x(n). \quad (2)$$

The study of the  $k$ -binomial complexity function has so far been studied for restricted families of words. For example, for  $k \geq 2$ , the  $k$ -binomial complexity of Sturmian words coincides with their factor complexity [43] (recalled here as Theorem 2.8) and the same property holds for the Tribonacci word [30]. For any  $k \geq 2$ , fixed points of Parikh-constant morphisms (see the next part for a definition) are known to have bounded  $k$ -binomial complexity [43]. Recently, the  $k$ -binomial complexities of the Thue–Morse word [29] (given in (3)) and the 2-binomial complexities of generalized Thue–Morse words was also computed [33]. That is the extent to which the notion has been studied.

We remark that a better understanding of the  $k$ -binomial complexity may give information about the language  $\mathcal{L}(x)$  of an infinite word  $x$  for which the reconstruction problem could be solved. The aim is to restrict the reconstruction problem to the language of an infinite word having a  $k$ -binomial complexity of the same order as its factor complexity. Indeed, if  $b_x^{(k)} = p_x$  for some  $k$ , then for any two distinct factors  $y, z$  of  $x$ , there exists a subword  $v$  of length  $k$  such that  $\binom{y}{v} \neq \binom{z}{v}$ .

Finally a parallel can be drawn between the  $k$ -abelian complexity introduced by Karhumäki et al. [26] and the  $k$ -binomial complexity. In both cases, we have a series of refinements (2) of the abelian equivalence. The fundamental difference is the following one. Two finite words  $u, v$  are *k-abelian equivalent* if, for each word  $w$  of length at most  $k$ , we count the same number of occurrences of the factor  $w$  in both words  $u$  and  $v$ . We thus make the important distinction between a *factor* and a *subword* of a word. Many properties of the  $k$ -abelian complexity have been recently and extensively studied such as growth and fluctuations,  $k$ -abelian palindromes, variation of Morse–Hedlund theorem, etc. [13, 12, 27]. This is to be contrasted with the limited knowledge we have on the  $k$ -binomial complexity function. Indeed, part of our motivation for this work stems from this rather limited state of the art as described above.

## 1.1 Our Results

We present three kinds of results: a new characterization of Parikh-collinear morphisms and links with bounded binomial complexities; a characterization of recurrent words with the same  $j$ -binomial complexities as the Thue–Morse word for  $j = 1, \dots, k$ ; study of the relationships existing between  $b_w^{(k)}$  and  $b_w^{(k+1)}$ . This paper improves upon the preliminary conference version [44]: not only do we provide proofs of results announced therein, but we also significantly extend them.

- Morphisms mapping all infinite words to words with bounded abelian complexity have been characterized in [14]. Such a morphism  $f: A^* \rightarrow B^*$  is said to be *Parikh-collinear*: for all letters  $a, b \in A$ , there is  $r_{a,b} \in \mathbb{Q}$  such that  $\Psi(f(b)) = r_{a,b} \Psi(f(a))$ , where  $\Psi(u)$  denotes the *Parikh vector* of a word  $u$  (see Section 2 for definitions). In Section 3, we obtain several new characterizations of Parikh-collinear morphisms. Connecting this with the series of inequalities (2), we show with Theorem 3.5 that a morphism is Parikh-collinear if and only if it maps all words with bounded  $k$ -binomial complexity to words with bounded  $(k+1)$ -binomial complexity.

It is known that any fixed point of a prolongable *Parikh-constant morphism*  $f: A^* \rightarrow A^*$ , i.e.,  $\Psi(f(a)) = \Psi(f(b))$  for all letters  $a, b \in A$ , has a bounded  $k$ -binomial complexity [43]. Any Parikh-constant morphism is obviously Parikh-collinear. As a direct consequence of our characterization of Parikh-collinear morphisms, Corollary 3.6 extends the previous result: bounded  $k$ -binomial complexity holds for any fixed point of a prolongable Parikh-collinear morphism.

- We then turn to words sharing their binomial complexities with the Thue–Morse word  $\mathbf{t}$ . From the above discussion (the Thue–Morse morphism  $\varphi$  is Parikh-constant), for all  $j \geq 1$ , the  $j$ -binomial complexity of  $\mathbf{t}$  is bounded by a constant depending on  $j$ . But more is known, the exact value of  $b_{\mathbf{t}}^{(j)}(n)$  computed in [29] is given by

$$b_{\mathbf{t}}^{(j)}(n) = \begin{cases} p_{\mathbf{t}}(n), & \text{if } n < 2^j; \\ 3 \cdot 2^j - 3, & \text{if } n \equiv 0 \pmod{2^j} \text{ and } n \geq 2^j; \\ 3 \cdot 2^j - 4, & \text{otherwise,} \end{cases} \quad (3)$$

where the factor complexity  $p_{\mathbf{t}}$  of  $\mathbf{t}$  is given by (1). Considering  $j = 1$  in (3), words having the same abelian complexity as the Thue–Morse words have been characterized in [40] as follows. The abelian complexity of an aperiodic word  $x \in \{0, 1\}^{\mathbb{N}}$  is, for  $n > 0$ ,  $b_x^{(1)}(n) = 3$  if  $n$  is even, and  $b_x^{(1)}(n) = 2$  if  $n$  is odd, if and only if there exists a word  $\mathbf{y}$  such that  $x = u\varphi(\mathbf{y})$  with  $u \in \{\varepsilon, 0, 1\}$ . Sections 4 and 5 are about binomial properties of iterates of  $\varphi$ . We generalize the latter result and obtain a characterization of words having the same  $j$ -binomial complexity as the Thue–Morse word  $\mathbf{t}$  for all  $j \leq k$ . Except for a remark in [18] (see Theorem 2.9), such a result together with Theorem 2.11 are the first where binomial complexity leads to the characterization of combinatorial families of words. In this paper, with Theorem 2.11, we observe that a word  $x$  is Sturmian if and only if  $b_x^{(2)}(n) = n + 1$ , for all  $n$ . We make the statements about words sharing the same  $j$ -binomial complexities as  $\mathbf{t}$  more precise.

Let  $k$  be an integer and let  $\mathbf{y}$  be an aperiodic binary word. With Theorem 4.2 we show that for  $x = u\varphi^k(\mathbf{y})$  we have, for all  $j \leq k$ ,  $b_x^{(j)} = b_{\mathbf{t}}^{(j)}$  which is given by (3), where  $u$  is a (possibly

empty) proper suffix of  $\varphi^k(0)$  or  $\varphi^k(1)$ . Conversely, with Theorem 5.2, if  $\mathbf{b}_x^{(j)} = \mathbf{b}_t^{(j)}$  for all  $j \leq k$  for a *recurrent* word  $\mathbf{x}$ , i.e., each factor of  $\mathbf{x}$  appears infinitely often, then  $\mathbf{x} = u\varphi^k(\mathbf{y})$  where  $u$  is a proper suffix of  $\varphi^k(0)$  or  $\varphi^k(1)$  and  $\mathbf{y}$  is some aperiodic binary word.

- In general, not much is known about the general behavior or fluctuations that can be expected for the  $k$ -binomial complexity of an infinite word. In particular, computing the  $k$ -binomial complexity of a particular infinite word remains quite challenging. It would also be desirable to compare in some ways the  $k$ - and  $(k+1)$ -binomial complexities of a word.

**Definition 1.4.** For two functions  $f, g: \mathbb{N} \rightarrow \mathbb{N}$ , we write  $f \prec g$  when the relation  $f(n) < g(n)$  holds for infinitely many  $n \in \mathbb{N}$ .

We define  $\prec$  this way because for some words, the 2-binomial complexity attains the factor complexity infinitely often while it is less than the factor complexity infinitely often. See end of Section 7.1 for a discussion.

As an example, a consequence of Proposition 4.17 is that  $\mathbf{b}_x^{(k)} \prec \mathbf{b}_x^{(k+1)}$  for  $\mathbf{x} = \varphi^k(\mathbf{y})$  with  $\mathbf{y}$  aperiodic. Our reflection is here driven by the following questions inspired by Lejeune's questions [28, pp. 115–117] that are natural to consider in view of (2).

**Question A.** Does there exist an infinite word  $\mathbf{w}$  such that, for all  $k \geq 1$ ,  $\mathbf{b}_w^{(k)}$  is unbounded and  $\mathbf{b}_w^{(k)} \prec \mathbf{b}_w^{(k+1)}$ ? If the answer is positive, can we find a (pure) morphic such word  $\mathbf{w}$ ?

From (2), notice that  $\mathbf{b}_w^{(k)}$  is unbounded, for all  $k \geq 1$ , if and only if the abelian complexity  $\mathbf{b}_w^{(1)}$  is unbounded. Even though the Thue–Morse word  $\mathbf{t}$  is such that, for all  $k \geq 1$ ,  $\mathbf{b}_t^{(k)} \prec \mathbf{b}_t^{(k+1)}$ ,  $\mathbf{b}_t^{(k)}$  remains bounded (3). So  $\mathbf{t}$  is not a satisfying answer to Question A. However, in Section 6, we provide several positive answers to this question.

**Question B.** For each  $\ell \geq 1$ , does there exist a word  $\mathbf{w}$  (depending on  $\ell$ ) such that  $\mathbf{b}_w^{(1)} \prec \mathbf{b}_w^{(2)} \prec \dots \prec \mathbf{b}_w^{(\ell-1)} \prec \mathbf{b}_w^{(\ell)} = p_w$ ? If the answer is positive, is there a (pure) morphic such word  $\mathbf{w}$ ?

Putting together results from Sections 4 and 7 we fully answer Question B: Theorem 4.2 and Proposition 4.17 provide a word  $\mathbf{x} = \varphi^k(\mathbf{y})$  for which  $\mathbf{b}_x^{(1)} \prec \mathbf{b}_x^{(2)} \prec \dots \prec \mathbf{b}_x^{(k-1)} \prec \mathbf{b}_x^{(k)} \prec \mathbf{b}_x^{(k+1)}$ , while assuming that  $\mathbf{y}$  above is Sturmian, we show that  $\mathbf{b}_x^{(k+2)} = p_x$ . We remark that iterates of  $\varphi$  applied to Sturmian words have been studied (among other words) in [21]. We observe that our construction leads to words with bounded abelian complexity. Question B is then strengthened in Section 7 where we ask for words with unbounded abelian complexity. We give a pure morphic answer when  $\ell = 3$ .

## 2 Preliminaries

Let us now give precise definitions and notation. For any integer  $k$ , we let  $A^k$  (resp.,  $A^{\leq k}$ ; resp.,  $A^{<k}$ ) denote the set of words of length exactly (resp., at most; resp., less than)  $k$  over  $A$ . We let  $A^*$  (resp.,  $A^+$ ) denote the semigroup of finite words (resp., non-empty finite words) over  $A$  equipped with concatenation. We let  $\varepsilon$  denote the empty word. The length of the word  $w$  is denoted by  $|w|$  and the number of occurrences of a letter  $a$  in  $w$  is denoted by  $|w|_a$ . For binary words  $u, v$  (always over  $\{0, 1\}$  in this note, unless otherwise stated), we refer to  $|u|_1$  as the *weight* of  $u$  and we say that  $u$  is *lighter* (resp., *heavier*) than  $v$  whenever  $|u|_1 < |v|_1$  (resp.,  $|u|_1 > |v|_1$ ). For instance, if  $\mathbf{b}_y^{(1)}(n) = 2$ , then there are only two kinds of factors in  $\mathbf{y}$ : the light ones and the heavy ones. A language  $L$  is said to be *balanced* if, for all words  $u, v \in L$  of the same length and all letters  $a$ , we have  $||u|_a - |v|_a| \leq 1$ . In particular, an infinite word  $\mathbf{z}$  is *balanced* if  $\mathcal{L}(\mathbf{z})$  is balanced.

We let  $\bar{\cdot}$  denote the (binary) complementation morphism defined by  $\bar{a} = 1 - a$ , for  $a \in \{0, 1\}$ . Writing  $A = \{a_1, \dots, a_k\}$  and fixing the order  $a_1 < a_2 < \dots < a_k$  on the letters, the *Parikh vector* of a word  $w \in A^*$  is defined as the column vector

$$\Psi(w) = (|w|_{a_1}, |w|_{a_2}, \dots, |w|_{a_k})^T.$$

Using a classical “length- $n$  sliding window” argument or extending factors of length  $n$  to factors of length  $n + 1$ , one has the following.

**Lemma 2.1** (Folklore). *For any binary word  $y$  over  $\{0, 1\}$ , for all  $n \geq 0$ , we have*

$$b_y^{(1)}(n) = 1 + \max_{u, v \in \mathcal{L}_n(y)} ||u|_1 - |v|_1| \quad \text{and} \quad |b_y^{(1)}(n+1) - b_y^{(1)}(n)| \leq 1.$$

## 2.1 Binomial Equivalence

We first collect some useful results on  $k$ -binomial equivalence. Note that  $\sim_k$  is a congruence, i.e., for  $u, v, x, y \in A^*$ ,  $u \sim_k v$  and  $x \sim_k y$  implies  $ux \sim_k vy$ . In particular,  $A^*/\sim_k$  is a monoid. In fact, it is a cancellative monoid (see [29, Lemma 10]; cancellativity also follows from  $A^*/\sim_k$  being isomorphic to a subsemigroup of the special linear group  $SL((k+1)n^k, \mathbb{Z})$  where  $n$  is the cardinality of the alphabet  $A$  [43]):

**Lemma 2.2** (Cancellation property). *Let  $u, v, w$  be words over  $A$ . We have*

$$v \sim_k w \Leftrightarrow uv \sim_k uw \quad \text{and} \quad v \sim_k w \Leftrightarrow vu \sim_k wu.$$

We will also need the following result characterizing  $k$ -binomial commutation among words of equal length.

**Theorem 2.3** ([48, Thm. 3.5]). *Let  $k \geq 2$  and  $x, y \in A^*$  such that  $|x| = |y|$ . Then  $xy \sim_k yx$  if and only if  $x \sim_{k-1} y$ .*

A proof of the next result can be conveniently found in [29, Lem. 30].

**Theorem 2.4** (Ochsenschläger [36]). *Let  $\phi: 0 \mapsto 01, 1 \mapsto 10$  be the Thue–Morse morphism. For all  $k \geq 1$ , we have  $\phi^k(0) \sim_k \phi^k(1)$  and  $\phi^k(0) \not\sim_{k+1} \phi^k(1)$ .*

The following result from [29, Lem. 31] will be of use. It can alternatively be proved using Theorem 2.3 combined with Ochsenschläger’s result.

**Lemma 2.5** (Transfer lemma). *Let  $k \geq 1$ . Let  $u, v, v'$  be three non-empty words such that  $|v| = |v'|$ . We have  $\phi^{k-1}(u)\phi^k(v) \sim_k \phi^k(v')\phi^{k-1}(u)$ .*

It is an exercise to see that, for an arbitrary morphism  $f: A^* \rightarrow B^*$ , we have, for all  $u \in A^*$ ,  $e \in B^*$ ,

$$\binom{f(u)}{e} = \sum_{\substack{a_1, \dots, a_\ell \in A \\ \ell \leq |e|}} \binom{u}{a_1 \cdots a_\ell} \sum_{\substack{e = e_1 \cdots e_\ell \\ e_i \in B^+}} \prod_{i=1}^{\ell} \binom{f(a_i)}{e_i}. \quad (4)$$

The next result will turn out to be useful in several places of the paper.

**Lemma 2.6.** *Let  $x, y \in A^*$  be two  $k$ -binomially equivalent words. For any integer  $n \geq 0$  and any word  $e \in A^*$  of length  $k + 1$ , we have*

$$\binom{x^n}{e} - \binom{y^n}{e} = n \left[ \binom{x}{e} - \binom{y}{e} \right].$$

*In particular, for all  $n \geq 1$ ,  $x \sim_{k+1} y$  if and only if  $x^n \sim_{k+1} y^n$ .*

*Proof.* For any words  $u, v, w \in A^*$ , we have

$$\binom{uv}{w} = \binom{u}{w} + \binom{v}{w} + \sum_{\substack{w = w_1 w_2 \\ w_i \neq \varepsilon}} \binom{u}{w_1} \binom{v}{w_2}.$$

To show the statement, we proceed by induction and we make use of the previous formula. The statement is trivially true for  $n \in \{0, 1\}$ . By the previous formula (with  $(u, v, w) = (x^n, x, e)$  and  $(u, v, w) = (y^n, y, e)$  respectively) and the induction hypothesis, we obtain

$$\binom{x^{n+1}}{e} - \binom{y^{n+1}}{e} = (n+1) \left[ \binom{x}{e} - \binom{y}{e} \right] + \sum_{\substack{e=e_1 e_2 \\ e_i \neq \varepsilon}} \left[ \binom{x^n}{e_1} \binom{x}{e_2} - \binom{y^n}{e_1} \binom{y}{e_2} \right].$$

Since  $x \sim_k y$  and  $|e| = k+1$ , the sum in the right-hand term is zero and we obtain the desired result.  $\square$

We recall the following lemma that appears in [48]; it is a straightforward generalization of an observation in [46]. We give a proof for the sake of completeness.

**Lemma 2.7.** *Let  $C \in A^*/\sim_1$  be an abelian equivalence class of non-empty words with Parikh vector  $(m_a)_{a \in A}$ . Then, for any word  $u \in A^*$ , we have  $\sum_{w \in C} \binom{u}{w} = \prod_{a \in A} \binom{|u|_a}{m_a}$ .*

*Proof.* The sum on the left counts the number of ways one can choose a subword  $w$  of  $u$  so that  $\Psi(w) = (m_a)_{a \in A}$ . On the other hand, for a vector  $(m_a)_{a \in A}$ , any choice of  $m_a$  many distinct  $a$ 's in  $u$  for each  $a \in A$  gives rise to a subword of  $u$  having Parikh vector  $(m_a)_{a \in A}$ . The number of distinct such choices is the product on the right.  $\square$

## 2.2 Binomial equivalence in Sturmian words

The following result links the factor complexity and the 2-binomial complexity of Sturmian words.

**Theorem 2.8** ([43, Thm. 7]). *For any Sturmian word  $s$ , we have  $b_s^{(2)} = p_s$ .*

In particular, the theorem implies that for two distinct equal-length factors  $u, v$  of a Sturmian word, we have either  $u \not\sim_1 v$ , or  $\binom{u}{01} \neq \binom{v}{01}$ . It further implies that  $b_s^{(k)}(n) = n+1$  for all  $k \geq 2$ . In the survey paper [18] on abelian combinatorics on words, Fici and Puzynina derive a characterization of Sturmian words from Theorem 2.8:

**Theorem 2.9** ([18, Rem. 80]). *Let  $x$  be an infinite word. The following are equivalent:*

1.  $x$  is Sturmian;
2. for all  $n \geq 1$  and some  $k \geq 2$ ,  $b_x^{(1)}(n) = 2$  and  $b_x^{(k)}(n) = n+1$ ;
3. for all  $n \geq 1$  and  $k \geq 2$ ,  $b_x^{(1)}(n) = 2$  and  $b_x^{(k)}(n) = n+1$ .

In fact, the second property in the above can be weakened to " $b_x^{(1)}(n) = 2$  for all  $n \geq 1$  and  $\sup_{k, n \in \mathbb{N}} b_x^{(k)}(n) = \infty$ " using the same arguments<sup>1</sup>. In the following we show that the assumption of balancedness can be removed from the second point; Sturmian words are characterized by their  $k$ -binomial complexity for any fixed integer  $k \geq 2$ . We first recall the following crucial observation, which can be found in part of [15, Lem. 4.02<sup>2</sup>]

**Lemma 2.10.** *Let  $z$  be an infinite binary word. Let  $N \geq 2$  be such that*

1.  $\mathcal{L}(z) \cap A^{<N}$  is balanced;
2.  $\mathcal{L}(z) \cap A^N$  is unbalanced;
3.  $p_z(N) = N+1$ ;

*Then  $z$  is ultimately periodic.*

<sup>1</sup>Indeed, the former property implies  $x$  is balanced and binary, and the latter implies that  $x$  is aperiodic.

<sup>2</sup>The statement has a fourth condition, which does not affect the conclusion appearing here.

**Theorem 2.11.** *Let  $\mathbf{z}$  be an infinite word such that for some  $k \geq 2$ ,  $b_{\mathbf{z}}^{(k)} = n + 1$  for all  $n$ . Then  $\mathbf{z}$  is Sturmian.*

*Proof.* Assume that  $b_{\mathbf{z}}^{(k)}(n) = n + 1$  for all  $n \geq 0$ . In particular,  $\mathbf{z}$  is binary and also aperiodic because  $p_{\mathbf{z}}(n) \geq b_{\mathbf{z}}^{(k)}(n) = n + 1$ . This also implies  $b_{\mathbf{z}}^{(1)}(n) \geq 2$  for all  $n \geq 1$ . To get a contradiction, assume that  $\mathbf{z}$  is unbalanced. Hence there exists a minimal integer  $N \geq 2$  such that  $b_{\mathbf{z}}^{(1)}(N) = 3$ . There is a pair  $(u, v)$  of factors of  $\mathbf{z}$  of length  $N$  such that  $|u|_1 - |v|_1 = 2$ . The minimality of  $N$  implies that this pair is unique and of the form  $(1w1, 0w0)$ . For details, see [15, Lem. 3.06].

Let  $|w|_1 = r$ . Let  $x \in \mathcal{L}_{N-2}(\mathbf{z})$ . If  $|x|_1 = r - 1$ , then  $0x$  and  $x0$  do not belong to  $\mathcal{L}_{N-1}(\mathbf{z})$ . Indeed,  $1w, w1$  belong to the latter set and  $|1w|_1 = |w1|_1 = r + 1$  but by minimality of  $N$ , the set  $\mathcal{L}_{N-1}(\mathbf{z})$  is balanced. In that case,  $x$  is preceded and followed by 1. Similarly, if  $y \in \mathcal{L}_{N-2}(\mathbf{z})$  and  $|y|_1 = r + 1$ , then  $y$  is preceded and followed by 0. This means that  $\mathcal{L}_N(\mathbf{z}) \setminus \{0w0, 1w1\}$  is a subset of

$$\{1x1 : |x|_1 = r - 1, |x| = N - 2\} \cup \{0y0 : |y|_1 = r + 1, |y| = N - 2\} \cup \bigcup_{a \in \{0,1\}} \{az\bar{a} : |z|_1 = r, |z| = N - 2\}$$

where all words have weight  $r + 1$ . Now, observe that  $\mathcal{L}(\mathbf{z}) \cap A^{<N}$  is a balanced set (by minimality of  $N$ ) and that is also the case of  $\mathcal{L}_N(\mathbf{z}) \setminus \{1w1\}$ . The union of these two sets is factorial. By [41, Thm. 3.1], there exists a Sturmian word  $\mathbf{s}$  such that

$$(\mathcal{L}(\mathbf{z}) \cap A^{<N}) \cup (\mathcal{L}_N(\mathbf{z}) \setminus \{1w1\}) \subset \mathcal{L}(\mathbf{s}).$$

As a consequence of Theorem 2.8, any two distinct words in the left-hand side set are not  $k$ -binomially equivalent. Also,  $1w1$  is not abelian (and thus not  $k$ -binomially) equivalent to any word in  $\mathcal{L}_N(\mathbf{z}) \setminus \{1w1\}$ . In particular, since  $b_{\mathbf{z}}^{(k)}(n) = n + 1$  for  $n \leq N$ ,

$$\mathcal{L}(\mathbf{z}) \cap A^{<N} = \mathcal{L}(\mathbf{s}) \cap A^{<N}$$

and  $\#(\mathcal{L}_N(\mathbf{z}) \setminus \{1w1\}) = N$ . Therefore,  $\#(\mathcal{L}_N(\mathbf{z})) = N + 1$ . The word  $\mathbf{z}$  now fulfills all the conditions of Lemma 2.10 implying the contradiction that  $\mathbf{z}$  is ultimately periodic.  $\square$

### 3 Parikh-Collinear Morphisms via Binomial Complexities

In this section, we obtain a new characterization of Parikh-collinear morphisms and show that, given an infinite fixed point of a prolongable Parikh-collinear morphism, its  $k$ -binomial complexity is bounded for each  $k$ . Note that the automaticity of such fixed points is discussed in [45].

**Definition 3.1** (Parikh-collinear morphisms). A morphism  $f: A^* \rightarrow B^*$  is said to be *Parikh-collinear* if, for all letters  $a, b \in A$ , there is  $r_{a,b} \in \mathbb{Q}$  such that  $\Psi(f(b)) = r_{a,b}\Psi(f(a))$ .

**Remark 3.2.** Given a morphism  $f: A^* \rightarrow B^*$ , its *adjacency matrix*  $M_f$  is the matrix of size  $\#B \times \#A$  defined by  $(M_f)_{b,a} = |f(a)|_b$  for all  $a \in A, b \in B$ . Observe that  $f$  is a Parikh-collinear morphism if and only if  $M_f$  has rank 1 (unless it is totally erasing). We observe that for any word  $u \in A^*$ , we have that  $\Psi(f(u)) = M_f\Psi(u)$ .

**Example 3.3.** The morphism  $f$  defined by  $0 \mapsto 000111; 1 \mapsto 0110$  is Parikh-collinear since  $\Psi(f(1)) = \frac{2}{3}\Psi(f(0))$ .

**Theorem 3.4** ([14, Thm. 11]). *A morphism  $f: A^* \rightarrow B^*$  is Parikh-collinear if and only if it maps all infinite words to words with bounded abelian complexity.*

We extend the above theorem to the following one. We say that a morphism  $f: A^* \rightarrow B^*$  satisfies  $P_k$  if  $f$  maps all words with bounded  $k$ -binomial complexity to words with bounded  $(k + 1)$ -binomial complexity. Note that, for  $k = 0$ , *0-binomial complexity* has to be understood as the “equal length” equivalence relation. So the 0-binomial complexity of an infinite word is the constant function 1 and Theorem 3.4 can be restated as  *$f$  is Parikh-collinear if and only if  $f$  satisfies  $P_0$ .*



**Theorem 3.5.** *Let  $f: A^* \rightarrow B^*$  be a morphism. The following are equivalent.*

- (i) *The morphism  $f$  is Parikh-collinear.*
- (ii) *For all  $k \geq 0$ ,  $f$  satisfies  $P_k$ .*
- (iii) *There exists an integer  $k \geq 0$  such that  $f$  satisfies  $P_k$ .*

Before proving this result in Section 3.2, let us mention a straightforward consequence, which generalizes [43, Thm. 13] from Parikh-constant to Parikh-collinear morphisms. For example, the Thue–Morse morphism is Parikh-constant and thus Parikh-collinear but the morphism of Example 3.3 is Parikh-collinear but not Parikh-constant.

**Corollary 3.6.** *Let  $\mathbf{z}$  be a fixed point of a Parikh-collinear morphism. For any  $k \geq 1$  there exists a constant  $C_{\mathbf{z},k} \in \mathbb{N}$  such that  $\mathbf{b}_{\mathbf{z}}^{(k)}(n) \leq C_{\mathbf{z},k}$  for all  $n \in \mathbb{N}$ .*

*Proof.* Let  $f: A^* \rightarrow A^*$  be a Parikh-collinear morphism whose fixed point is  $\mathbf{z}$ . Since  $f(\mathbf{z}) = \mathbf{z}$ , Theorem 3.4 implies that  $\mathbf{z}$  has bounded abelian complexity. For any  $k \geq 1$ , we have that  $\mathbf{z} = f(f^{k-1}(\mathbf{z}))$  implying that  $\mathbf{z}$  has bounded  $k$ -binomial complexity by induction and the previous theorem.  $\square$

**Remark 3.7.** We cannot relax the (implicit) assumption on the rank of the adjacency matrix  $M_f$  in Corollary 3.6. For example, the morphism  $f: \{0, 1, 2\}^* \rightarrow \{0, 1, 2\}^*$  defined by  $0 \mapsto 0^3 2^3$ ,  $1 \mapsto 0^3 1^3 2$ ,  $2 \mapsto 2^4 0^6 1^3$  has an adjacency matrix of rank 2. The fixed point  $\mathbf{x}$  starting with 0 is aperiodic as  $f^n(0)$  is readily seen to be right special for all  $n \geq 0$ . Yet, its adjacency matrix has eigenvalues  $\theta_1 = 5 + \sqrt{13}$ ,  $\theta_2 = 5 - \sqrt{13}$ , and 0, and the former two are greater than 1. This means that the word has unbounded abelian complexity. Indeed, a deep result of Adamczewski on balances in primitive pure morphic words [1, Thm. 13(ii)] implies that the “lim sup”-growth of the function

$$n \mapsto \max_{a \in \Sigma, u, v \in \mathcal{L}_n(\mathbf{x})} \{ |u|_a - |v|_a \}$$

grows as  $\Theta(n^{\log_{\theta_1} \theta_2})$ , where  $\log_{\theta_1} \theta_2 \approx 0,15448$ . It follows (see, e.g., [40, Lem. 2.2]) that  $\mathbf{b}_{\mathbf{x}}^{(k)}$  is unbounded for each  $k \geq 1$ .

### 3.1 An Intermediate Characterization of Parikh-Collinearity

To prove Theorem 3.5, we give further characterizations of Parikh-collinear morphisms. To this end, we require the following lemma where we define a map  $g_e$  which is constant on any abelian equivalence class. Notice that such a map appears within (4).

**Lemma 3.8.** *Let  $A, B$  be finite alphabets with  $\#A \geq 2$ . Let  $f: A^* \rightarrow B^*$  be a Parikh-collinear morphism. For a word  $e = e_1 \cdots e_n$  of length  $n$  over  $B$ , define  $g_e: A^n \rightarrow \mathbb{N}$  by*

$$g_e(a_1 \cdots a_n) := \prod_{i=1}^n \binom{f(a_i)}{e_i}.$$

*Then, for all words  $w, w' \in A^n$  with  $w \sim_1 w'$ , we have  $g_e(w) = g_e(w')$ .*

*Proof.* Write  $w = a_1 \cdots a_n$  with  $a_i \in A$  for all  $i \in \{1, \dots, n\}$ . For all  $\alpha \in A$  and  $\beta \in B$ , define  $I(\alpha, \beta) := \{i \in \{1, \dots, n\} \mid a_i = \alpha \text{ and } e_i = \beta\}$ . We get

$$g_e(w) = \prod_{\substack{\alpha \in A \\ \beta \in B}} \prod_{i \in I(\alpha, \beta)} \binom{f(\alpha)}{\beta}.$$

The claim is trivial if  $f$  maps all words to  $\varepsilon$ , so let  $0 \in A$  be a letter for which  $|f(0)| \neq 0$ . Since the morphism  $f$  is Parikh-collinear, for all  $\alpha \in A$  and all  $\beta \in B$ , there exists  $r_\alpha \in \mathbb{Q}$  such that  $\binom{f(\alpha)}{\beta} = r_\alpha \binom{f(0)}{\beta}$ . We now get

$$\begin{aligned} g_e(w) &= \prod_{\substack{\alpha \in A \\ \beta \in B}} \prod_{i \in I(\alpha, \beta)} \binom{f(\alpha)}{\beta} = \prod_{\substack{\alpha \in A \\ \beta \in B}} \prod_{i \in I(\alpha, \beta)} r_\alpha \binom{f(0)}{\beta} \\ &= \left( \prod_{\substack{\alpha \in A \\ \beta \in B}} \prod_{i \in I(\alpha, \beta)} \binom{f(0)}{\beta} \right) \left( \prod_{\substack{\alpha \in A \\ \beta \in B}} \prod_{i \in I(\alpha, \beta)} r_\alpha \right). \end{aligned}$$

For any letter  $\beta \in B$ , the definition of  $I(\alpha, \beta)$  gives

$$\prod_{\alpha \in A} \prod_{i \in I(\alpha, \beta)} \binom{f(0)}{\beta} = \binom{f(0)}{\beta}^{|\varepsilon|_\beta}.$$

Similarly, for any letter  $\alpha \in A$ , the definition of  $I(\alpha, \beta)$  yields

$$\prod_{\beta \in B} \prod_{i \in I(\alpha, \beta)} r_\alpha = r_\alpha^{|\varepsilon|_\alpha}.$$

Thus

$$g_e(w) = \left( \prod_{\beta \in B} \binom{f(0)}{\beta}^{|\varepsilon|_\beta} \right) \left( \prod_{\alpha \in A} r_\alpha^{|\varepsilon|_\alpha} \right).$$

Observe that the first factor in this product only depends on (the Parikh vector of)  $e$  — in particular, not on  $w$  — as the morphism  $f$  is fixed. Similarly, the second factor in the product depends solely on the Parikh vector of  $w$ , not on the word  $w$  itself. The desired result follows.  $\square$

We now characterize Parikh-collinear morphisms by means of binomial complexities.

**Proposition 3.9.** *Let  $f: A^* \rightarrow B^*$  be a morphism. The following are equivalent.*

- (i) *The morphism  $f$  is Parikh-collinear.*
- (ii) *For all  $k \geq 2$  and  $u, v \in A^*$ ,  $u \sim_{k-1} v$  implies  $f(u) \sim_k f(v)$ .*
- (iii) *There exists an integer  $k \geq 2$  such that for all  $u, v \in A^*$ ,  $u \sim_{k-1} v$  implies  $f(u) \sim_k f(v)$ .*
- (iv) *For all  $u, v \in A^*$ ,  $u \sim_1 v$  implies  $f(u) \sim_2 f(v)$ .*

*Proof.* Clearly (ii) implies (iii). We show that (iii) implies (iv). There is nothing to prove if (iii) holds for  $k = 2$ , so assume that  $k \geq 3$ . We show that  $f$  also satisfies (iii) with  $k - 1$  instead of  $k$ , and hence, by repeating the argument,  $f$  satisfies (iii) with  $k = 2$ . Assume to the contrary that there exists a pair  $u, v$  such that  $u \sim_{k-2} v$  but  $f(u) \not\sim_{k-1} f(v)$ . Since  $u$  and  $v$  are abelian equivalent ( $k - 2 \geq 1$ ) they have equal length, so by Theorem 2.3, we have that  $uv \sim_{k-1} vu$ . Then, since  $f$  has the property for  $k$ , we have  $f(u)f(v) \sim_k f(v)f(u)$ . Furthermore,  $f(u)$  and  $f(v)$  have the same length (due to  $u \sim_1 v$ ). This implies that  $f(u) \sim_{k-1} f(v)$  by the converse part of Theorem 2.3, contrary to what was assumed.

Assuming (iv), we show that (i) holds. Let  $x, y$  be distinct letters from  $A$ . Since  $xy \sim_1 yx$ , we have  $f(xy) \sim_2 f(yx)$  by assumption. In other words, for all  $s, t \in B$  we have, applying (4),

$$\begin{aligned} 0 &= \binom{f(xy)}{st} - \binom{f(yx)}{st} = \sum_{\substack{\alpha_1, \dots, \alpha_\ell \in A \\ \ell \leq 2}} \left[ \binom{xy}{\alpha_1 \cdots \alpha_\ell} - \binom{yx}{\alpha_1 \cdots \alpha_\ell} \right] \sum_{\substack{st = b_1 \cdots b_\ell \\ b_i \in B^+}} \prod_{i=1}^{\ell} \binom{f(\alpha_i)}{b_i} \\ &= \sum_{\alpha_1, \alpha_2 \in A} \left( \binom{xy}{\alpha_1 \alpha_2} - \binom{yx}{\alpha_1 \alpha_2} \right) \binom{f(\alpha_1)}{s} \binom{f(\alpha_2)}{t} = \binom{f(x)}{s} \binom{f(y)}{t} - \binom{f(y)}{s} \binom{f(x)}{t}, \end{aligned}$$

where in the third equality we use  $\binom{xy}{a} = \binom{yx}{a}$  for all  $a \in A$  (since  $xy \sim_1 yx$ ). Summing over  $s \in B$ , we get  $|f(x)|\binom{f(y)}{t} = |f(y)|\binom{f(x)}{t}$  for all  $t \in B$ . Now  $x$  and  $y$  were chosen arbitrarily from the alphabet  $A$ . If  $|f(x)| = 0$  for all  $x \in A$ , then  $f$  is clearly Parikh-collinear. If there is a letter  $x$  for which  $|f(x)| > 0$ , we may write  $\left(\binom{f(y)}{t}\right)_{t \in B} = \frac{|f(y)|}{|f(x)|} \left(\binom{f(x)}{t}\right)_{t \in B}$  for each  $y \in A$ . In other words,  $f$  is Parikh-collinear.

To complete the proof, we show that (i) implies (ii). So let  $f$  be a Parikh-collinear morphism and  $u \sim_{k-1} v$  with  $k \geq 2$ . We again apply (4): for any word  $e \in B^*$ , we have

$$\binom{f(u)}{e} - \binom{f(v)}{e} = \sum_{\substack{a_1, \dots, a_\ell \in A \\ \ell \leq |e|}} \left( \binom{u}{a_1 \cdots a_\ell} - \binom{v}{a_1 \cdots a_\ell} \right) \sum_{\substack{e = e_1 \cdots e_\ell \\ e_i \in B^+}} \prod_{i=1}^{\ell} \binom{f(a_i)}{e_i}.$$

Notice that for words  $e \in B^{<k}$ , we have  $\binom{u}{a_1 \cdots a_\ell} = \binom{v}{a_1 \cdots a_\ell}$  since  $u \sim_{k-1} v$ , which in turn gives  $\binom{f(u)}{e} = \binom{f(v)}{e}$ . So to show that  $f(u) \sim_k f(v)$ , it suffices to consider words  $e \in B^k$ . By assumption, for  $\ell < k$ , we again have  $\binom{u}{a_1 \cdots a_\ell} = \binom{v}{a_1 \cdots a_\ell}$ . Therefore, we have  $\binom{f(u)}{e} = \binom{f(v)}{e}$  if and only if

$$\sum_{a_1, \dots, a_k \in A} \binom{u}{a_1 \cdots a_k} \prod_{i=1}^k \binom{f(a_i)}{e_i} = \sum_{a_1, \dots, a_k \in A} \binom{v}{a_1 \cdots a_k} \prod_{i=1}^k \binom{f(a_i)}{e_i}. \quad (5)$$

Observe here that  $\prod_{i=1}^k \binom{f(a_i)}{e_i} = g_e(a_1 \cdots a_k)$  as defined in Lemma 3.8. Let  $\mathcal{C}$  be an abelian equivalence class in  $A^k / \sim_1$ . By Lemma 3.8,  $g_e(\cdot)$  is constant on  $\mathcal{C}$ , so write  $g_e(w) = g_{\mathcal{C}, e}$  for all words  $w \in \mathcal{C}$ . For each  $w \in \mathcal{C}$  we may write  $\Psi(w) = (m_{\mathcal{C}, a})_{a \in A}$ . We now have

$$\sum_{w \in A^k} \binom{u}{w} g_e(w) = \sum_{\mathcal{C} \in A^k / \sim_1} \sum_{w \in \mathcal{C}} \binom{u}{w} g_e(w) = \sum_{\mathcal{C} \in A^k / \sim_1} g_{\mathcal{C}, e} \sum_{w \in \mathcal{C}} \binom{u}{w} = \sum_{\mathcal{C} \in A^k / \sim_1} g_{\mathcal{C}, e} \prod_{a \in A} \binom{|u|_a}{m_{\mathcal{C}, a}},$$

where the last equality is from Lemma 2.7. One obtains the same formula by replacing  $u$  with  $v$ , and equality indeed holds in (5) as  $|u|_a = |v|_a$  for each letter  $a \in A$ . This concludes the proof.  $\square$

**Remark 3.10.** In [34, Lem. 5], the authors show that, for a morphism  $f$  such that  $f(a) \sim_h f(b)$  for all  $a, b \in A$ , for all words  $u, v \in A^*$  with  $u \sim_k v$  we have that  $f(u) \sim_{k+h} f(v)$ . Towards the converse, assume that  $f$  is a morphism for which the conclusion holds (for all  $k \geq 1$  but fixed  $h \geq 1$ ). Then we necessarily have  $f(a)^m f(b)^n \sim_{h+1} f(b)^n f(a)^m$  for all  $a, b \in A$ ,  $m, n \geq 1$ . From this we infer that, e.g.,  $f(a)^{|f(b)|} \sim_h f(b)^{|f(a)|}$  (as a corollary of Theorem 2.3). In particular, if  $f$  is uniform, we have  $f(a) \sim_h f(b)$  for all  $a, b \in A$ . It would be interesting to characterize the non-uniform morphisms  $f$  with this property. For example, one can take any Parikh-collinear morphism  $g$ ; then  $f = g^h$  is such a morphism. We highly suspect that these are not the only such morphisms.

## 3.2 Proof of Theorem 3.5

We require the following technical result, which essentially appears in the proof of [14, Thm. 12]. We give a proof here for the sake of completeness.

**Lemma 3.11.** *Let  $\mathbf{x}$  be an infinite word over  $A$  with bounded abelian complexity. Let  $f : A^* \rightarrow B^*$  be a morphism and assume  $\mathbf{y} = f(\mathbf{x})$  is an infinite word. Then for all  $c \in \mathbb{N}$  there exists  $D_{\mathbf{x}, c} \in \mathbb{N}$  such that if  $\|f(u) - f(v)\| \leq c$ , for some  $u, v \in \mathcal{L}(\mathbf{x})$ , then  $\|u - v\| \leq D_{\mathbf{x}, c}$ .*

*Proof.* Assume without loss of generality that  $|u| \geq |v|$  and write  $u = u'v'$  with  $|v'| = |v|$ . Let  $M_f$  be the adjacency matrix of  $f$ . If  $\|f(u) - f(v)\| \leq c$ , we have by the reverse triangle inequality

$$c \geq \|f(u') - f(v) + f(v')\| \geq \|f(u') - f(v') - f(v)\| = \|f(u') - \langle M_f(\Psi(v') - \Psi(v)), \vec{1} \rangle\|,$$

where  $\langle \cdot, \cdot \rangle$  denotes the inner product of vectors, and  $\vec{1}$  is the all-ones-vector. Recall that  $x$  has bounded abelian complexity if and only if it is  $C$ -balanced for some  $C$  [40]. Hence, as  $v$  and  $v'$  are factors of the same length,  $\Psi(v') - \Psi(v)$  attains finitely many distinct integer points (in particular, belonging to  $[-C, C]^{\#A}$ ). So does  $M_f(\Psi(v') - \Psi(v))$ . We therefore obtain  $|f(u')| \leq D$  for some  $D \in \mathbb{N}$ . We deduce that  $u'$  is bounded in length as well: indeed, let  $a \in A$  be a letter occurring infinitely often in  $x$  and for which  $f(a) \neq \varepsilon$  (such a letter exists because  $f(x)$  is infinite). Since  $x$  is balanced, we deduce that all long enough factors of  $x$  contain more than  $|u'|$  occurrences of  $a$ . We let  $D_{x,c}$  be this bound on  $|u'|$  to conclude the proof.  $\square$

We are now ready to prove the main result of this section, characterizing Parikh-collinear morphisms in terms the property  $P_k$  defined at the beginning of Section 3.

*Proof of Theorem 3.5.* Let us first show that (i) implies (ii). Assume thus that  $f$  is Parikh-collinear. Theorem 3.4 implies that  $f$  maps all words (i.e., all words with bounded 0-binomial complexity) to words with bounded 1-binomial complexity. Let  $k \geq 1$  and let  $x$  be a word with bounded  $k$ -binomial complexity. Let  $n \in \mathbb{N}$ . Any length- $n$  factor of  $f(x)$  can be written as  $pf(u)s$ , where the word  $u$  is a factor of  $x$ ,  $p$  is a suffix of  $f(a)$  and  $s$  is a prefix of  $f(b)$  for some letters  $a, b \in A$ . Here  $n - 2m < |f(u)| \leq n$ , where  $m := \max_{a \in A} |f(a)|$ . The  $(k+1)$ -binomial equivalence class of  $pf(u)s$  is completely determined by the words  $p, s$ , and the  $k$ -binomial equivalence class of  $f(u)$ , which itself is determined by the abelian equivalence class of  $u$  by Proposition 3.9.

The former two words  $p$  and  $s$  are drawn from a finite set, as their lengths are bounded by the constant  $m$  (depending on  $f$ ). The length of  $u$  can be chosen from an interval whose length is uniformly bounded in  $n$ . Indeed, assume we have equal length factors  $w = pf(u)s$  and  $w' = p'f(v)s'$ . As observed above,  $n \geq |f(u)|$  and  $|f(v)| > n - 2m$ , so that  $||f(u)| - |f(v)|| < 2m$ . Applying Lemma 3.11 (by assumption,  $x$  has bounded  $k$ -binomial complexity and thus,  $x$  has bounded abelian complexity by (2)) there exists a bound  $D$  such that  $||u| - |v|| \leq D$  uniformly in  $n$ . Since the number of  $k$ -binomial equivalence classes in  $x$  of each length is uniformly bounded by assumption, and the number of admissible lengths for  $u$  above is bounded, we conclude that the number of choices for the  $k$ -binomial equivalence class of  $u$  is bounded. We have shown that the number of  $(k+1)$ -binomial equivalence classes among factors of length  $n$  in  $f(x)$  is determined from a bounded amount of information (not depending on  $n$ ), as was to be shown. Consequently,  $f$  satisfies  $P_k$ .

Notice that (ii) trivially implies (iii).

Let us turn to the last implication, namely (iii) implies (i). Assume (iii) holds, that is, for some integer  $k \geq 0$ ,  $f$  satisfies  $P_k$ . If  $k = 0$ , then  $f$  maps all words to words with bounded 1-binomial complexity, so  $f$  is Parikh-collinear by Theorem 3.4. Assume that  $k \geq 1$ , and towards a contradiction, assume further that  $f$  is not Parikh-collinear. By Proposition 3.9, there exist words  $u, v$  with  $u \sim_k v$  and  $f(u) \not\sim_{k+1} f(v)$ . Write  $U = f(u)$  and  $V = f(v)$ . Now define the word  $x = uvu^2v^2u^3v^3 \dots u^n v^n \dots$  and consider

$$f(x) = UVU^2V^2U^3V^3 \dots U^n V^n \dots .$$

Below we show that  $x$  has bounded  $k$ -binomial complexity, while  $f(x)$  has unbounded  $(k+1)$ -binomial complexity, which is enough to contradict (iii).

Since  $u$  and  $v$  are  $k$ -binomially equivalent,  $b_x^{(k)}$  is bounded. (To see this, one may apply arguments similar to those developed in the first part of the proof.) Let us prove the second. For each integer  $n$ ,  $f(x)$  contains the factors  $U^r V^{n-r}$  with  $r \in \{0, \dots, n\}$ . These factors are actually all  $(k+1)$ -binomially inequivalent. Indeed, assume towards a contradiction that  $U^r V^{n-r} \sim_{k+1} U^s V^{n-s}$  for some  $r, s \in \{0, \dots, n\}$  with  $r > s$ . By the Cancellation property (Lemma 2.2), we obtain  $U^{r-s} \sim_{k+1} V^{r-s}$ . Lemma 2.6 then implies that  $U \sim_{k+1} V$ , which is a contradiction. Consequently,  $b_{f(x)}^{(k+1)}$  is unbounded, as desired.  $\square$

## 4 Binomial Properties of the Thue–Morse Morphism, Part I

In this section, we consider binomial complexities of iterates of the Thue–Morse morphism  $\varphi$  on aperiodic binary words. The section is split into three subsections. To state the main result, we define the following.

**Definition 4.1.** Let  $\mathbf{x}$  be a binary word and  $k \geq 1$  an integer. We say that  $\mathbf{x}$  has property  $\mathcal{TM}\mathcal{B}(k)$  if, for all  $1 \leq j \leq k$ , we have  $\mathbf{b}_{\mathbf{x}}^{(j)} = \mathbf{b}_{\mathbf{t}}^{(j)}$ .

Recall that the exact values for  $\mathbf{b}_{\mathbf{t}}^{(j)}$  were computed in [29, Thm. 6] (and are given by (3)). The main result of Section 4.2 is the following theorem, which can be seen as a generalization of the aforementioned result.

**Theorem 4.2.** Let  $k$  be an integer and let  $\mathbf{y}$  be an aperiodic binary word. Then the word  $\mathbf{x} = \varphi^k(\mathbf{y})$  (and any of its suffixes) has property  $\mathcal{TM}\mathcal{B}(k)$ .

The application of  $\varphi^k$  to a word changes the  $j$ -binomial complexity,  $j \leq k$ , to that of the Thue–Morse word. Putting this bluntly, the binomial complexities of the original word play no role in the  $j$ -binomial complexities of the image word (for small  $j$ ).

The topic of Section 4.3 is to characterize the  $k$ - and  $(k + 1)$ -binomial equivalence among factors of words of the form  $\varphi^k(\mathbf{y})$  (Theorem 4.12 and Proposition 4.17). In the latter, we see that structure of  $\mathbf{y}$  already appears to affect the  $(k + 1)$ -binomial complexity of  $\varphi^k(\mathbf{y})$ . This allows to conclude, for example, that  $\mathbf{b}^{(k)} \prec \mathbf{b}^{(k+1)}$  (Corollary 4.18) for words of this form. Throughout the rest of this section we fix  $\mathbf{x}$  and  $\mathbf{y}$  to be as in Theorem 4.2.

We begin with a subsection introducing a convenient tool, called abelian Rauzy graphs, which we use throughout the current and the following section.

### 4.1 Abelian Rauzy graphs

For an infinite word  $\mathbf{z} \in A^{\mathbb{N}}$ , consider a sliding window of length  $n$ : as the window shifts, one goes from heavier factors to lighter factors and vice versa. One can consider a directed labeled graph  $G = (V, E)$  capturing its progress: the vertices are the Parikh vectors of factors of length  $n$ , and there is an edge from  $\vec{x}$  to  $\vec{y}$  labeled with  $(a, b) \in A \times A$ , if there exists  $aub \in \mathcal{L}_{n+1}(\mathbf{z})$  such that  $\Psi(au) = \vec{x}$  and  $\Psi(ub) = \vec{y}$ . We call  $G$  the *abelian Rauzy graph (of order  $n$ )*. Such graphs were considered already in [39].

**Remark 4.3.** The abelian Rauzy graph  $G = (V, E)$  defined here is a quotient of the usual Rauzy graph of  $\mathbf{z}$  of order  $n$ . The latter one is defined as  $R = (V', E')$  where  $V' = \mathcal{L}_n(\mathbf{z})$  and there is an edge from  $au$  to  $ub$  of label  $(a, b)$  whenever  $aub \in \mathcal{L}_{n+1}(\mathbf{z})$ . The Parikh map  $\Psi : V' \rightarrow V$  is a morphism of graphs: any labeled path in  $R$  is mapped to a path in  $G$  with same label.

We describe some properties of abelian Rauzy graphs.

**Observation 4.4.** Let  $G = (V, E)$  be the abelian Rauzy graph of order  $n$  of an infinite word  $\mathbf{z} \in A^{\mathbb{N}}$ .

- The number of vertices is  $\#V = \mathbf{b}_{\mathbf{z}}^{(1)}(n)$ ;
- An edge with label  $(a, b)$  corresponds to an increase in weight if and only if  $ab = 01$ ;
- An edge with label  $(a, b)$  corresponds to a decrease in weight if and only if  $ab = 10$ ;
- An edge is a loop if and only if  $a = b$ .
- A right special factor of length  $n$  gives rise to a *right special vertex*: a vertex with two outgoing edges for which the labels have the same first component. If  $\mathbf{z}$  is binary, then one of the two edges is a loop.

- Each vertex has at least one outgoing edge (possibly a loop). Furthermore, if the word  $z$  is aperiodic, each vertex has at least one outgoing edge that is not a loop, and there is at least one vertex with two outgoing edges. In particular, the graph has at least  $\#V + 1$  edges.

**Example 4.5.** Let us consider the abelian Rauzy graphs of the Thue–Morse word. For a fixed  $n$ , we identify the vertices of  $G$ , that is, the Parikh vectors of factors of length  $n$ , with their second components. Indeed, for a binary word  $u$  of a fixed length  $n$ , we have  $\Psi(u) = (n - |u|_1, |u|_1)^\top$ . We show that, for all  $m \geq 1$ ,  $G_{2m}$  and  $G_{2m+1}$  take forms as depicted in Fig. 1.

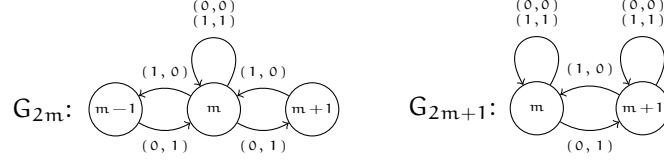


Figure 1: The abelian Rauzy graphs of the Thue–Morse word at even and odd lengths.

Let us write  $\mathbf{t} = 01101001 \cdots = a_0 a_1 a_2 \cdots$ . We first consider  $G_{2m}$ . It is well-known and plain to see that  $\mathbf{t}$  is closed under complementation; for all  $u \in \mathcal{L}(\mathbf{t})$ , we have  $\bar{u} \in \mathcal{L}(\mathbf{t})$ . Let  $au$ ,  $a \in \{0, 1\}$ ,  $au \in \mathcal{L}_m(\mathbf{t})$ , be right special in  $\mathbf{t}$ . Consequently,  $\varphi(au)a \in \mathcal{L}(\mathbf{t})$ ; hence  $G_{2m}$  contains the loop  $m \xrightarrow{(a,a)} m$ . By complementing such a factor, we also find  $m \xrightarrow{(\bar{a}, \bar{a})} m$  in  $G_{2m}$ . To see that, e.g., there is no loop at  $m-1$ , we note that any  $m-1$ -factor is of the form  $0\varphi(u')0$ , and appears at an odd position in  $\mathbf{t}$ . Hence there is certainly no loop with label  $(1, 1)$  at  $m-1$ . Neither can there be a loop with label  $(0, 0)$ , as  $00$  is not the image of a letter.

We then inspect  $G_{2m+1}$ . The following statements are easy to prove using, e.g., the automatic prover Walnut [35]. For a comprehensive take on the usage of Walnut, we recommend the book [47].

- For all  $m \geq 1$  there exists a length- $(2m+2)$  factor  $u$  starting at an even index in  $\mathbf{t}$ , and which begins and ends with 1.  
The following Walnut formula returns 'TRUE':  
`eval OddL11 "Am (m>0) => Ej T[2*j]=@1 & T[2*j+2*m+1]=@1";`
- For all  $m \geq 0$  there exists a length- $(2m+2)$  factor starting at an odd index in  $\mathbf{t}$ , and which begins and ends with 0.  
The following Walnut formula returns 'TRUE':  
`eval OddL00 "Am Ej j>0 & T[2*j-1]=@0 & T[2*j+2*m]=@0";`

Then the first item implies that  $G_{2m+1}$  contains the loop  $m \xrightarrow{(1,1)} m$ . Indeed, the length- $(2m+1)$  prefix of  $u$  in the first item has weight  $m$  (because  $u = \varphi(u')01$  for some  $|u'| = m$ ). Similarly the second item implies that  $m$  has a loop with label  $(0, 0)$ . Recalling that  $\mathbf{t}$  is closed under complementation,  $G_{2m+1}$  is seen to be of the claimed form.

We show the structure of the abelian Rauzy graphs of Sturmian words in Proposition 7.8.

We shall make use of the following general lemma, a part of which appears as [40, Lem. 3.2.]

**Lemma 4.6.** *Let  $z$  be an aperiodic binary word. Then, for every  $n \geq 1$ , the set of edge labels of the abelian Rauzy graph  $G_n$  contains  $(0, 1)$ ,  $(1, 0)$ , and either  $(0, 0)$  or  $(1, 1)$ . Furthermore, if  $(a, a)$  does not appear as a label of a loop in  $G_n$ , then  $G_{n+1}$  contains a loop with the label  $(\bar{a}, \bar{a})$ .*

*Proof.* Since  $z$  is aperiodic, it must have a connected component containing a loop (a right special vertex) and at least two vertices (by a theorem of Coven and Hedlund [15]). An edge from a lighter vertex to a heavier one has label  $(0, 1)$ , and  $(1, 0)$  appears as the label from the heavier to the lighter one.

Assume that  $G_n$  does not contain a loop with label  $(0, 0)$ . Consider the walk  $W$  in  $G_n$  defined by  $z$ : in particular, consider the strongly connected subgraph of  $G_n$  comprising the vertices and

edges that  $W$  traverses infinitely many times. There is an edge from the second lightest vertex to the lightest one 1, labeled with  $(1, 0)$ . This means that the factor of length  $n + 1$  corresponding to this edge begins with 1 and ends with 0. Since  $(0, 0)$  does not appear as a label in  $G_n$  and from the lightest vertex there is no outgoing edge with label  $(1, 0)$ , the edge that  $W$  takes from 1 has label  $(\cdot, 1)$ . Thus the factor of length  $n + 2$  begins and ends with 1. This gives an edge in  $G_{n+1}$  with label  $(1, 1)$ .  $\square$

## 4.2 The First $k$ Binomial Complexities

We begin by defining the notion of  $\varphi^j$ -factorizations of factors of  $\mathbf{x}$ . This will be used throughout this and the next section.

**Definition 4.7.** For any factor  $u$  of  $\varphi^j(\mathbf{y})$  of length at least  $2^j - 1$  there exist  $a, b \in \{0, 1\}$  and  $z \in \{0, 1\}^*$  with  $azb \in \mathcal{L}(\mathbf{y})$  such that  $u = p\varphi^j(z)s$  for some proper suffix  $p$  of  $\varphi^j(a)$  and some proper prefix  $s$  of  $\varphi^j(b)$ . (Note that  $z$  could be empty.) The triple  $(p, \varphi^j(z), s)$  is called a  $\varphi^j$ -factorization<sup>3</sup> of  $u$ . The word  $azb$  (resp.,  $zb$ ;  $az$ ;  $z$ ) is said to be the corresponding  $\varphi^j$ -ancestor of  $u$  when  $p, s$  are non-empty (resp.,  $p = \varepsilon$  and  $s \neq \varepsilon$ ;  $p \neq \varepsilon$  and  $s = \varepsilon$ ;  $p = s = \varepsilon$ ).

Since the words  $\varphi^j(0)$  and  $\varphi^j(1)$  begin with different letters, we notice that if  $s \neq \varepsilon$  in a  $\varphi^j$ -factorization of a word, then the letter  $b$  is uniquely determined. Similarly the  $j$ th images of the letters end with distinct letters (for  $j$  fixed), whence the letter  $a$  is uniquely determined once  $p \neq \varepsilon$ .

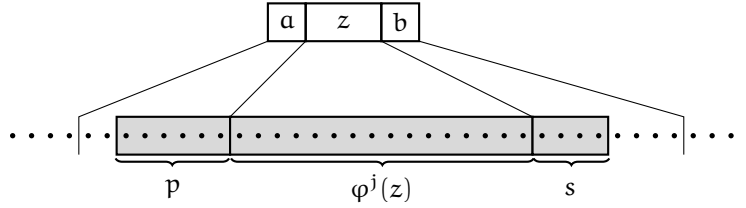


Figure 2: A  $\varphi^j$ -factorization and its  $\varphi^j$ -ancestor.

The following lemma says that an aperiodic word of the form  $\varphi^k(\mathbf{y})$  has the same short factors as the Thue–Morse word.

**Lemma 4.8.** For an aperiodic binary word  $\mathbf{y}$  and integer  $k$ , we have  $\mathcal{L}_n(\varphi^k(\mathbf{y})) = \mathcal{L}_n(\mathbf{t})$  for all  $n \leq 2^k$ .

*Proof.* Let  $\mathbf{x} = \varphi^k(\mathbf{y})$ . The claim is trivial for  $k = 0$  so assume  $k \geq 1$ . Any factor of  $\mathbf{x}$  of length at most  $2^k$  appears as a factor of  $\varphi^k(v)$ , where  $v \in \mathcal{L}_2(\mathbf{y})$ . Since  $\mathcal{L}_2(\mathbf{t}) = \{0, 1\}^2$ , all such factors appear in  $\varphi^k(\mathbf{t}) = \mathbf{t}$ . In particular we have shown  $\mathcal{L}_n(\mathbf{x}) \subseteq \mathcal{L}_n(\mathbf{t})$  for  $n \leq 2^k$ .

Since  $\mathbf{y}$  is aperiodic it contains both 01 and 10 and either of 00 and 11. If  $\mathbf{y}$  contains all factors of length 2, then clearly the considered languages are equal. So assume without loss of generality that 11 does not appear in  $\mathbf{y}$ . Consider the factors of length at most  $2^k$  of  $\varphi^k(11) = \varphi^{k-1}(1)\varphi^k(0)\varphi^{k-1}(0)$ . Such a factor is either a factor of  $\varphi^{k-1}(1)\varphi^k(0)$  or  $\varphi^k(0)\varphi^{k-1}(0)$ , both of which are factors of  $\varphi^k(00)$ . Hence all factors of length at most  $2^k$  appearing in  $\mathbf{t}$  appear in  $\mathbf{x}$  as well; this concludes the proof.  $\square$

We are in the position to prove the main result of this subsection.

*Proof of Theorem 4.2.* Let  $j \in \{1, \dots, k\}$ . We first prove the claim for  $\mathbf{x}$  (and afterwards the claim for any of its suffixes). As the factors of length at most  $2^j$  of  $\mathbf{x} = \varphi^j(\varphi^{k-j}(\mathbf{y}))$  coincide with those

<sup>3</sup>We warn the reader that the term  $\varphi$ -factorization has a different meaning in [29]. Our  $\varphi^j$ -factorization corresponds to their “factorization of order  $j$ ”.

of  $\mathbf{t}$  by the above lemma, the  $j$ -binomial complexity of  $\mathbf{x}$  coincides with that of the Thue–Morse word's for  $n < 2^j$ .

In the remaining of the proof we let  $n \geq 2^j$ . We show that  $\mathcal{L}_n(\mathbf{t})/\sim_j = \mathcal{L}_n(\mathbf{x})/\sim_j$  by double inclusion, which suffices for the claim since Theorem 4.2 holds true for  $\mathbf{x} = \mathbf{t}$ .

Let  $\mathbf{u} \in \mathcal{L}(\mathbf{x})$ ; we show that there exists  $\mathbf{v} \in \mathcal{L}(\mathbf{t})$  such that  $\mathbf{u} \sim_j \mathbf{v}$ . To this end, let  $\mathbf{z} = \varphi^{k-j}(\mathbf{y})$  so that  $\mathbf{x} = \varphi^j(\mathbf{z})$ . Let  $\mathbf{u}$  have  $\varphi^j$ -factorization  $p\varphi^j(\mathbf{u}')$ s with  $\varphi^j$ -ancestor  $\mathbf{au}'\mathbf{b} \in \mathcal{L}(\mathbf{z})$ . The Thue–Morse word contains a factor  $\mathbf{av}'\mathbf{b}$ , where  $|\mathbf{v}'| = |\mathbf{u}'|$  (see, e.g., [29, Prop. 33] or the abelian Rauzy graphs in Example 4.5). It follows that  $\mathbf{t}$  contains the factor  $\mathbf{v} := p\varphi^j(\mathbf{v}')$ s. Now  $\mathbf{u} \sim_j \mathbf{v}$  because  $\varphi^j(\mathbf{u}') \sim_j \varphi^j(\mathbf{v}')$  by Theorem 2.4.

Let then  $\mathbf{u} \in \mathcal{L}(\mathbf{t})$  have  $\varphi^j$ -factorization  $p\varphi^j(\mathbf{u}')$ s with  $\varphi^j$ -ancestor  $\mathbf{au}'\mathbf{b} \in \mathcal{L}(\mathbf{t})$ . As before we show that there exists  $\mathbf{v} \in \mathcal{L}(\mathbf{x})$  such that  $\mathbf{u} \sim_j \mathbf{v}$ . As a consequence of Lemma 4.6,  $\mathbf{z}$  contains, at each length, factors from both the languages  $0A^*1$  and  $1A^*0$ . Hence, if  $\mathbf{a}$  and  $\mathbf{b}$  above are distinct, we may argue as in the previous paragraph to obtain the desired conclusion. Assume thus that  $\mathbf{a} = \mathbf{b}$ . Again Lemma 4.6 implies that  $\mathbf{z}$  contains a factor of length  $|\mathbf{u}'| + 2$  in the language  $1A^*1 \cup 0A^*0$ . Assume without loss of generality that it contains a factor from  $0A^*0$ . Then, if  $\mathbf{a} = \mathbf{b} = 0$ , we may again argue as in the previous paragraph. So assume now that  $\mathbf{a} = \mathbf{b} = 1$  and  $\mathcal{L}_{|\mathbf{u}'|+2}(\mathbf{z}) \cap 1A^*1 = \emptyset$ . Notice that Lemma 4.6 implies that  $\mathcal{L}_{|\mathbf{u}'|+2}(\mathbf{z}) \cap 0A^*0 \neq \emptyset$  and, further,  $\mathcal{L}_{|\mathbf{u}'|+2\pm 1}(\mathbf{z}) \cap 0A^*0 \neq \emptyset$ . To conclude with the claim for  $\mathbf{x}$ , we have four cases to consider depending on the length of  $p$  and  $s$  which can be less or equal, or greater than  $2^{j-1}$ .

**Case 1:** Assume that  $p$  is a suffix of  $\varphi^{j-1}(0)$  and  $s$  is a prefix of  $\varphi^{j-1}(1)$ . For all  $\mathbf{v}'$  such that  $|\mathbf{v}'| = |\mathbf{u}'| - 1$ ,  $\varphi^j(\mathbf{u}') \sim_j \varphi^j(\mathbf{v}'1)$  by Theorem 2.4. By the Transfer Lemma (Lemma 2.5),  $\varphi^j(\mathbf{v}'1) \sim_j \varphi^{j-1}(1)\varphi^j(\mathbf{v}')\varphi^{j-1}(0)$ . Consequently

$$\mathbf{u} \sim_j p\varphi^{j-1}(1)\varphi^j(\mathbf{v}')\varphi^{j-1}(0)s =: \mathbf{v},$$

where  $p\varphi^{j-1}(1)$  is a suffix of  $\varphi^j(0)$  and  $\varphi^{j-1}(0)s$  is a prefix of  $\varphi^j(0)$ . Hence  $\mathbf{v}$  is a factor of  $\varphi^j(0\mathbf{v}'0)$ . Recall that a factor of the form  $0\mathbf{v}'0$  appears in  $\mathbf{z}$  by assumption, and thus  $\varphi^j(0\mathbf{v}'0)$  appears in  $\mathbf{x}$ . To recap, we have shown a factor  $\mathbf{v}$  of  $\mathbf{x}$   $j$ -binomially equivalent to  $\mathbf{u}$ .

**Case 2:** Assume that  $p = p'\varphi^{j-1}(0)$  where  $p'$  is a suffix of  $\varphi^{j-1}(1)$  and  $s$  is a prefix of  $\varphi^{j-1}(1)$ . For all  $\mathbf{v}'$  such that  $|\mathbf{u}'| = |\mathbf{v}'|$ , applying Theorem 2.4 and Lemma 2.5,

$$\mathbf{u} \sim_j p'\varphi^j(\mathbf{v}')\varphi^{j-1}(0)s =: \mathbf{v}.$$

Hence  $\mathbf{v}$  is a factor of  $\varphi^j(0\mathbf{v}'0)$ , and such a factor appears in  $\mathbf{z}$  by assumption. We conclude as above.

**Case 3:** Assume that  $p$  is a suffix of  $\varphi^{j-1}(0)$  and  $s = \varphi^{j-1}(1)s'$  where  $s'$  is a prefix of  $\varphi^{j-1}(0)$ . For all  $\mathbf{v}'$  such that  $|\mathbf{u}'| = |\mathbf{v}'|$ , applying Theorem 2.4 and Lemma 2.5, we have that  $\mathbf{u} \sim_j p\varphi^{j-1}(1)\varphi^j(\mathbf{v}')s' =: \mathbf{v}$  and the conclusion is the same as in the previous case.

**Case 4:** Assume that  $p = p'\varphi^{j-1}(0)$  and  $s = \varphi^{j-1}(1)s'$  where  $p'$  is a suffix of  $\varphi^{j-1}(1)$  and  $s'$  is a prefix of  $\varphi^{j-1}(0)$ . For all  $\mathbf{v}'$  such that  $|\mathbf{v}'| = |\mathbf{u}'| + 1$ , applying Theorem 2.4 and Lemma 2.5,

$$\mathbf{u} \sim_j p'\varphi^{j-1}(0)\varphi^{j-1}(1)\varphi^j(\mathbf{u}')s' \sim_j p'\varphi^j(\mathbf{w}')s' =: \mathbf{v},$$

Hence  $\mathbf{v}$  is a factor of  $\varphi^j(0\mathbf{w}'0)$  and the conclusion is similar to Case 1.

To conclude the proof, we consider the case of a suffix  $\mathbf{w}$  of  $\mathbf{x}$ . Now  $\mathbf{w}$  has a suffix of the form  $\varphi^k(\mathbf{y}')$ , where  $\mathbf{y}'$  is a suffix of  $\mathbf{y}$ . Notice now that  $\mathcal{L}(\mathbf{x}) \supseteq \mathcal{L}(\mathbf{w}) \supseteq \mathcal{L}(\varphi^k(\mathbf{y}'))$ . The theorem applies to both  $\mathbf{x}$  and  $\varphi^k(\mathbf{y}')$  by the previous part. Hence  $\mathbf{w}$  has property  $\mathcal{TM}\mathcal{B}(k)$  also.  $\square$

**Remark 4.9.** If  $\mathbf{y}$  is an aperiodic infinite word, then for all  $\mathbf{a}, \mathbf{b} \in \{0, 1\}$  and  $n \geq 2$  we have  $\mathcal{L}_n(\varphi(\mathbf{y})) \cap \mathbf{aA}^*\mathbf{b} \neq \emptyset$ . Indeed, for  $\mathbf{a} \neq \mathbf{b}$  the claim follows from Lemma 4.6. For  $\mathbf{a} = \mathbf{b}$ , we observe the following: for even length factors  $n = 2\ell$ ,  $\ell \geq 1$ , a factor  $\bar{\mathbf{a}}\mathbf{y}\mathbf{a}$  of  $\mathbf{y}$  of length  $\ell + 1$  (which exists by Lemma 4.6) gives a factor  $\bar{\mathbf{a}}\mathbf{a}\varphi(\mathbf{y})\bar{\mathbf{a}}$  in  $\mathbf{z}$ , hence we have the factor  $\mathbf{aza}$  with  $|\mathbf{z}| = 2\ell - 2$ .



For odd length factors  $n = 2\ell + 1$ ,  $\ell \geq 1$ , we have that a factor of the form  $cyc$ ,  $|y| = \ell - 1$ , of  $\mathbf{y}$  (such a factor exists for some  $c \in \{0, 1\}$  by Lemma 4.6) gives  $c\bar{c}\varphi(y)c\bar{c}$ . Consequently  $\mathbf{z}$  contains a factor in  $aA^*a$  of length  $n$  as well.

Applying this observation in the above proof to  $\mathbf{z}$  when  $j < k$ , we have  $\mathcal{L}_n(\mathbf{z}) \cap 1A^*1 \neq \emptyset$  for all  $n \geq 2$ , and thus case analysis at the end of the proof is only necessary for the case  $j = k$ .

### 4.3 On the $k$ - and $(k + 1)$ -Binomial Equivalence

The previous subsection was dealing with the  $j$ -binomial equivalence in  $\mathbf{x} = \varphi^k(\mathbf{y})$ , where  $\mathbf{y}$  is an aperiodic binary word and  $j \leq k$ . Here, we are concerned with the  $(k + 1)$ -binomial equivalence in such words. To this end, we need to have more control on the  $k$ -binomial equivalence in  $\mathbf{x}$ . First, we have a closer look at the  $\varphi^j$ -factorizations of a word and in particular at the associated prefixes and suffixes.

**Definition 4.10** ([29, Def. 43]). Let  $j \geq 1$ . Let us define the equivalence relation  $\equiv_j$  on  $A^{<2^j} \times A^{<2^j}$  by  $(p_1, s_1) \equiv_j (p_2, s_2)$  whenever there exists  $a \in A$  such that one of the following situations occurs:

1.  $|p_1| + |s_1| = |p_2| + |s_2|$  and
  - (a)  $(p_1, s_1) = (p_2, s_2)$ ;
  - (b)  $(p_1, \varphi^{j-1}(a)s_1) = (p_2\varphi^{j-1}(a), s_2)$ ;
  - (c)  $(p_2, \varphi^{j-1}(a)s_2) = (p_1\varphi^{j-1}(a), s_1)$ ;
  - (d)  $(p_1, s_1) = (s_2, p_2) = (\varphi^{j-1}(a), \varphi^{j-1}(\bar{a}))$ ;
2.  $||p_1| + |s_1| - (|p_2| + |s_2|)| = 2^j$  and
  - (a)  $(p_1, s_1) = (p_2\varphi^{j-1}(a), \varphi^{j-1}(\bar{a})s_2)$ ;
  - (b)  $(p_2, s_2) = (p_1\varphi^{j-1}(a), \varphi^{j-1}(\bar{a})s_1)$ .

The next lemma is essentially [29, Lem. 40 and 41] (except that with an arbitrary word  $\mathbf{y}$  instead of the Thue–Morse word  $\mathbf{t}$ , we cannot use the fact that  $\mathbf{t}$  is overlap-free, so factors such as 10101 may appear in  $\mathbf{y}$ ). To each  $\varphi^j$ -factorization there is a natural corresponding  $\varphi^{j-1}$ -factorization, though two  $\varphi^j$ -factorizations may correspond to the same  $\varphi^{j-1}$ -factorization. The next lemma also describes how such factorizations are related.

**Lemma 4.11.** *Let  $j \geq 1$ . Let  $u$  be a factor of  $\varphi^j(\mathbf{y})$  such that  $|u| \geq 2^j - 1$ . Then  $u$  has at most two  $\varphi^j$ -factorizations. Let further  $u$  have a  $\varphi^j$ -factorization of the form  $(p, \varphi^j(z), s)$  and  $z_0zz_{n+1}$  being the corresponding  $\varphi^j$ -ancestor (where according to Definition 4.7  $z_0, z_{n+1}$  or  $z$  could be empty). The factor  $u$  has a unique  $\varphi^j$ -factorization if and only if the word  $z_0zz_{n+1}$  contains both letters 0 and 1. Moreover, if there is another  $\varphi^j$ -factorization  $(p', \varphi^j(z'), s')$  with  $\varphi^j$ -ancestor  $z'_0z'z'_{m+1}$ , then  $(p, s) \equiv_j (p', s')$  with  $||p| - |p'|| = ||s| - |s' || = 2^{j-1}$ ,  $z_0zz_{n+1} = a^{n+2}$ , and  $z'_0z'z'_{m+1} = \bar{a}^{m+2}$  for some  $a \in \{0, 1\}$ .*

Otherwise stated, the  $\varphi^j$ -factorization is not unique if and only if  $u$  is a factor of  $\varphi^{j-1}(x)$  with  $x \in (01)^* \cup (10)^* \cup 1(01)^* \cup 0(10)^*$ .

*Proof.* Since  $|u| \geq 2^j - 1$ ,  $u$  has a factor of the form  $\varphi^{j-1}(a)$  and thus at least one  $\varphi^j$ -factorization of the prescribed form exists with  $z = z_1 \cdots z_n$  and  $n \geq 0$  ( $n = 0$  if  $z = \varepsilon$ ).

We first prove the claim for uniqueness by induction on  $j$ . For  $j = 1$ , assume that  $u = z_0\varphi(z_1) \cdots \varphi(z_n)z_{n+1}$  with  $z_0, z_{n+1} \in \{0, 1, \varepsilon\}$ . Suppose, as in the statement, that both letters 0 and 1 occur in  $z_0 \cdots z_{n+1}$ . Then we have  $z_i z_{i+1} = 01$  (or similarly 10) for some  $i$ . This means that  $u$  contains the factor 11 forcing uniqueness of this kind of a factorization:  $11 \notin \{\varphi(0), \varphi(1)\}$ . Assume that the property holds true up to  $j - 1$  and prove it for  $j \geq 2$ . Let  $u = p\varphi^j(z_1) \cdots \varphi^j(z_n)s$  be a  $\varphi^j$ -factorization and assume that  $z_i z_{i+1} = 01$  for some  $i$ . To this factorization, we have a corresponding factorization of the form

$$u = p\varphi^{j-1}(z_1)\varphi^{j-1}(\bar{z}_1) \cdots \varphi^{j-1}(z_n)\varphi^{j-1}(\bar{z}_n)s.$$

Notice that  $p$  is a suffix of  $\varphi^{j-1}(\bar{z}_0)$  if  $|p| < 2^{j-1}$  and otherwise,  $p = p' \varphi^{j-1}(\bar{z}_0)$  with  $p'$  a suffix of  $\varphi^{j-1}(z_0)$ . Similarly,  $s$  is a prefix of  $\varphi^{j-1}(z_{n+1})$  if  $|s| < 2^{j-1}$  and otherwise,  $s = \varphi^{j-1}(z_{n+1})s'$  with  $s'$  a prefix of  $\varphi^{j-1}(\bar{z}_{n+1})$ . Observe that  $z_i \bar{z}_i z_{i+1} \bar{z}_{i+1} = 0110$ . So by the induction hypothesis, the  $\varphi^{j-1}$ -factorization of  $u$  is unique. There are at most two  $\varphi^j$ -factorizations corresponding to a  $\varphi^{j-1}$ -factorization. But since  $\varphi^{j-1}(1)\varphi^{j-1}(1) \notin \{\varphi^j(0), \varphi^j(1)\}$ , the claimed uniqueness follows.

We then prove the claim for non-unique factorizations. Assume that  $z_0 = z_1 = \dots = z_{n+1} = 0$ . Then

$$u = p\varphi^j(0) \dots \varphi^j(0)s = p\varphi^{j-1}(0)\varphi^{j-1}(1) \dots \varphi^{j-1}(0)\varphi^{j-1}(1)s$$

with  $p$  (resp.,  $s$ ) a suffix (resp., prefix) of  $\varphi^j(0)$ . If  $|p| \geq 2^{j-1}$ , then  $p = p' \varphi^{j-1}(1)$  with  $p'$  a suffix of  $\varphi^{j-1}(0)$  (and thus, a suffix of  $\varphi^j(1)$ ), otherwise set  $p' = p\varphi^{j-1}(0)$ . Similarly, if  $|s| \geq 2^{j-1}$ , then  $s = \varphi^{j-1}(0)s'$  with  $s'$  a prefix of  $\varphi^{j-1}(1)$ , otherwise  $s' = \varphi^{j-1}(1)s$ . Notice that the corresponding  $\varphi^{j-1}$ -factorization of  $u$  is unique since the  $\varphi^{j-1}$ -ancestor is not a power of a letter: if  $n \neq 0$  then the claim is clear. Otherwise  $|u| = |ps| \geq 2^j - 1$ ; this implies that either  $|p| \geq 2^{j-1}$  or  $|s| \geq 2^{j-1}$ . Assuming the latter (the other case being symmetric), we have that  $s = \varphi^{j-1}(0)s'$  with  $s'$  a prefix of  $\varphi^{j-1}(1)$ . If  $s' \neq \varepsilon$ , then the  $\varphi^{j-1}$ -ancestor contains both letters. If  $s' = \varepsilon$ , then  $p \neq \varepsilon$ , and then again the  $\varphi^{j-1}$ -factorization contains both letters.

Now  $u$  can also be written as

$$p' \varphi^{j-1}(1)\varphi^{j-1}(0) \dots \varphi^{j-1}(1)\varphi^{j-1}(0)s' = p' \varphi^j(1) \dots \varphi^j(1)s'.$$

There are no other  $\varphi^j$ -factorizations due to the uniqueness of the  $\varphi^{j-1}$  factorization of  $u$ . To conclude the claim in this case, a straightforward case analysis shows that  $(p, s) \equiv_j (p', s')$  with  $||p| - |p'|| = ||s| - |s'|| = 2^{j-1}$ :

If  $|p| \geq 2^{j-1}$  and if  $|s| \geq 2^{j-1}$ , then  $(p, s) = (p' \varphi^{j-1}(1), \varphi^{j-1}(0)s')$ .

If  $|p| \geq 2^{j-1}$  and if  $|s| < 2^{j-1}$ , then  $(p, \varphi^{j-1}(1)s) = (p' \varphi^{j-1}(1), s')$ .

If  $|p| < 2^{j-1}$  and if  $|s| \geq 2^{j-1}$ , then  $(p\varphi^{j-1}(0), s) = (p', \varphi^{j-1}(0)s')$ .

If  $|p| < 2^{j-1}$  and if  $|s| < 2^{j-1}$ , then  $(p\varphi^{j-1}(0), \varphi^{j-1}(1)s) = (p', s')$ .  $\square$

We have the following theorem, the proof of which is essentially the proof of [29, Thm. 48]. Indeed, the lemmas in [29] leading to its proof do not require that the factors  $u$  and  $v$  are from the Thue–Morse word, only that they have  $\varphi^j$ -factorizations. We note that [29, Thm. 48] is stated for  $j \geq 3$ . However, the statement holds also for  $j = 1$  (trivially) and for  $j = 2$  as it is essentially a restatement of [29, Thm. 34] obtained by closely inspecting its proof.

**Theorem 4.12.** *Let  $\mathbf{y}$  be an aperiodic binary word. Let  $k \geq j \geq 1$ . Let  $u$  and  $v$  be equal-length factors of  $\mathbf{x} = \varphi^k(\mathbf{y})$  with  $\varphi^j$ -factorizations  $u = p_1 \varphi^j(z) s_1$  and  $v = p_2 \varphi^j(z') s_2$ . Then  $u \sim_j v$  if and only if  $(p_1, s_1) \equiv_j (p_2, s_2)$ .*

We then turn to the  $(k+1)$ -binomial equivalence in  $\mathbf{x}$ . A straightforward consequence of (4) together with the identities  $\sum_{x \in A^\ell} \binom{u}{x} = \binom{|u|}{\ell}$ ,  $\ell \geq 1$ , is the following observation.

**Lemma 4.13.** *Let  $u \in \{0, 1\}^*$ . Then*

$$\binom{\varphi(u)}{0} = |u|; \quad \binom{\varphi(u)}{01} = |u|_0 + \binom{|u|}{2}; \quad \binom{\varphi(u)}{011} = \binom{u}{01} + \binom{|u|_0}{2} + \binom{|u|}{3}.$$

*Proof.* Observe that  $\binom{\varphi(a)}{011} = 0 = \binom{\varphi(a)}{111}$  for both  $a \in \{0, 1\}$ . Similarly  $\binom{\varphi(a)}{b} = 1$  for letters  $a, b \in \{0, 1\}$ . Therefore

$$\begin{aligned} \binom{\varphi(u)}{011} &= \sum_{x_1, x_2 \in A} \binom{u}{x_1 x_2} \sum_{\substack{011 = e_1 e_2 \\ e_i \in A^+}} \binom{\varphi(x_1)}{e_1} \binom{\varphi(x_2)}{e_2} + \sum_{|x|=3} \binom{u}{x} \\ &= \binom{u}{00} + \binom{u}{01} + \binom{|u|}{3}. \end{aligned}$$

and the claim follows.  $\square$

The next technical lemma has an important role in studying the  $(k+1)$ -binomial equivalence.

**Lemma 4.14.** *Let  $u, v$  be two binary words of equal length. For  $k \geq 1$ , we have*

$$\binom{\varphi^k(u)}{01^k} - \binom{\varphi^k(v)}{01^k} = 2^{(k-1)(k-2)/2} (|u|_0 - |v|_0).$$

*In particular,  $u \not\sim_1 v$  implies  $\varphi^k(u) \not\sim_{k+1} \varphi^k(v)$ . Moreover, if  $u \sim_1 v$ , for  $k \geq 1$ , we have*

$$\binom{\varphi^k(u)}{01^{k+1}} - \binom{\varphi^k(v)}{01^{k+1}} = 2^{(k-1)(k-2)/2} \left( \binom{u}{01} - \binom{v}{01} \right).$$

*In particular,  $u \not\sim_2 v$  implies  $\varphi^k(u) \not\sim_{k+2} \varphi^k(v)$ .*

*Proof.* The case  $k = 1$  is deduced from Lemma 4.13. Then assume  $k \geq 2$ . We encourage the reader to refer to [29] for details that would be too long to reproduce here. From [29, Rem. 23], we have the following expression

$$\binom{\varphi^k(u)}{01^k} - \binom{\varphi^k(v)}{01^k} = \sum_{x \in f^k(01^k)} m_{f^k(01^k)}(x) \left[ \binom{u}{x} - \binom{v}{x} \right],$$

where the map  $f$  is defined to take into account the multiple ways factors  $01$  or  $10$  may occur in a word:  $f(u)$  is a multiset of words of length shorter than  $u$ ; see [29, Def. 15 and 17]. We let the coefficient  $m_{f^k(01^k)}(x)$  denote the multiplicity of  $x$  as an element of the multiset  $f^k(01^k)$ . It can be shown that the multiset  $f^k(01^k)$  only contains the elements  $0$  and  $1$ . Therefore we obtain

$$\binom{\varphi^k(u)}{01^k} - \binom{\varphi^k(v)}{01^k} = m_{f^k(01^k)}(0) (|u|_0 - |v|_0) + m_{f^k(01^k)}(1) (|u|_1 - |v|_1).$$

To conclude the proof, we use two facts. The first is that  $|u|_1 - |v|_1 = -(|u|_0 - |v|_0)$  since  $u, v$  have equal length. The second is that

$$m_{f^k(01^k)}(0) - m_{f^k(01^k)}(1) = m_{f^{k-1}(01^k)}(01) - m_{f^{k-1}(01^k)}(10) = 2^{(k-1)(k-2)/2},$$

which follows from [29, Prop. 28]. For the second part, the same reasoning may be applied to obtain

$$\binom{\varphi^k(u)}{01^{k+1}} - \binom{\varphi^k(v)}{01^{k+1}} = \sum_{x \in f^k(01^{k+1})} m_{f^k(01^{k+1})}(x) \left[ \binom{u}{x} - \binom{v}{x} \right].$$

The multiset  $f^k(01^{k+1})$  only contains  $0, 1, 00, 01, 10, 11$ . But since it is assumed that  $u \sim_1 v$ , the only (potentially) non-zero terms in the sum correspond to  $x \in \{01, 10\}$ . Then the observation  $\binom{u}{01} - \binom{v}{01} = \binom{v}{10} - \binom{u}{10}$  following from Lemma 2.7 suffices to conclude.  $\square$

Next we consider the structure of factors of the image of an arbitrary binary word  $y$ .

**Definition 4.15.** For  $n \geq 1$  we let  $\mathcal{S}(n) = \mathcal{L}_n(y)$ . Further, for all  $a, b \in \{\varepsilon, 0, 1\}$  such that  $ab \neq \varepsilon$ , we define  $\mathcal{S}_{a,b}(n) = \mathcal{L}_{n+|ab|}(y) \cap aA^*b$ . We call these sets *factorization classes of order  $n$* .

Consider now a factor  $u$  of  $\varphi(y)$ . We associate with  $u$  some factorization classes as follows. Let  $a\varphi(u')b$  be the  $\varphi$ -factorization of  $u$  with  $\varphi$ -ancestor  $au'b \in \mathcal{L}(y)$ . If  $ab = \varepsilon$ , we associate the factorization class  $\mathcal{S}(|u'|)$ . For  $ab \neq \varepsilon$ , we have that  $u$  is a factor of  $\varphi(\bar{a}u'b)$ . In this case we associate the factorization class  $\mathcal{S}_{\bar{a},b}(|u'|)$ . If  $u$  is associated with a factorization class  $\mathcal{T}$ , we write  $u \models \mathcal{T}$ , otherwise we write  $u \not\models \mathcal{T}$ .

Observe that  $u \models \mathcal{S}(n)$  implies that  $|u| = 2n$ . Also, for  $ab \neq \varepsilon$ ,  $u \models \mathcal{S}_{a,b}(n)$  implies that  $|u| = 2n + |ab|$ . Notice also that a factor  $u$  of  $\varphi(y)$  can be associated with several factorization classes: take, e.g.,  $(10)^\ell 1 = 1(01)^\ell$  which is associated with both  $\mathcal{S}_{\varepsilon,1}(\ell)$  and  $\mathcal{S}_{0,\varepsilon}(\ell)$ , or  $(01)^{\ell+1} = 0(10)^\ell 1$  which is associated with both  $\mathcal{S}(\ell+1)$  and  $\mathcal{S}_{1,1}(\ell)$ .

**Lemma 4.16.** *For two 2-binomially equivalent factors  $u, v \in \mathcal{L}(\varphi(\mathbf{y}))$ , if  $u \models \mathcal{T}$  for some factorization class  $\mathcal{T}$ , then  $v \models \mathcal{T}$ . Furthermore, a factor  $u$  of  $\mathbf{y}$  is associated with distinct factorization classes if and only if  $u \in L = (01)^* \cup (10)^* \cup 1(01)^* \cup 0(10)^*$ .*

*Proof. Even-length factors.* Let  $u \sim_2 v$  with  $|u| = 2n$ . If  $u \models \mathcal{S}_{\bar{a}, a}(n-1)$  with  $a \in \{0, 1\}$ , then  $u$  is of the form  $a\varphi(x)a$  with  $|x| = n-1$ , whence  $|u|_a = n+1$ . Factors  $v' \not\models \mathcal{S}_{\bar{a}, a}(n-1)$  of length  $2n$  have  $|v'|_a \leq n$  by inspection. Hence also  $v \models \mathcal{S}_{\bar{a}, a}(n-1)$ . The above arguments also show that  $u$  is associated with exactly one factorization class. For the latter claim, we note that  $u$  has even length and begins and ends with the same letter, so it cannot appear in the language  $L$ .

Assume then that  $u \not\models \mathcal{S}_{\bar{a}, a}(n-1)$ ,  $a \in \{0, 1\}$ . Then  $v \not\models \mathcal{S}_{\bar{a}, a}(n-1)$ ,  $a \in \{0, 1\}$  by the previous observation. Notice that we may assume  $n \geq 2$  as otherwise we have  $|u| = 2$  and the claim is trivial (2-binomial equivalence is equality in this case). We compare the values of  $\binom{y}{01}$  for  $y$  associated with  $\mathcal{S}_{1,1}(n-1)$ ,  $\mathcal{S}_{0,0}(n-1)$ , and  $\mathcal{S}(n)$ , respectively.

**Case 1:**  $y \models \mathcal{S}_{1,1}(n-1)$ . We have  $\binom{y}{01} \geq \binom{n}{2} + n$ , and equality holds for  $y = (01)^n$ . Indeed, say  $y = 0\varphi(x)1$  for some  $x \in \{0, 1\}^{n-1}$ . Then we have by Lemma 4.13

$$\binom{y}{01} = \binom{\varphi(x)}{01} + |\varphi(x)|_0 + |\varphi(x)|_1 = |x|_0 + \binom{|x|}{2} + 2|x| + 1 = |x|_0 + \binom{n}{2} + n,$$

since  $|x| = n-1$ . Equality now holds when  $|x|_0 = 0$ , i.e.,  $x = 1^{n-1}$ .

**Case 2:**  $y \models \mathcal{S}_{0,0}(n-1)$ . We have  $\binom{y}{01} \leq \binom{n}{2}$ , and equality holds when  $y = (10)^n$ . Indeed, say  $y = 1\varphi(x)0$  for some  $x \in \{0, 1\}^{n-1}$ . Then

$$\binom{y}{01} = \binom{\varphi(x)}{01} = |x|_0 + \binom{|x|}{2} = |x|_0 + \binom{n}{2} - (n-1).$$

Since  $|x| = n-1$ , we have  $\binom{y}{01} \leq \binom{n}{2}$ . Equality holds when  $x = 0^{n-1}$ .

**Case 3:**  $y \models \mathcal{S}(n)$ . We have  $\binom{n}{2} \leq \binom{y}{01} \leq \binom{n}{2} + n$ . The former equality is attained with  $y = (10)^n$  and the latter with  $y = (01)^n$ . Indeed, say  $y = \varphi(x')$  for some  $x' \in \{0, 1\}^n$ . We have  $\binom{y}{01} = \binom{n}{2} + |x'|_0$  from Lemma 4.13. Therefore,  $\binom{n}{2} \leq \binom{y}{01} \leq \binom{n}{2} + n$ . The former equality is attained with  $x' = 1^n$  and the latter with  $x' = 0^n$ .

We conclude that  $u$  and  $v$  are associated with a common factorization class. In fact, the latter claim is also implied from the above: a word can be associated with two (and only two) factorization classes if and only if it appears in  $L$ . This concludes the proof in the case of even length factors.

**Odd-length factors.** Assume without loss of generality that  $u \models \mathcal{S}_{a, \varepsilon}(n)$  with  $u = a\varphi(u')$  of length  $2n+1$ . Recalling that  $|\varphi(u')|_0 = |u'| = n$ , if  $u \sim_2 v$  with  $u$  and  $v$  associated with distinct factorization classes, then necessarily  $v \in \mathcal{S}_{\varepsilon, a}$ , say  $v = \varphi(v')a$ . We show that this is impossible, unless  $u = v \in L$ .

Indeed, assuming that we have 2-binomial equivalence, we have

$$\binom{a\varphi(u')}{01} = \binom{\varphi(u')}{01} + \delta_0(a) \binom{\varphi(u')}{1} = |u'|_0 + \binom{n}{2} + \delta_0(a)n \quad (6)$$

which is equal to

$$\binom{\varphi(v')a}{01} = \binom{\varphi(v')}{01} + \delta_1(a) \binom{\varphi(v')}{0} = |v'|_0 + \binom{n}{2} + \delta_1(a)n \quad (7)$$

where  $\delta_a(b) = 1$  if  $a = b$ , otherwise  $\delta_a(b) = 0$ . Rearranging, we get  $|u'|_0 - |v'|_0 = (\delta_1(a) - \delta_0(a))n \in \{\pm n\}$ . This implies, without loss of generality, that  $u' = 0^n$ ,  $v' = 1^n$ , and  $a = 1$ . But then  $u = 1(01)^n = (10)^n 1 = v \in L$ , as claimed.  $\square$

The next result characterizes  $(k+1)$ -binomial equivalence in  $x = \varphi^k(\mathbf{y})$  when  $\mathbf{y}$  is an arbitrary binary word.

**Proposition 4.17.** *Let  $u$  and  $v$  be factors of length at least  $2^k - 1$  of  $\mathbf{x} = \varphi^k(\mathbf{y})$  with the  $\varphi^k$ -factorizations  $u = p_1 \varphi^k(z) s_1$  and  $v = p_2 \varphi^k(z') s_2$ . Then  $u \sim_{k+1} v$  and  $u \neq v$  if and only if  $z \sim_1 z'$ ,  $z' \neq z$ , and  $(p_1, s_1) = (p_2, s_2)$ .*

Notice that the proposition claims that those factors of  $\mathbf{x}$  having at least two  $\varphi^k$ -factorizations are  $(k+1)$ -binomially equivalent only to themselves (in  $\mathcal{L}(\mathbf{x})$ ).

*Proof.* The “if”-part of the statement follows by a repeated application of Proposition 3.9 on the Thue–Morse morphism together with the fact that the morphism is injective.

Let us assume that  $u \sim_{k+1} v$  for some distinct factors. It follows that  $u \sim_k v$ , which implies that  $(p_1, s_1) \equiv_k (p_2, s_2)$  by Theorem 4.12. Next we show that  $(p_1, s_1) = (p_2, s_2)$  and  $z \sim_1 z'$ . We have the following case distinction from Definition 4.10:

(1)(a): We have that  $(p_1, s_1) = (p_2, s_2)$ . By deleting the common prefix  $p_1$  and suffix  $s_1$ , we are left with the equivalent statement  $\varphi^k(z) \sim_{k+1} \varphi^k(z')$ . If  $z \not\sim_1 z'$ , then we have a contradiction with Lemma 4.14. The desired result follows in this case.

In the remaining cases, we assume towards a contradiction that  $(p_1, s_1) \neq (p_2, s_2)$ .

(1)(b): Suppose that  $(p_1, s_1) = (p_2 \varphi^{k-1}(a), \varphi^{k-1}(a) s_1)$ . Deleting the common prefixes  $p_2$  and suffixes  $s_1$ , we are left with  $\varphi^{k-1}(a \varphi(z)) \sim_{k+1} \varphi^{k-1}(a \varphi(z'))$ . Now  $a \varphi(z) \sim_1 \varphi(z') a$ , but  $a \varphi(z) \not\sim_2 \varphi(z') a$  by Lemma 4.16 (otherwise  $a \varphi(z) = \varphi(z') a$  and thus  $u = v$  contrary to the assumption). Lemma 4.14 then implies that  $\varphi^{k-1}(a \varphi(z)) \not\sim_{k+1} \varphi^{k-1}(a \varphi(z'))$ , which is a contradiction.

(1)(c): Suppose that  $(p_2, \varphi^{k-1}(a) s_2) = (p_1 \varphi^{k-1}(a), s_1)$ . This is symmetric to the previous case.

(1)(d): Suppose that  $(p_1, s_1) = (s_2, p_2) = (\varphi^{k-1}(a), \varphi^{k-1}(\bar{a}))$ . We thus have directly

$$\varphi^{k-1}(a \varphi(z) \bar{a}) \sim_{k+1} \varphi^{k-1}(\bar{a} \varphi(z) a).$$

The claim follows by an argument similar to that of in Case (1)(b).

(2)(a): Suppose that  $(p_1, s_1) = (p_2 \varphi^{k-1}(a), \varphi^{k-1}(\bar{a}) s_2)$ . After removing common prefixes and suffixes, we are left with  $\varphi^{k-1}(a \varphi(z) \bar{a}) \sim_{k+1} \varphi^{k-1}(\varphi(z'))$ . We have that  $a \varphi(z) \bar{a} \sim_1 \varphi(z')$ , but by Lemma 4.16  $a \varphi(z) \bar{a} \not\sim_2 \varphi(z')$  (otherwise  $z = \bar{a}^\ell$  and  $z' = a^{\ell+1}$ , implying that  $u = v$ , a contradiction). This is again a contradiction by Lemma 4.14.

(2)(b): Suppose that  $(p_2, s_2) = (p_1 \varphi^{j-1}(a), \varphi^{j-1}(\bar{a}) s_1)$ . This is symmetric to the previous case.  $\square$

Notice that Theorem 4.2 and Proposition 4.17 have the following corollary:

**Corollary 4.18.** *Let  $\mathbf{x} = \varphi^k(\mathbf{y})$ , where  $\mathbf{y}$  is an arbitrary aperiodic binary word. We have*

$$\mathbf{b}_x^{(1)} \prec \mathbf{b}_x^{(2)} \prec \dots \prec \mathbf{b}_x^{(k)} \prec \mathbf{b}_x^{(k+1)}.$$

*Proof.* Recall that  $\mathbf{y}$  contains arbitrarily long factors of the form  $\bar{a} z a$ ,  $a \in \{0, 1\}$ , by Lemma 4.6. Therefore  $\mathbf{x}$  contains the  $k$ -binomially equivalent (by Lemma 2.5) factors  $\varphi^{k-1}(a) \varphi^k(z)$  and  $\varphi^k(z) \varphi^{k-1}(a)$ . However, by Proposition 4.17 these factors are either not  $(k+1)$ -binomially equivalent, or  $\varphi^{k-1}(a) \varphi^k(z) = \varphi^k(z) \varphi^{k-1}(a)$ . The latter happens when  $\varphi^k(z) = \varphi^{k-1}(a)^\ell$  for some  $\ell \geq 0$ , and thus only when  $\ell = 0$  and  $z = \varepsilon$ . (Indeed, it is not hard to prove that if  $w$  is primitive so is  $\varphi(w)$ .) This observation suffices for showing  $\mathbf{b}_x^{(k)} \prec \mathbf{b}_x^{(k+1)}$ . The rest of the claim follows by Theorem 4.2.  $\square$

## 5 Binomial Properties of the Thue–Morse Morphism, Part II

In this section we consider a complementary result to Theorem 4.2, which partially extends the following theorem of Richomme, Saari, and Zamponi [40].

**Theorem 5.1** ([40, Thm 3.3]). *Let  $\mathbf{x}$  be an aperiodic word. Then  $\mathbf{b}_x^{(1)} = \mathbf{b}_t^{(1)}$  if and only if there exists a binary word  $\mathbf{y}$  and  $a \in \{\varepsilon, 0, 1\}$  such that  $\mathbf{x} = a \varphi(\mathbf{y})$ .*

Notice that Theorem 4.2 is a generalization of the “if”-direction. The following Theorem 5.2 is a partial generalization in the other direction. It is proved in Section 5.2. Recall Definition 4.1: For an integer  $k \geq 1$ , a binary word  $x$  has property  $\mathcal{TM}\mathcal{B}(k)$  if, for all  $1 \leq j \leq k$ , we have  $b_x^{(j)} = b_t^{(j)}$ .

**Theorem 5.2.** *Let  $x$  be a recurrent binary word having property  $\mathcal{TM}\mathcal{B}(k)$  for some  $k \geq 1$ . Then there exists a binary word  $y$  such that  $x = u\varphi^k(y)$ , where  $u$  is a proper suffix of  $\varphi^k(0)$  or  $\varphi^k(1)$ .*

To prove the theorem, we first derive a formula for counting  $(k + 1)$ -binomial equivalence classes of words that are of the form  $\varphi^k(y)$  for  $y$  aperiodic in Section 5.1.

## 5.1 A formula for counting $(k + 1)$ -binomial complexities

For a binary word  $y$  we define

$$\begin{aligned} X_y(n) &:= \{(a, \Psi(u), b) : a, b \in \{0, 1\}, aub \in \mathcal{L}_{n+1}(y)\}, \\ Y_{y,L}(n) &:= \{(a, \Psi(u)) : a \in \{0, 1\}, au \in \mathcal{L}_{n+1}(y)\}, \\ Y_{y,R}(n) &:= \{(\Psi(u), a) : a \in \{0, 1\}, ua \in \mathcal{L}_{n+1}(y)\}, \\ Y_y(n) &:= Y_{y,L} \cup Y_{y,R}. \end{aligned}$$

**Observation 5.3.** Let  $G_n = (V, E)$  be the abelian Rauzy graph of  $y$  (of order  $n$ ).

- $E$  is in one-to-one correspondence with  $X_y(n)$ , namely,  $\vec{x} \xrightarrow{(a,b)} \vec{y}$  is identified with  $(a, \vec{x} - \Psi(a), b)$ . In particular,  $\#E = \#X_y(n)$ .
- The set  $Y_{y,R}(n)$  is in one-to-one correspondence with  $E/\equiv_R$ , where  $\equiv_R$  is the equivalence relation defined by the (surjective) mapping  $E \rightarrow Y_{y,R}(n)$  (meaning, the equivalence classes are the full preimages of elements of  $Y_{y,R}(n)$ ),

$$\left(\vec{x} \xrightarrow{(a,b)} \vec{y}\right) \mapsto (\vec{x}, b).$$

In particular,  $\#Y_{y,R}(n) \leq \#E$ . Similarly  $Y_{y,L}(n)$  is in one-to-one correspondence with  $E/\equiv_L$ , where  $\equiv_L$  is the equivalence relation defined by the (surjective) mapping

$$\left(\vec{x} \xrightarrow{(a,b)} \vec{y}\right) \mapsto (a, \vec{y}).$$

In particular,  $\#Y_{y,L}(n) \leq \#E$ .

- Each equivalence class in  $E/\equiv_L$  contains at most two elements: two edges are equivalent if their target vertices and the first components of the labels are equal. Hence the equivalence relation can only identify a non-loop edge with a loop.
- Note that any loop at a vertex  $v$  with label  $(0, 0)$  can only be equivalent under either  $\equiv_L$  or  $\equiv_R$  to an edge between  $v$  and a lighter vertex. Similarly a loop with label  $(1, 1)$  can only be equivalent to an edge between  $v$  and a heavier vertex. In particular, a loop with  $(0, 0)$  (resp.,  $(1, 1)$ ) on the lightest (resp., heaviest) vertex is not equivalent to any other edge under either  $\equiv_L$  or  $\equiv_R$ .

**Example 5.4.** Recall the abelian Rauzy graphs of the Thue–Morse word from Example 4.5. The edges correspond exactly to  $X_t(n)$ . The equivalence classes of  $E/\equiv_R$  (resp.,  $E/\equiv_L$ ) corresponding to  $Y_{t,R}(n)$  (resp.,  $Y_{t,L}(n)$ ) containing at least two elements are listed below:

$$\begin{aligned} n = 2m: Y_{t,R} &: \left\{ m \xrightarrow{(0,0)} m, m \xrightarrow{(1,0)} m-1 \right\}, \quad \left\{ m \xrightarrow{(1,1)} m, m \xrightarrow{(0,1)} m+1 \right\}; \\ Y_{t,L} &: \left\{ m \xrightarrow{(0,0)} m, m-1 \xrightarrow{(0,1)} m \right\}, \quad \left\{ m \xrightarrow{(1,1)} m, m+1 \xrightarrow{(1,0)} m \right\}. \end{aligned}$$

$$n = 2m - 1: Y_{t,R} : \left\{ m \xrightarrow{(1,1)} m, m \xrightarrow{(0,1)} m+1 \right\}, \quad \left\{ m+1 \xrightarrow{(0,0)} m+1, m+1 \xrightarrow{(1,0)} m \right\};$$

$$Y_{t,L} : \left\{ m \xrightarrow{(1,1)} m, m+1 \xrightarrow{(1,0)} m \right\}, \quad \left\{ m+1 \xrightarrow{(0,0)} m+1, m \xrightarrow{(0,1)} m+1 \right\}.$$

We may now establish a formula for counting the  $(k+1)$ -binomial complexity of the  $k$ th image of a word  $\mathbf{y}$  under the Thue–Morse morphism. This will turn out to be key in proving a converse to Theorem 4.2.

**Proposition 5.5.** *Let  $\mathbf{y}$  be an infinite binary word and let  $\mathbf{x} = \varphi^k(\mathbf{y})$  with  $k \geq 1$ . Let  $m = \max\{n \in \mathbb{N} : 0^n \text{ and } 1^n \in \mathcal{L}(\mathbf{y})\}$  and  $m' = \max\{n \in \mathbb{N} : 0^n \text{ or } 1^n \in \mathcal{L}(\mathbf{y})\}$ , where we allow  $m$  and  $m'$  to equal  $\infty$ . We have  $\mathbf{b}_x^{(k+1)}(r) = \mathbf{p}_t(r)$  for all  $0 \leq r < 2^k$ . Setting  $Z(n, 0) := (2^k - 1)\#X_y(n) + \mathbf{b}_y^{(1)}(n)$ , for all  $n \geq 1$  we have*

$$\mathbf{b}_x^{(k+1)}(2^k n) = Z(n, 0) - \begin{cases} 2^k, & \text{if } n < m; \\ 1, & \text{if } n = m < m'; \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

For all  $n \geq 1$  and  $0 < r < 2^k$ , setting  $Z(n, r) := (r-1)\#X_y(n+1) + (2^k - r - 1)\#X_y(n) + \#Y_y(n)$ , we have

$$\mathbf{b}_x^{(k+1)}(2^k n + r) = Z(n, r) - \begin{cases} 2^k, & \text{if } n+1 < m; \\ (2^k - r + 1), & \text{if } n+1 = m < m'; \\ (2^k - 2(r-1)), & \text{if } n+1 = m = m' \text{ and } r \leq 2^{k-1}; \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

*Proof.* Lemma 4.8 implies the formula for lengths less than  $2^k$ . The proof strategy to establish formulas (8) and (9) is as follows. For  $n \geq 1$  and  $0 \leq r < 2^k$ , we first obtain an upper bound  $Z(n, r)$  on  $\mathbf{b}_x^{(k+1)}(2^k n + r)$  by counting the different  $\varphi^k$ -factorizations up to the equivalence implied by Proposition 4.17. Then we establish the exact formula by subtracting the number of  $(k+1)$ -binomial classes that admit several  $\varphi^k$ -factorizations as counted above. To do so, we use the following argument. By Proposition 4.17 and the observation made right after its statement, those factors of  $\mathbf{x}$  that admit several  $\varphi^k$ -factorizations are  $(k+1)$ -binomially equivalent only to themselves. In fact, such factors are well-understood by Lemma 4.11; they only admit two distinct  $\varphi^k$ -factorizations. Hence, counting the number of factors that have two  $\varphi^k$ -factorizations and subtracting that number from the term  $Z(n, r)$  gives the number of  $(k+1)$ -binomial equivalence classes.

We first prove formula (8) by inspecting factors of length  $2^k n$  for some  $n \geq 1$ . They are of the following two forms: either  $\varphi^k(u)$ , with  $|u| = n$ , or  $p\varphi^k(v)s$ , with  $|v| = n-1$ ,  $p$  and  $s$  non-empty. Each abelian equivalence class in  $\mathcal{L}_n(\mathbf{y})/\sim_1$  gives a  $(k+1)$ -equivalence class of factors of the first form by Proposition 3.9 (recall that  $\varphi$  is Parikh-collinear). Hence the term  $\mathbf{b}_y^{(1)}(n)$  in  $Z(n, 0)$ . For the factors of the second form, we notice the following. Such a factor has the  $\varphi^k$ -ancestor  $a\psi b$ , with  $(a, \Psi(v), b) \in X_y(n)$ . On the other hand, any  $(a, \Psi(v), b) \in X_y(n)$  gives rise to  $(2^k - 1)$   $(k+1)$ -binomial equivalence classes, namely, those represented by the words

$$\text{suff}_i(\varphi^k(a)) \varphi^k(v) \text{pref}_{2^k-i}(\varphi^k(b)), \quad 1 \leq i < 2^k,$$

where, for a word  $w$  and  $i \in \{1, \dots, |w|\}$ , we let  $\text{pref}_i(w)$  (resp.,  $\text{suff}_i(w)$ ) denote the length- $i$  prefix (resp., suffix) of  $w$ . Hence the term  $(2^k - 1)\#X_y(n)$  in the formula. Therefore we have established the upper bound  $\mathbf{b}_x^{(k+1)}(2^k n) \leq Z(n, 0)$ , with  $Z(n, 0) = (2^k - 1)\#X_y(n) + \mathbf{b}_y^{(1)}(n)$ .

As explained at the beginning of the proof, we now examine factors admitting several  $\varphi^k$ -factorizations and subtract their number from  $Z(n, 0)$  to establish formula (8). Let  $x$  be such a factor. Then it has, by Lemma 4.11, exactly two  $\varphi^k$ -factorizations, and we may write

$$p\varphi^k(u)s = x = p'\varphi^k(u')s'. \quad (10)$$

Here we note that  $|ps| = 2^k$  if and only if  $p$  or  $s$  non-empty. Moreover, the corresponding  $\varphi^k$ -ancestors are powers of letters, and given in Table 2.

fact.	ancest.	conditions on $p, s$	2nd fact.	ancest.
$p\varphi^k(a^n)s$	$a^n$	$ p  = 0 =  s $	$\varphi^{k-1}(a)\varphi^k(\bar{a}^{n-1})\varphi^{k-1}(\bar{a})$	$\bar{a}^{n+1}$
$p\varphi^k(a^{n-1})s$	$a^{n+1}$	$p = \varphi^{k-1}(\bar{a}), s = \varphi^{k-1}(a)$	$\varphi^k(\bar{a}^n)$	$\bar{a}^n$
$p\varphi^k(a^{n-1})s$	$a^{n+1}$	$ p  > 2^{k-1}, 0 <  s  < 2^{k-1}$	$p'\varphi^k(\bar{a}^{n-1})\varphi^{k-1}(\bar{a})s$	$\bar{a}^{n+1}$
$p\varphi^k(a^{n-1})s$	$a^{n+1}$	$0 <  p  < 2^{k-1},  s  > 2^{k-1}$	$p\varphi^{k-1}(a)\varphi^k(\bar{a}^{n-1})s'$	$\bar{a}^{n+1}$

Table 2: Factors of length  $2^k n$  of  $x$  admitting two  $\varphi^k$ -factorizations and their potential  $\varphi^k$ -ancestors. In the third row  $p'$  is defined by  $p = p'\varphi^{k-1}(\bar{a})$ , and in the fourth row,  $s'$  is defined by  $s = \varphi^{k-1}(a)s'$ .

In particular, for  $x$  to have two  $\varphi^k$ -factorizations (and so for a class to have been counted twice),  $a^n$  and  $\bar{a}^n$  both must appear in  $y$ , and at least one of  $a^{n+1}$  and  $\bar{a}^{n+1}$  has to also appear in the word. We divide the proof into three cases.

**Case 1.** If  $n > m$  or when  $n = m = m'$ , then as concluded above, there is no factor having several  $\varphi^k$ -factorizations, and the formula holds.

**Case 2.** Assume now that  $n < m$ , so both  $a^{n+1}$  and  $\bar{a}^{n+1}$  appear in  $y$ . Reusing Table 2 above, we see that any equivalence class corresponding to a factor having  $\varphi^k$ -ancestor  $a^{n+1}$  (or  $a^n$ ) has been counted twice (and corresponds to a word with  $\varphi^k$ -ancestor  $\bar{a}^{n+1}$  or  $\bar{a}^n$ ). There are  $2^k$  of those, whence the formula for  $n < m$ .

**Case 3.** Assume finally that  $n = m$  and  $m' > m$ . Assume without loss of generality that  $a^{m+1}$  appears in  $y$ . Therefore  $\bar{a}^{m+1}$  does not appear in  $y$ . Thus, if  $x$  has two  $\varphi^k$ -ancestors, one of them is  $\bar{a}^m$ . There is only one such factor, and this proves the remaining case in the formula.

We now turn to the proof of formula (9), and consider factors of the length  $2^k n + r$ , with  $n \geq 1$  and  $0 < r < 2^k$ . Let us first establish the upper bound  $Z(n, r)$  on  $b_x^{(k+1)}(2^k n + r)$ . Each element of  $Y_y(n)$  gives rise to a unique  $(k+1)$ -factorization; for example,  $(a, \Psi(u))$  gives the class represented by  $\text{suff}_r(\varphi^k(a))\varphi^k(u)$ . Hence the term  $\#Y_y(n)$  in  $Z(n, r)$ . Each element of  $X_y(n+1)$  gives rise to  $r-1$  many  $(k+1)$ -factorizations as follows:  $(a, \Psi(u), b)$  gives

$$\text{suff}_i(\varphi^k(a))\varphi^k(u)\text{pref}_{r-i}(\varphi^k(b)), \quad 1 \leq i < r.$$

Similarly each element of  $X_y(n)$  gives  $2^k - r - 1$  elements, namely  $(a, \Psi(u'), b)$  gives

$$\text{suff}_i(\varphi^k(a))\varphi^k(u')\text{pref}_{2^k+r-i}(\varphi^k(b)) \quad r < i < 2^k.$$

Hence  $b_x^{(k+1)}(2^k n + r) \leq Z(n, r)$  with  $Z(n, r) = (r-1)\#X_y(n+1) + (2^k - r - 1)\#X_y(n) + \#Y_y(n)$ .

We again face the problem of over-counting. As previously, we count the number of factors that have two  $\varphi^k$ -factorizations. Let again  $x$  have two  $\varphi^k$ -factorizations as in (10), where now  $ps$  and  $p's'$  are both non-empty. The  $\varphi^k$ -ancestor of each factorization is given in Table 3. We conclude that, for  $x$  to have two  $\varphi^k$ -factorizations (and so for a class to have been counted twice), we must have  $n+1 \leq m$ . We divide the proof into three cases.

**Case 1.** Assume that  $n+1 > m$ . As concluded above, there is no factor having several  $\varphi^k$ -factorizations, and the formula holds.

**Case 2.** Assume that  $n+1 < m$ . In this case we have  $1^{n+2}$  and  $0^{n+2}$  appearing in  $y$ . We claim that any factor with a  $\varphi^k$ -ancestor  $a^{n+2}$  or  $a^{n+1}$  has also a  $\varphi^k$ -ancestor  $\bar{a}^{n+2}$  or  $\bar{a}^{n+1}$ , and we show there are  $2^k$  such factors. Hence the formula follows.



fact.	ancestor	conds. on p, s
$p\varphi^k(u)s$	$a^{n+2}$	if $ ps  = r$ and $p, s \neq \varepsilon$
	$a^{n+1}$	if $ ps  = r$ with $p = \varepsilon$ or $s = \varepsilon$ , or if $ ps  > r$
$p'\varphi^k(u')s'$	$\bar{a}^{n+2}$	if $ p's'  = r$ and $p', s' \neq \varepsilon$
	$\bar{a}^{n+1}$	if $ p's'  = r$ with $p' = \varepsilon$ or $s' = \varepsilon$ , or if $ p's'  > r$

Table 3: Factors of  $x$  of length  $2^k n + r$ ,  $0 < r < 2^k$ , admitting two  $\varphi^k$ -factorizations together with their potential  $\varphi^k$ -ancestors.

conditions on p, s	$\varphi^k$ -fact.	$\varphi^k$ -ancestor
$ p ,  s  < 2^{k-1}$	$p'\varphi^k(\bar{a}^{n-1})s'$	$\bar{a}^{n+1}$
$ p  \geq 2^{k-1}$ and $ s  < 2^{k-1}$	$p'\varphi^k(\bar{a}^n)\varphi^{k-1}(\bar{a})s$	$\bar{a}^{n+2}$
$ p  < 2^{k-1}$ and $ s  \geq 2^{k-1}$	$p\varphi^{k-1}(a)\varphi^k(\bar{a}^n)s'$	$\bar{a}^{n+2}$

Table 4: The accompanying second  $\varphi^k$ -factorization of a factor  $p\varphi^k(a^n)s$ ,  $|ps| = r$ , having two  $\varphi^k$ -factorization. Here  $p'$  and  $s'$  are suitably chosen; for example, in the first row we have  $p' = p\varphi^{k-1}(a)$  and  $s' = \varphi^{k-1}(\bar{a})s$ .

First, if  $x = p\varphi^k(a^n)s$ , with  $|ps| = r$ , we have  $x = p\varphi^{k-1}(a)\varphi^k(\bar{a}^{n-1})\varphi^{k-1}(\bar{a})s$ . The other  $\varphi^k$ -factorization of  $x$  is given in Table 4. (Observe that  $|p| \geq 2^{k-1}$  and  $|s| \geq 2^{k-1}$  cannot simultaneously hold as  $|ps| = r < 2^k$ .)

Second, if  $x = p\varphi^k(a^{n-1})s$  with  $|ps| = 2^k + r$ ,  $r < |p| < 2^k$ , the other  $\varphi^k$ -factorization of  $x$  is given in Table 5, where  $p'$  and  $s'$  are again suitably chosen. This concludes the proof for this part,

conditions on p, s	$\varphi^k$ -fact.	$\varphi^k$ -ancestor
$ p  \geq 2^{k-1},  s  < 2^{k-1}$	$p'\varphi^k(\bar{a}^{n-1})\varphi^{k-1}(\bar{a})s$	$\bar{a}^{n+1}$
$ p  = 2^{k-1},  s  = 2^{k-1} + r$	$\varphi^k(\bar{a}^n)s'$	$\bar{a}^{n+1}$
$ p ,  s  > 2^{k-1}$	$p'\varphi^k(\bar{a}^n)s'$	$\bar{a}^{n+2}$
$ p  = 2^{k-1} + r,  s  = 2^{k-1}$	$p'\varphi^k(\bar{a}^n)$	$\bar{a}^{n+1}$
$ p  < 2^{k-1},  s  \geq 2^{k-1}$	$p\varphi^{k-1}(a)\varphi^k(\bar{a}^{n-1})s'$	$\bar{a}^{n+1}$

Table 5: The accompanying second  $\varphi^k$ -factorization of a factor  $p\varphi^k(a^{n-1})s$ ,  $|ps| = 2^k + r$ ,  $r < |p| < 2^k$ , having two  $\varphi^k$ -factorizations. Here  $p'$  and  $s'$  are again suitably chosen and can be inferred from the other  $\varphi^k$ -ancestor and the length constraints.

as we have exhibited  $2^k$  distinct factors, and there are no other possibilities. (Indeed, there are  $r - 1$  factors having  $a^{n+2}$  as a  $\varphi^k$ -ancestor, and  $2^k - r + 1$  factors having  $a^{n+1}$  as such.)

**Case 3.** Assume finally that  $n + 1 = m$ . We divide the proof into two subcases.

**Case 3.1.** Assume that  $m < m'$ . Let us assume that  $a^{n+1}$  appears in  $y$  but  $a^{n+2}$  does not. Then  $\bar{a}^{n+2}$  does under the assumption. Notice that in the previous case there were exactly  $r - 1$  factors having  $a^{n+2}$  as a  $\varphi^k$ -ancestor. Under our current assumption, these factors do not have this ancestor, but have instead the ancestor  $\bar{a}^{n+2}$ . They thus have only one  $\varphi^k$ -factorization. However, as before, the  $2^{k-1} - r + 1$  factors with  $\varphi^k$ -ancestor  $a^{n+1}$  have a second  $\varphi^k$ -factorization. We conclude that the formula holds also in this case.

**Case 3.2.** Finally assume that  $m' = m$ . Then we have that neither  $a^{n+2}$  nor  $\bar{a}^{n+2}$  appears in  $y$ , while both  $a^{n+1}$  and  $\bar{a}^{n+1}$  do. We thus need to count those factors that have both  $a^{n+1}$

and  $\bar{a}^{n+1}$  as  $\varphi^k$ -ancestors. Looking at the previous table, only the center row gives  $\bar{a}^{n+2}$  as a  $\varphi^k$ -ancestor. Such factors appear when  $|ps| = 2^k + r$  with  $2^{k-1} < |p| < 2^{k-1} + r$ , i.e., there are  $r - 1$  of them whenever  $r \leq 2^{k-1}$  (recall that  $|p|, |s| < 2^k$ ). Symmetric arguments apply to factors having  $\varphi^k$ -ancestors  $\bar{a}^{n+1}$  (i.e., exchanging the role of  $a$  and  $\bar{a}$ ). We conclude that the number of factors having two  $\varphi^k$ -factorizations is  $2^k - 2(r - 1)$  when  $r \leq 2^{k-1}$ , as is claimed in the formula.

We are left with the case that  $r > 2^{k-1}$ . Here we show that no factor has two  $\varphi^k$ -factorizations with respect to  $\mathbf{y}$ . Since  $r > 2^{k-1}$ , we have  $|ps| = 2^k + r > 2^k + 2^{k-1}$  with  $|p|, |s| < 2^k$ . It follows that for such  $\varphi^k$ -factorizations we must have  $|p|, |s| > 2^{k-1}$ , which only leaves the center row of the previous table. But, we already discarded these factors, so the proof is completed.  $\square$

## 5.2 A converse to Theorem 4.2

As announced at the beginning of the section, we now obtain a partial converse statement to Theorem 4.2. Before giving the proof, which is quite long and technical, we give a brief sketch of it. The proof is by induction on  $k$ . The induction hypothesis allows to conclude that  $\mathbf{x}$  in the statement is essentially the  $k$ th image of a recurrent word  $\mathbf{z}$ . We then show that  $\mathbf{z}$  has property  $\mathcal{TM}\mathcal{B}(1)$  using several times the formulas established in Proposition 5.5. The word  $\mathbf{z}$  having property  $\mathcal{TM}\mathcal{B}(1)$  allows to show that  $\mathbf{x}$  is essentially the  $(k + 1)$ st image of another binary word  $\mathbf{y}$  which then suffices for the claim by Theorem 5.2.

*Proof of Theorem 5.2.* Observe first that  $\mathbf{x}$  is aperiodic; we shall implicitly use this fact throughout the proof. Indeed, if it was not aperiodic, it would be purely periodic by the recurrence assumption. However, purely periodic words have  $b^{(1)}(n) = 1$  for infinitely many  $n$ . This would contradict the assumption that  $\mathbf{x}$  has property  $\mathcal{TM}\mathcal{B}(1)$ .

We shall prove the claim by induction. So let first  $k = 1$ . Then Theorem 5.1 asserts that there exist  $a \in \{\varepsilon, 0, 1\}$  and a binary word  $\mathbf{y}$  such that  $\mathbf{x} = a\varphi(\mathbf{y})$ , which was to be proven. Assume then that the claim holds for some  $k$  and assume further that  $\mathbf{x}$  has property  $\mathcal{TM}\mathcal{B}(k + 1)$ . It follows that  $\mathbf{x} = u'\varphi^k(\mathbf{z})$ , where  $u'$  is a proper (possibly empty) suffix of  $\varphi^k(0)$  or  $\varphi^k(1)$ .

**Claim 1.** *If  $\mathbf{z}$  has property  $\mathcal{TM}\mathcal{B}(1)$ , then  $\mathbf{x}$  is of the form  $\mathbf{x} = u\varphi^{k+1}(\mathbf{y})$  with  $u$  a suffix of  $\varphi^{k+1}(0)$  or  $\varphi^{k+1}(1)$ .*

*Proof of claim:* The assumption implies that  $\mathbf{z} = b\varphi(\mathbf{y})$  for some  $b \in \{\varepsilon, 0, 1\}$ , whence  $\mathbf{x} = u'\varphi^k(b)\varphi^{k+1}(\mathbf{y})$ . Let  $y$  be a prefix of  $\mathbf{y}$  that contains both letters 0 and 1. Then the factor  $Y = \varphi^{k+1}(y)$  has a unique  $\varphi^{k+1}$ -factorization by Lemma 4.11. Now  $u'\varphi^k(b)Y$  appears also in  $\varphi^{k+1}(\mathbf{y})$  due to  $\mathbf{x}$  being recurrent. In particular, it admits a  $\varphi^{k+1}$ -factorization, and since  $Y$  has a unique  $\varphi^{k+1}$ -factorization, we conclude that  $u'\varphi^k(b)$  must be the suffix of  $\varphi^{k+1}(0)$  or  $\varphi^{k+1}(1)$ .  $\blacksquare$

To prove the theorem, it is thus enough to show that  $\mathbf{z}$  has  $\mathcal{TM}\mathcal{B}(1)$ . Indeed, then  $\mathbf{x}$  is a suffix of the word of the form  $\varphi^{k+1}(\mathbf{y})$  with  $\mathbf{y}$  aperiodic, and Theorem 4.2 gives the claim. Notice that  $b_x^{(k+1)} = b_{\varphi^k(\mathbf{z})}^{(k+1)}$  again due to recurrence of  $\mathbf{x}$ . This fact is again used throughout the rest of the proof.

Notice now that  $\mathbf{z}$  is also recurrent: if it is not recurrent, then it has a prefix  $w$ , containing both letters, which appears only once in  $\mathbf{z}$ . Let us write  $\mathbf{z} = w\mathbf{z}'$ . However,  $\varphi^k(w)$  appears in  $\varphi^k(\mathbf{z}')$  by the recurrence of  $\mathbf{x}$ . Since  $w$  contains both letters, the  $\varphi^k$ -factorization of  $\varphi^k(w)$  is unique. But, since  $\varphi^k$  is injective, we must find  $w$  in  $\mathbf{z}'$ , a contradiction.

Let us assume towards a contradiction that  $\mathbf{z}$  does not have  $\mathcal{TM}\mathcal{B}(1)$ , and let  $n$  be the least integer for which

$$b_z^{(1)}(n) \neq b_t^{(1)}(n) = \begin{cases} 2, & \text{if } n \text{ is odd;} \\ 3, & \text{if } n \text{ is even.} \end{cases}$$

We now divide the proof into two cases, depending on the parity of  $n$ . As it appears, the case where  $n$  is even is easier to handle.

### 5.2.1 $n$ is even

By definition of  $n$ ,  $b_z^{(1)}(n-1) = b_t^{(1)}(n-1) = 2$  and  $b_z^{(1)}(n) \neq 3$ . Note however that  $b_z^{(1)}(n) \neq 1$  because  $\mathbf{z}$  is aperiodic. Since  $b_z^{(1)}$  can increase or decrease by at most 1 between consecutive values, we conclude that  $b_z^{(1)}(n) = 2$ .

**Claim 2.** *We have  $n > m$ , where  $m$  is as in Proposition 5.5.*

*Proof of claim:* If  $n > 2$  (i.e.,  $n \geq 4$ ), then  $m = 2$  (and  $m' = 2$ ) because  $b_z^{(1)}(2) = 3$  implies that  $00, 11 \in \mathcal{L}(\mathbf{z})$  while  $b_z^{(1)}(3) = 2$  implies that  $000, 111 \notin \mathcal{L}(\mathbf{z})$ . If  $n = 2$  then  $m = 1$ . ■

We next show that  $\#X_z(n) \leq 5$ . To this end, let  $M_n = \max\{|u|_1 : u \in \mathcal{L}_n(\mathbf{z})\}$ , i.e., the maximum weight among length- $n$  factors of  $\mathbf{y}$ . Assume first that  $M_{n-1} = M_n$ . Then any factor  $v \in \mathcal{L}_{n-1}(\mathbf{z})$  with  $|v|_1 = M_{n-1}$  is followed and preceded by 0 (except possibly for the prefix, which is still followed by 0), as otherwise  $M_n \neq M_{n-1}$ . We conclude that

$$X_z(n) \subseteq \{(0, \Psi(v), 0)\} \cup (\{0, 1\} \times \{\Psi(v) + (1, -1)\} \times \{0, 1\}),$$

so the claim follows.

Assume second that  $M_{n-1} = M_n - 1$ . Then each factor  $v \in \mathcal{L}_{n-1}(\mathbf{z})$  with  $|v|_1 = M_{n-1} - 1$  is followed and preceded by 1; this is because  $b_z^{(1)}(n-1) = 2 = b_z^{(1)}(n)$ . In this case

$$X_z(n) \subseteq \{(1, \Psi(v), 1)\} \cup (\{0, 1\} \times \{\Psi(v) + (-1, 1)\} \times \{0, 1\}).$$

We have shown  $\#X_z(n) \leq 5$ . This however leads to a contradiction: applying formula (8), we find

$$3 \cdot 2^{k+1} - 3 = b_t^{(k+1)}(2^k n) = b_x^{(k+1)}(2^k n) = b_{\varphi^k(\mathbf{z})}^{(k+1)}(2^k n) \leq 5 \cdot (2^k - 1) + 2 = 3 \cdot 2^{k+1} - 2^k - 3,$$

where the leftmost equality follows from  $n$  being even and (3). We hence move to the case where  $n$  is odd.

### 5.2.2 $n$ is odd

We have that  $n \geq 3$  is odd (as  $\mathbf{z}$  is binary). Since  $b_z^{(1)}(n-1) = 3$ , we have  $b_z^{(1)}(n) \in \{3, 4\}$  arguing as in the case when  $n$  was even.

**Claim 3.** *We have  $b_z^{(1)}(n) = 3$ .*

*Proof of claim:* Assume for a contradiction that  $b_z^{(1)}(n) = 4$ . As  $\mathbf{z}$  recurrent, the abelian Rauzy graph  $G_n$  has a subgraph of the form depicted in Fig. 3. Further, since  $\mathbf{z}$  is also aperiodic, it must

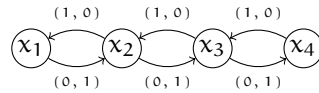


Figure 3: A subgraph of the order- $n$  abelian Rauzy graph of  $\mathbf{z}$  in Claim 3.

have at least one loop (at a right special vertex). We conclude that  $G_n$  has at least seven edges, that is,  $\#X_z(n) \geq 7$ . Since  $n$  is odd we have, using (8) and recalling that  $b_x^{(k+1)} = b_{\varphi^k(\mathbf{z})}^{(k+1)}$

$$3 \cdot 2^{k+1} - 4 = b_t^{(k+1)}(2^k n) = b_x^{(k+1)}(2^k n) \geq \#X_z(n)(2^k - 1) + b_z^{(1)}(n) - 2^k \geq 6 \cdot 2^k - 3 = 3 \cdot 2^{k+1} - 3,$$

which is absurd. ■

Recall the entities  $m$  and  $m'$  from Proposition 5.5.

**Claim 4.** *We have  $m = 2$ . If  $n \geq 5$ , then  $m' = 2$  also. Otherwise  $n = 3$  and  $m' > 2$ .*

*Proof of claim:* For  $n \geq 5$  one may proceed as in the proof of Claim 2. When  $n = 3$ , we still have that  $b_z^{(1)}(2) = 3$  implies that both 00 and 11 appear in the word. However,  $b_z^{(1)}(3) = 3$  implies that one of 000 or 111 appears in  $z$  while the other does not. ■

**Claim 5.** *We have that  $k = 1$  and  $\#X_z(n) = 5$ .*

*Proof of claim:* Consider  $b_x^{(k+1)}(2^k n)$ ; applying (8) (using the previous claim) we have

$$(2^k - 1)\#X_z(n) + b_z^{(1)}(n) = b_x^{(k+1)}(2^k n) = b_t^{(k+1)}(2^k n) = 3 \cdot 2^{k+1} - 4$$

because  $n$  is odd. By Claim 3, this is equivalent to

$$(6 - \#X_z(n))2^k + \#X_z(n) = 7.$$

Since  $z$  is aperiodic and recurrent and the abelian Rauzy graph  $G_n$  has three vertices,  $G_n$  must have at least five edges, i.e.,  $\#X_z(n) \geq 5$ . The only way to satisfy the above equality is when  $k = 1$  and  $\#X_z(n) = 5$ . Indeed, the function  $x \mapsto (6 - x)2^k + x$  is strictly decreasing (as  $k \geq 1$ ), and for  $x = 6$ , it yields 6 which is less than the right-hand side in the above equation. Therefore we must have  $\#X_z(n) \leq 5$ . We conclude that  $\#X_z(n) = 5$ . Plugging this into the above equation, we find that  $k = 1$ , as claimed. ■

The previous claim shows that  $G_n$  is a graph with three vertices and five edges. Since  $z$  is recurrent and aperiodic,  $G_n$  can be obtained, without loss of generality, by adding one loop to the graph depicted in Fig. 4.

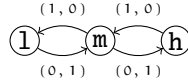


Figure 4: A subgraph of the order- $n$  abelian Rauzy graph  $G_n$  of  $z$ ;  $G_n$  is assumed to have three vertices and five edges. Here the leftmost vertex  $1$  corresponds to the lightest abelian equivalence class, the rightmost vertex  $h$  to the heaviest, and the center vertex  $m$  to the remaining class.

For any vertex  $v \in \{1, m, h\}$  of  $G_n$ , we shall refer to factors of length  $n$  having their Parikh-vector corresponding to  $v$  as  $v$ -factors.

Recall that  $M_n$  is defined as the maximum weight among factors of length  $n$ .

**Lemma 5.6.** *The graph  $G_n$  contains either the loop  $h \xrightarrow{(0,0)} h$  or the loop  $1 \xrightarrow{(1,1)} 1$ .*

*Proof.* Assume first that  $M_n = M_{n-1}$ . This implies that all the heaviest factors of length  $n - 1$  are surrounded by 0s in  $z$  (meaning, preceded and followed by 0; notice that by assumption that the prefix is not a heavy factor). Hence, there is a factor  $0u0$  in  $z$  which corresponds to the loop  $\Psi(0u) \xrightarrow{(0,0)} \Psi(u0)$ , where  $0u$  is an  $h$ -factor, because  $|u0|_1 = |u|_1 = M_{n-1} = M_n$ .

Assume then that  $M_n = M_{n-1} + 1$ . Since  $b_z^{(1)}(n) = b_z^{(1)}(n - 1) = 3$ , we must have that the minimum weight of factors of length  $n$  is one greater than that of factors of length  $n - 1$ ; thus all length- $(n - 1)$  minimum-weight factors are surrounded by 1s in  $z$ . (The only exception is the prefix, which is still followed by 1.) So any non-prefix occurrence of such a factor (recall  $z$  is recurrent) gives a loop on  $1$  with label  $(1, 1)$  similar to the above. □

Applying (9) with  $r = 1$  and  $n + 1 \geq 4 > 2 = m$  (by Claim 4), we have

$$\#Y_z(n) = b_{\varphi(z)}^{(2)}(2n + 1) = b_x^{(2)}(2n + 1) = b_t^{(2)}(2n + 1) = 8.$$

We show that this is impossible. Recall that we add either of the loops  $h \xrightarrow{(0,0)} h$  or  $1 \xrightarrow{(1,1)} 1$  to the graph  $G_n$ . In either case we note the following: all  $1$ -factors are followed by 1 and all  $h$ -factors are followed by 0. In particular, inspecting factors of length  $n + 1$ , we only have two distinct

Parikh-vectors. Therefore, the graph  $G_{n+1}$  has only two vertices, i.e.,  $b_z^{(1)}(n+1) = 2$ . The number of edges of such a graph is at most six: both vertices can have two loops and one outgoing edge to the other vertex. However, applying (8) (with  $n+1 > m$ ), we find  $\#X_z(n+1) + 2 = 9$  because  $n+1$  is even. But then  $\#X_z(n+1) = 7$ , which is impossible. This final contradiction proves that  $n$  cannot be odd either. This concludes the proof of Theorem 5.2.  $\square$

## 6 Several Answers to Question A

In this section, we are interested in Question A. Namely, does there exist an infinite word  $\mathbf{w}$  such that, for all  $k \geq 1$ ,  $b_w^{(k)}$  is unbounded and  $b_w^{(k)} \prec b_w^{(k+1)}$ ? If the answer is positive, can we find a (pure) morphic such word  $\mathbf{w}$ ?

One can give a rather direct answer to this question. Indeed, let  $\mathbf{c}$  be the binary Champernowne word, that is, the concatenation of the binary representations of the non-negative integers:  $0, 1, 10, 11, 100, 101, 110, 111, \dots$ . Notice that  $\mathbf{c}$  contains all binary words. For each  $k$ , there exist two binary words  $u, v$  such that  $u \sim_k v$  and  $u \not\sim_{k+1} v$  (see, for instance, Theorem 2.4). Therefore, the same properties hold for  $ux$  and  $vx$ , for all  $x \in \{0, 1\}^*$ , thus  $b_c^{(k)} \prec b_c^{(k+1)}$  for all  $k$ . Clearly  $b_c^{(1)}(n) = n+1$  is unbounded and so is  $b_c^{(k)}$  for  $k \geq 2$ .

Observe that  $\mathbf{c}$  is not morphic, nor *uniformly recurrent* (a word  $\mathbf{x}$  is uniformly recurrent if for each  $x \in \mathcal{L}(\mathbf{x})$  there exists  $N \in \mathbb{N}$  such that  $x$  appears in all factors in  $\mathcal{L}_N(\mathbf{x})$ ). Therefore in the rest of the section we provide more “structured” words answering Question A.

### 6.1 A Non-Binary Pure Morphic Answer

Consider the morphism  $g: \{a, 0, 1, b\}^* \rightarrow \{a, 0, 1, b\}^*$  defined by

$$a \mapsto a0b, 0 \mapsto \varphi(0), 1 \mapsto \varphi(1), b \mapsto b^2$$

where  $\varphi$  is the Thue–Morse morphism. We have  $\mathbf{g} = g^\omega(a) = a \prod_{j=0}^{\infty} \varphi^j(0)b^{2^j}$ . We show that the word  $\mathbf{g}$  answers Question A:

**Proposition 6.1.** *The abelian complexity of  $\mathbf{g}$  is unbounded and  $b_g^{(k)} \prec b_g^{(k+1)}$  for all  $k \geq 1$ .*

*Proof.* The abelian complexity of  $\mathbf{g}$  is (at least) linear, since

$$\{ |u|_b : u \in \mathcal{L}_n(\mathbf{g}) \} = \{0, \dots, n\}.$$

Furthermore, for each  $k \in \mathbb{N}$  there exist infinitely many words  $u_n, v_n \in \mathcal{L}(\mathbf{g})$  such that  $u_n \sim_k v_n$  but  $u_n \not\sim_{k+1} v_n$ : by Theorem 2.4, take  $u_n = \varphi^k(0)b^n$  and  $v_n = \varphi^k(1)b^n$ . Consequently  $b_g^{(k)} \prec b_g^{(k+1)}$  for all  $k \geq 1$ .  $\square$

### 6.2 A Binary Morphic Answer

Consider the word  $\tau(\mathbf{g})$ , where  $\mathbf{g}$  is the word defined in the previous subsection, and  $\tau$  is the coding  $a \mapsto \varepsilon, 0 \mapsto 0, 1 \mapsto 1$ , and  $b \mapsto 1$ . We have the following:

**Proposition 6.2.** *The abelian complexity of  $\tau(\mathbf{g})$  is unbounded and  $b_{\tau(\mathbf{g})}^{(k)} \prec b_{\tau(\mathbf{g})}^{(k+1)}$  for all  $k \geq 1$ .*

*Proof.* The word  $\tau(\mathbf{g})$  has unbounded abelian complexity: it contains arbitrarily long words  $u$  for which  $|u|_1 = \lfloor |u|/2 \rfloor$  (take factors of the Thue–Morse word for instance). Similarly it contains arbitrarily long powers of 1. Consequently, the word has unbounded abelian complexity (recall Lemma 2.1).

To show  $b_{\tau(\mathbf{g})}^{(k)} \prec b_{\tau(\mathbf{g})}^{(k+1)}$  for all  $k$ , we notice that the same arguments as in the case of  $\mathbf{g}$  can be applied verbatim with  $\tau(u_n)$  and  $\tau(v_n)$ .  $\square$

### 6.3 A Binary Uniformly Recurrent Answer

We note that none of the above words are uniformly recurrent. A natural candidate for such a word is one that has relatively high factor complexity. Uniformly recurrent words having positive *topological entropy*<sup>4</sup> were studied by Grillenberger in [23]. A construction for uniformly recurrent positive entropy words appears in [9, §4.4.3]; this construction is simpler than that of Grillenberger's, though some properties are lost (see [9, §4.4.3] for a discussion). We recall this construction here. To attain a word with entropy between 0 and  $\log d$ , define  $D_0 = \{0, 1, \dots, d-1\}$  and let  $(q_k)_{k \geq 0}$  be a sequence of positive integers. Assuming  $D_k$  is constructed, let  $u_k$  be the product of words of  $D_k$  in lexicographic order (assuming, e.g.,  $0 < 1 < \dots < d-1$ ). Define then  $D_{k+1} := u_k D_k^{q_k}$ . The sequence  $(u_k)_{k \in \mathbb{N}}$  converges to a uniformly recurrent word  $\mathbf{u}$  having, with a suitable choice of  $(q_k)$ , the prescribed entropy. We consider the word with  $d = 2$  and  $q_k = 2$  for all  $k$  (and are not interested in the entropy). Hence for us  $\mathbf{u} = 0100010101100111 \dots$ .

**Lemma 6.3.** *Let  $k \geq 1$ . If, for some  $j \geq 0$ ,  $D_j$  contains two words  $u, v$ , such that  $u \sim_k v$  and  $u \not\sim_{k+1} v$ , then  $D_{j+1}$  contains words  $x, y, z$  and  $w$  such that*

- $x \sim_k y$  but  $x \not\sim_{k+1} y$ ;
- $z \sim_{k+1} w$  but  $z \not\sim_{k+2} w$ .

*Proof.* By definition, the set  $D_{j+1}$  contains the words  $x = u_j u u$ ,  $y = u_j v v$ ,  $z = u_j u v$ , and  $w = u_j v u$ .

We first consider the pair  $x, y$ . Since  $\sim_k$  is a congruence,  $x \sim_k y$ . To see that  $x \not\sim_{k+1} y$ , assume the contrary, so that this equivalence reduces to  $u u \sim_{k+1} v v$  by Lemma 2.2. Lemma 2.6 implies  $u \sim_{k+1} v$ , a contradiction.

Next we have  $u v \sim_{k+1} v u$  by Theorem 2.3, and thus  $z = u_j u v \sim_{k+1} u_j v u = w$  by Lemma 2.2. Similarly  $z \sim_{k+2} w$  would imply  $u v \sim_{k+2} v u$  and thus  $u \sim_{k+1} v$  by Theorem 2.3, a contradiction. The claim follows.  $\square$

**Theorem 6.4.** *The abelian complexity of  $\mathbf{u}$  is unbounded and  $\mathbf{b}_u^{(k)} \prec \mathbf{b}_u^{(k+1)}$  for all  $k \geq 1$ .*

*Proof.* First we show that  $\mathbf{b}_u^{(1)}$  is unbounded. Assume, for some  $j \geq 0$ , that  $D_j$  contains words  $u, v$  with  $|u|_0 - |v|_0 = 2^j$  (this holds for  $j = 0$ ). Then by definition  $D_{j+1}$  contains the words  $x = u_j u u$  and  $y = u_j v v$ , for which  $|x|_0 - |y|_0 = 2(|u|_0 - |v|_0) = 2^{j+1}$ . This observation suffices for the claim by Lemma 2.1.

We then prove the second part of the statement. Observe that  $D_1$  contains the words 0101 and 0110, which are abelian equivalent, but not 2-binomially equivalent (as  $\binom{0101}{01} = 3$  and  $\binom{0110}{01} = 2$ ). The above lemma then implies that for all  $k \geq 1$  and for all  $j \geq k$ , the set  $D_j$  contains words that are  $k$ -binomially equivalent, but not  $(k+1)$ -binomially equivalent. The claim follows.  $\square$

**Remark 6.5.** It can be shown that the word  $\mathbf{u}$  above has topological entropy equal to 0. By modifying the arguments above suitably, the statement of the above theorem holds for any choice of  $d$  and  $(q_k)$  in the construction—as long as  $q_k > 1$  for infinitely many  $k$ . Note that to attain a word with positive entropy, the sequence  $(q_k)$  must satisfy this property. Hence we have: *For any positive real number  $h$  there is a uniformly recurrent  $d$ -ary word (with  $d = \max\{2, \lfloor h \rfloor + 1\}$ ) having entropy  $h$ , unbounded  $\mathbf{b}^{(1)}$ , and  $\mathbf{b}^{(k)} \prec \mathbf{b}^{(k+1)}$  for all  $k$ .*

## 7 Answer to Question B and Beyond

In this section, we are interested in Question B. Namely, for each  $\ell \geq 1$ , does there exist a word  $\mathbf{w}$  (depending on  $\ell$ ) such that  $\mathbf{b}_w^{(1)} \prec \mathbf{b}_w^{(2)} \prec \dots \prec \mathbf{b}_w^{(\ell-1)} \prec \mathbf{b}_w^{(\ell)} = \mathbf{p}_w$ ? If the answer is positive, is there a (pure) morphic such word  $\mathbf{w}$ ?

The word  $0^\omega$  gives  $\mathbf{b}^{(1)} = \mathbf{p}$ . The Fibonacci word  $\mathbf{f} = 0100101001001010010 \dots$ , the fixed point of the morphism  $0 \mapsto 01, 1 \mapsto 0$ , is a pure morphic word such that  $2 = \mathbf{b}_f^{(1)} \prec \mathbf{b}_f^{(2)} = \mathbf{p}_f$  by Theorem 2.8.

<sup>4</sup>The topological entropy of a word  $x$  is defined as the quantity  $\lim_{n \rightarrow \infty} \frac{\log p_x(n)}{n}$ , which exists for any  $x$  (see [9, §4.3.2]).

**Remark 7.1.** We notice that  $b_x^{(1)} = p_x$  cannot be attained for an aperiodic word  $x$  (indeed, there must exist a factor  $ava$ , with  $a \in A$  and  $v$  containing a letter different to  $a$ , whence  $av \sim_1 va$  with  $av \neq va$ ). In fact, the only ultimately periodic words over an  $m$ -letter alphabet  $\{a_1, \dots, a_m\}$  for which the equality holds are of the form  $a_1^{n_1} a_2^{n_2} \dots a_m^{n_m}$ ,  $n_i \in \mathbb{N}$  (up to permutation of the letters).

To answer Question B for larger values of  $k$ , we take images of a Sturmian word  $s$  by a power of the Thue–Morse morphism  $\varphi$  and we prove the following result.

**Theorem 7.2.** *Let  $\varphi$  be the Thue–Morse morphism. Let  $s$  be a Sturmian word. For each  $k \geq 0$ , the word  $s_k := \varphi^k(s)$  has*

$$b_{s_k}^{(1)} \prec b_{s_k}^{(2)} \prec \dots \prec b_{s_k}^{(k+1)} \prec b_{s_k}^{(k+2)} = p_{s_k}.$$

*In particular, putting the Fibonacci word for  $s$  gives a morphic positive answer to Question B.*

*Proof.* Observe that  $s_k$  has bounded  $(k+1)$ -binomial complexity as a straightforward application of Theorem 3.5 (because  $s$  has bounded abelian complexity), and thus  $b_{s_k}^{(k+1)} \prec p_{s_k}$ . By Corollary 4.18, we need only to show that  $b_{s_k}^{(k+2)} = p_{s_k}$ .

Let  $u$  and  $v$  be distinct factors of  $s_k$ . Assume they are  $(k+2)$ -binomially equivalent. By Proposition 4.17, we have that  $u = p\varphi^k(z)s$ ,  $v = p\varphi^k(z')s$  with  $z \sim_1 z'$ . If  $z \neq z'$ , then  $z \not\sim_2 z'$  by Theorem 2.8. But then Lemma 4.14 implies that  $\varphi^k(z) \not\sim_{k+2} \varphi^k(z')$ , contradicting the assumption. Hence we deduce that  $z = z'$ , but then  $u = v$  contrary to the assumption.  $\square$

**Remark 7.3.** In the above proof, since  $s$  is Sturmian, Theorem 2.8 says distinct factors are not 2-binomially equivalent. This means that Theorem 7.2 applies to and only to aperiodic words  $s$  such that  $b_s^{(2)} = p_s$ . The “only if”-part of the statement follows by a repeated application of Proposition 3.9 on the Thue–Morse morphism together with the fact that the morphism is injective.

## 7.1 Strengthening Question B

We answered Question B by providing a word with bounded abelian complexity. We can therefore strengthen the question with the following extra requirement.

**Question C.** For each  $\ell \geq 1$ , does there exist a word  $w$  (depending on  $\ell$ ) such that  $b_w^{(1)}$  is unbounded and

$$b_w^{(1)} \prec b_w^{(2)} \prec \dots \prec b_w^{(\ell-1)} \prec b_w^{(\ell)} = p_w?$$

If the answer is positive, can we find a (pure) morphic such word  $w$ ?

The following word answers the question for  $\ell = 3$  in the positive.

**Theorem 7.4.** *The word  $h = 0112122122212222122222 \dots$  fixed point of the morphism  $0 \mapsto 01, 1 \mapsto 12$ , and  $2 \mapsto 2$  is such that its abelian complexity  $b_h^{(1)}$  is unbounded and  $b_h^{(1)} \prec b_h^{(2)} \prec b_h^{(3)} = p_h$ .*

We obtain the previous theorem by combining the following two results.

**Proposition 7.5.** *The abelian complexity  $b_h^{(1)}$  of  $h$  is unbounded and  $b_h^{(1)}(n) < b_h^{(2)}(n) < p_h(n)$  for all  $n \geq 6$ .*

*Proof.* We claim that  $b_h^{(1)}$  is of the order  $\Theta(\sqrt{n})$ . Clearly it suffices to show the claim for the word  $h' = 0^{-1}h$ , as removing the first zero always removes exactly one abelian equivalence class: the only one that contains a zero. The resulting word  $h'$  is effectively a binary word; it is evident that the maximal number of 1’s in a word of length  $n$  is attained by the prefix of  $h'$ . This value equals the maximal  $m$  for which  $\sum_{i=1}^m i = \binom{m+1}{2} \leq n$ . Clearly  $m = \Theta(\sqrt{n})$ . By Lemma 2.1, we conclude that the abelian complexity of  $h$  is  $\Theta(\sqrt{n})$ .

Since the abelian complexity of  $\mathbf{h}$  is unbounded, so is its 2-binomial complexity. However, the 2-binomial complexity does not equal the factor complexity at lengths  $n \geq 6$ :  $\mathbf{h}$  contains both the factors  $12^{n-2}1$  and  $212^{n-4}12$  which are readily seen to be 2-binomially equivalent. (One may also invoke a result from [20] for binary alphabets.)

Finally observe that the abelian complexity does not coincide with the 2-binomial complexity either: the factors  $2^x 12^y$  with  $x + y = n - 1$  are abelian equivalent but not 2-binomially equivalent. This ends the proof.  $\square$

**Proposition 7.6.** *We have  $b_{\mathbf{h}}^{(3)} = p_{\mathbf{h}}$ .*

*Proof.* We may again discard the first 0 of  $\mathbf{h}$ , as the prefix is the only factor containing a zero. Assume to the contrary that there exist 3-binomially equivalent distinct factors  $u_1$  and  $u_2$  in  $\mathbf{h}' = 0^{-1}\mathbf{h}$ . The two factors must contain the same number of 1's, and hence at least one under the assumption that they are distinct. If the factors are of the form  $u_i = 2^{x_i} 12^{y_i}$  with  $x_1 \neq x_2$ , then the factors are not even 2-binomially equivalent. So the words contain at least two 1's. By the structure of  $\mathbf{h}$ , we may write  $u_i = 2^{x_i} 12^{a_i} 12^{a_i+1} 1 \dots 12^{a_i+t} 12^{y_i}$  for some  $t \geq 0$ ,  $a_i \in \mathbb{N}$ ,  $x_1 < a_1$  and  $y_i \leq a_i + t + 1$  for all  $i \in \{1, 2\}$ . If  $a_1 = a_2$ , then  $x_1 \neq x_2$ , and we again deduce that the factors are not even 2-binomially equivalent. So we must have  $a_1 < a_2$  without loss of generality. We show that in this case the factors are not 3-binomially equivalent. Indeed, consider the coefficient  $\binom{\cdot}{121}$ . For  $i = 1, 2$ , we clearly have

$$\binom{u_i}{121} = \binom{v_i}{121}, \quad (11)$$

where  $v_i = 12^{a_i} 12^{a_i+1} 1 \dots 12^{a_i+t+1}$  is obtained from  $u_i$  by deleting a prefix and a suffix. But, since  $a_1 < a_2$ , notice now that  $v_1$  is a proper subword of  $v_2$ , meaning that each occurrence of 121 in  $v_1$  has a corresponding occurrence in  $v_2$ . Clearly  $v_2$  will have more occurrences of 121. This combined with (11) gives the claim.  $\square$

**Remark 7.7.** Consider the morphic words  $\mathbf{h}_k = f_k^\omega(0)$ , where  $f_k: \{0, \dots, k\} \rightarrow \{0, \dots, k\}^*$  is defined by  $i \mapsto i(i+1)$  for  $i < k$ , and  $k \mapsto k$ . It can be shown that each of the words  $\mathbf{h}_k$ ,  $k \geq 2$ , has the property  $b^{(1)} \prec b^{(2)} \prec b^{(3)} = p$  (this was shown for  $\mathbf{h}_2$  in the above). Indeed, one may proceed by induction on  $k$ ; we find all factors of  $f_k$  in  $f_{k+1}$  (up to renaming the letters  $i \mapsto i+1$ ), so the first two relations hold immediately. To show the equality between  $b^{(3)}$  and  $p$ , one may proceed as in Proposition 7.6. If two factors were 3-binomially equivalent, they should contain at least one occurrence of 1 (otherwise they are factors of  $1 \prod_{i=0}^m f_k^i(2)$  for some large enough  $m$ , and are factors of  $f_{k-1} = 0 \prod_{i=0}^\omega f_{k-1}^i(1)$  after renaming letters). If they contain exactly one occurrence of 1, then inspecting the coefficients  $\binom{\cdot}{1j}$ ,  $j \geq 2$ , reveals that the maximal suffixes not containing a 1 must be of the same length. Since they are both prefixes of the same word ( $\prod_{i=0}^m f^i(2)$  for some large enough  $m$ ), they must be equal. Hence, by cancellativity, we may remove the common suffix starting from 1 and keep 3-binomial equivalence, which is a contradiction. If the words contained at least two 1s, then inspecting the coefficients  $\binom{\cdot}{121}$  would reveal that the first occurrence of 1 in both words must correspond to the same occurrence of 1 in  $\mathbf{h}_k$  (much like in the proof of Proposition 7.6). Inspecting the coefficients  $\binom{\cdot}{1j}$ , one can proceed as in the case where only one 1 was assumed to conclude with a contradiction.

A complete answer to Question C is far from obvious; especially if one wishes to obtain a pure morphic word. Conversely, for a non-periodic morphic word  $\mathbf{w}$  which is not the fixed point of a Parikh-collinear morphism, one can wonder about the existence of a minimal value  $m$  for which the binomial and factor complexities would coincide. Does there exist  $m \in \mathbb{N}$  such that  $b_{\mathbf{w}}^{(m)} = p_{\mathbf{w}}$ ?

Even with an apparently simple situation, it is far from obvious. As stated in the introduction, computing the  $k$ -binomial complexity of a particular infinite word remains challenging. The period-doubling word  $\mathbf{pd} = 01000101010001 \dots$ , the fixed point of  $\sigma: 0 \mapsto 01, 1 \mapsto 00$ , can be proved to have the following properties. Its abelian complexity  $b_{\mathbf{pd}}^{(1)}$  is unbounded [27, Lem. 4].



For the 2-binomial complexity, we have  $b_{pd}^{(2)}(2^n) = \rho_{pd}(2^n)$  for all  $n$ , but  $b_{pd}^{(2)}(n) < \rho_{pd}(n)$  for all  $n \neq 2^m$  [28, Prop. 4.5.1]. Otherwise stated,  $b_{pd}^{(1)} < b_{pd}^{(2)} < \rho_{pd}$ . Computer experiments show that  $b_{pd}^{(3)} < \rho_{pd}$  and suggest that  $b_{pd}^{(4)} = \rho_{pd}$ .

## 7.2 Completing the Binomial Complexities of $\varphi^k$ Applied to a Sturmian word

For any  $k \geq 1$ , the results presented so far imply that we have the exact  $j$ -binomial complexity function of  $\varphi^k(s)$ , with  $s$  a Sturmian word, for each  $j \neq k + 1$ . As a bonus, we compute the  $(k + 1)$ -binomial complexity in Proposition 7.9. We first analyze the abelian Rauzy graphs of Sturmian words, after which we may apply Proposition 5.5 to obtain the exact  $(k + 1)$ -binomial complexity as well.

A Sturmian word  $s$  has  $b_s^{(1)}(n) = 2$  for all  $n \geq 1$ . Hence its abelian Rauzy graph has two vertices.

**Proposition 7.8.** *Let  $G_n = (V_n, E_n)$ . We have  $\#E_1 = 3$  and  $\#E_n = 4$  for all  $n \geq 2$ . For all  $n \geq 1$ , we have  $\#E_n/\equiv_L + \#E_n/\equiv_R = 6$ .*

*Proof.* Sturmian words are aperiodic, so the graph  $G_n$  is always strongly connected. It also always has a right special vertex.

The claim is plain to verify for  $G_1$ . Indeed, only one of the vertices can have a loop, and this loop can only be labeled with  $(a, a)$ , where  $a$  is the letter corresponding to the vertex in question. The second claim is straightforward to check in this case.

We next consider  $G_n$ ,  $n \geq 2$ . Since  $G_n$  is strongly connected and has a right special vertex, we conclude that  $G_n$  is obtained by adding (possibly zero) loops to one of the graphs in Fig. 5

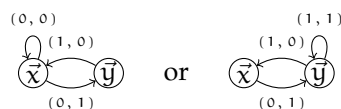


Figure 5: Possible subgraphs of a Sturmian word's abelian Rauzy graph.

By symmetry, we may assume it is the one on the left. We claim that we add exactly one loop to this graph. More precisely, the added loop is either  $\vec{x} \xrightarrow{(1,1)} \vec{x}$  or  $\vec{y} \xrightarrow{(0,0)} \vec{y}$ .

First off, we cannot add the loop  $\vec{y} \xrightarrow{(1,1)} \vec{y}$ ; otherwise we have the factors  $x0$  and  $y1$ , where  $x$  is an  $\vec{x}$ -factor and  $y$  is a  $\vec{y}$ -factor, for which we have  $|y1|_1 - |x0|_1 = 2$  contradicting balancedness at  $n + 1$ . Towards a contradiction, we consider whether we add neither or both of the remaining admissible loops.

Assume first that we add both loops. Then we have that length  $n + 1$  factors  $1x1$  and  $0y0$ , where  $1x, x1$  are  $\vec{x}$ -factors and  $0y, y0$  are  $\vec{y}$ -factors, having  $|1x1|_1 = |0y0|_1$ . However, we then have  $|y|_1 - |x|_1 = 2$ , which gives a contradiction with balancedness at  $n - 1$ .

To complete the proof of the first claim, suppose we add neither of the loops. Inspecting  $G_n$ , we see that a factor of length  $n$  is an  $\vec{x}$ -factor if and only if it begins with  $0$ . Consider the right special factor  $0v$  of length  $n$  (i.e.,  $|v| = n - 1 \geq 1$ ) (it begins with  $0$  by the form of the graph). Since  $v1$  is a  $\vec{y}$ -factor, we deduce that  $v$  begins with a  $1$ . But then  $v0$  is a  $\vec{x}$ -factor beginning with  $1$ , a contradiction.

The second part of the claim is now straightforward. The edges in the left-hand graph in Fig. 5 are all pairwise inequivalent under both  $\equiv_L$  and  $\equiv_R$ . Both of the admissible loops to be added to the graph to obtain  $G_n$  are equivalent to the non-loop edges of the graph. Hence  $\#E_n/\equiv_L + \#E_n/\equiv_R = 6$ .  $\square$

For a Sturmian word  $s$ , we have, by Proposition 7.8,  $X_s(1) = 3$ ,  $X_s(n) = 4$  for all  $n \geq 2$ , and  $Y_s(n) = 6$  for all  $n \geq 1$ . We also have that  $1 = m < m'$  using the notation of Proposition 5.5. Hence, applying the proposition, we have:

**Proposition 7.9.** For all  $n \geq 0$  and  $0 \leq r < 2^k$

$$\mathbf{b}_{\varphi^k(\mathbf{s})}^{(k+1)}(2^k n + r) = \begin{cases} \mathbf{p}_t(r), & \text{if } n = 0 \text{ and } 0 \leq r < 2^k; \\ 3 \cdot 2^k - 2, & \text{if } n = 1 \text{ and } r = 0; \\ 3 \cdot 2^k + r - 1, & \text{if } n = 1 \text{ and } r > 0; \\ 2^{k+2} - 2, & \text{otherwise.} \end{cases} \quad (12)$$

We thus conclude the following. For any  $k \geq 1$  and a Sturmian word  $\mathbf{s}$ , we have  $\mathbf{b}_{\varphi^k(\mathbf{s})}^{(j)} = \mathbf{b}_t^{(j)}$  if  $1 \leq j \leq k$  (Theorem 4.2);  $\mathbf{b}_{\varphi^k(\mathbf{s})}^{(k+1)}$  is as in (12); and  $\mathbf{b}_{\varphi^k(\mathbf{s})}^{(j)} = \mathbf{p}_{\varphi^k(\mathbf{s})}$  if  $j \geq k + 2$  (Theorem 7.2). The exact value for  $\mathbf{p}_{\varphi^k(\mathbf{s})}(n)$  is given by  $\mathbf{p}_t(n)$  when  $n \leq 2^k$  (Lemma 4.8), and by  $n + 2^{k+1} - 1$  for  $n > 2^k$ . The latter can be deduced by using the methods described in [21, §4.1].

## Acknowledgments

We thank the reviewers for their comments improving the paper.

## References

- [1] Boris Adamczewski. Balances for fixed points of primitive substitutions. *Theoret. Comput. Sci.*, 307(1):47–75, 2003. doi:10.1016/S0304-3975(03)00092-6.
- [2] Boris Adamczewski and Yann Bugeaud. On the complexity of algebraic numbers. I. Expansions in integer bases. *Ann. of Math. (2)*, 165(2):547–565, 2007. doi:10.4007/annals.2007.165.547.
- [3] Jean-Paul Allouche. Thue, combinatorics on words, and conjectures inspired by the Thue-Morse sequence. *J. Théor. Nombres Bordeaux*, 27(2):375–388, 2015. URL: [http://jtnb.cedram.org/item?id=JTNB\\_2015\\_\\_27\\_2\\_375\\_0](http://jtnb.cedram.org/item?id=JTNB_2015__27_2_375_0).
- [4] Jean-Paul Allouche and Jeffrey Shallit. The ubiquitous Prouhet–Thue–Morse sequence. In C. Ding, T. Helleseth, and H. Niederreiter, editors, *Sequences and their Applications*, pages 1–16, London, 1999. Springer London. doi:10.1007/978-1-4471-0551-0\_1.
- [5] Jean-Paul Allouche and Jeffrey Shallit. *Automatic sequences: Theory, applications, generalizations*. Cambridge University Press, Cambridge, 2003. doi:10.1017/CB09780511546563.
- [6] Adrian Atanasiu, Carlos Martín-Vide, and Alexandru Mateescu. On the injectivity of the Parikh matrix mapping. *Fund. Inform.*, 49(4):289–299, 2002.
- [7] Sergei V. Avgustinovich, Dmitrii G. Fon-Der-Flaass, and Anna E. Frid. Arithmetical complexity of infinite words. In *Words, languages & combinatorics, III (Kyoto, 2000)*, pages 51–62. World Sci. Publ., River Edge, NJ, 2003. doi:10.1142/9789812704979\_0004.
- [8] Jean Berstel, Maxime Crochemore, and Jean-Éric Pin. Thue-Morse sequence and p-adic topology for the free monoid. *Discrete Math.*, 76(2):89–94, 1989. doi:10.1016/0012-365X(89)90302-6.
- [9] Valérie Berthé and Michel Rigo, editors. *Combinatorics, automata and number theory*, volume 135 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, Cambridge, 2010. doi:10.1017/CB09780511777653.
- [10] Julien Cassaigne, Gabriele Fici, Marinella Sciortino, and Luca Q. Zamboni. Cyclic complexity of words. *J. Comb. Theory, Ser. A*, 145:36–56, 2017. doi:10.1016/j.jcta.2016.07.002.

- [11] Julien Cassaigne, Anna E. Frid, Svetlana Puzynina, and Luca Q. Zamboni. A characterization of words of linear complexity. *Proc. Amer. Math. Soc.*, 147(7):3103–3115, 2019. doi:10.1090/proc/14440.
- [12] Julien Cassaigne, Juhani Karhumäki, and Svetlana Puzynina. On  $k$ -abelian palindromes. *Inf. Comput.*, 260:89–98, 2018. doi:10.1016/j.ic.2018.04.001.
- [13] Julien Cassaigne, Juhani Karhumäki, and Alekski Saarela. On growth and fluctuation of  $k$ -abelian complexity. *Eur. J. Comb.*, 65:92–105, 2017. doi:10.1016/j.ejc.2017.05.006.
- [14] Julien Cassaigne, Gwénaél Richomme, Kalle Saari, and Luca Q. Zamboni. Avoiding Abelian powers in binary words with bounded Abelian complexity. *Int. J. Found. Comput. S.*, 22(4):905–920, 2011. doi:10.1142/S0129054111008489.
- [15] Ethan M. Coven and G. A. Hedlund. Sequences with minimal block growth. *Math. Syst. Theory*, 7(2):138–153, 1973. doi:10.1007/BF01762232.
- [16] Xavier Droubay and Giuseppe Pirillo. Palindromes and Sturmian words. *Theoret. Comput. Sci.*, 223(1-2):73–85, 1999. doi:10.1016/S0304-3975(97)00188-6.
- [17] Paul Erdős. Some unsolved problems. *Michigan Math. J.*, 4:291–300, 1958.
- [18] Gabriele Fici and Svetlana Puzynina. Abelian combinatorics on words: a survey. *Comput. Sci. Rev.*, 47:Paper No. 100532, 21, 2023. doi:10.1016/j.cosrev.2022.100532.
- [19] Pamela Fleischmann, Marie Lejeune, Florin Manea, Dirk Nowotka, and Michel Rigo. Reconstructing words from right-bounded-block words. *Int. J. Found. Comput. Sci.*, 32(6):619–640, 2021. doi:10.1142/S0129054121420016.
- [20] Stéphane Fossé and Gwénaél Richomme. Some characterizations of Parikh matrix equivalent binary words. *Inform. Process. Lett.*, 92(2):77–82, 2004. doi:10.1016/j.ip1.2004.06.011.
- [21] Anna Frid. Applying a uniform marked morphism to a word. *Discrete Math. Theor. Comput. Sci.*, 3(3):125–139, 1999. doi:10.46298/dmtcs.255.
- [22] Paweł Gawrychowski, Maria Kosche, Tore Koß, Florin Manea, and Stefan Siemer. Efficiently testing Simon’s congruence. In *38th International Symposium on Theoretical Aspects of Computer Science*, volume 187 of *LIPICs. Leibniz Int. Proc. Inform.*, pages Art. No. 34, 18. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern, 2021.
- [23] Christian Grillenberger. Constructions of strictly ergodic systems. I. Given entropy. *Z. Wahrscheinlichkeit.*, 25:323–334, 1973. doi:10.1007/BF00537161.
- [24] L. I. Kalašnik. The reconstruction of a word from fragments. In *Numerical mathematics and computer technology, No. IV (Russian)*, pages 56–57, 137. Akad. Nauk Ukrain. SSR Fiz.-Tehn.-Inst. Nizkih Temperatur, Kharkov, 1973.
- [25] Teturo Kamae and Luca Zamboni. Sequence entropy and the maximal pattern complexity of infinite words. *Ergodic Theory Dynam. Systems*, 22(4):1191–1199, 2002. doi:10.1017/S0143385702000585.
- [26] Juhani Karhumäki, Alekski Saarela, and Luca Q. Zamboni. On a generalization of abelian equivalence and complexity of infinite words. *J. Comb. Theory, Ser. A*, 120(8):2189–2206, 2013. doi:10.1016/j.jcta.2013.08.008.
- [27] Juhani Karhumäki, Alekski Saarela, and Luca Q. Zamboni. Variations of the Morse-Hedlund theorem for  $k$ -abelian equivalence. *Acta Cybernet.*, 23(1):175–189, 2017. doi:10.14232/actacyb.23.1.2017.11.

- [28] Marie Lejeune. *On the  $k$ -binomial equivalence of finite words and  $k$ -binomial complexity of infinite words*. PhD thesis, Univ. of Liège, 2021. URL: <http://hdl.handle.net/2268/259266>.
- [29] Marie Lejeune, Julien Leroy, and Michel Rigo. Computing the  $k$ -binomial complexity of the Thue–Morse word. *J. Comb. Theory, Ser. A*, 176:44, 2020. doi:10.1016/j.jcta.2020.105284.
- [30] Marie Lejeune, Michel Rigo, and Matthieu Rosenfeld. Templates for the  $k$ -binomial complexity of the Tribonacci word. *Adv. Appl. Math.*, 112:26, 2020. doi:10.1016/j.aam.2019.101947.
- [31] Julien Leroy, Michel Rigo, and Manon Stipulanti. Generalized Pascal triangle for binomial coefficients of words. *Adv. Appl. Math.*, 80:24–47, 2016. doi:10.1016/j.aam.2016.04.006.
- [32] M. Lothaire. *Combinatorics on Words*. Cambridge Mathematical Library. Cambridge University Press, 1997. doi:10.1017/CB09780511566097.
- [33] Xiao-Tao Lü, Jin Chen, Zhi-Xiong Wen, and Wen Wu. On the 2-binomial complexity of the generalized Thue–Morse words, 2021. (preprint). doi:10.48550/ARXIV.2112.05347.
- [34] Bennet Manvel, Aaron D. Meyerowitz, Allen J. Schwenk, Kenneth W. Smith, and Paul K. Stockmeyer. Reconstruction of sequences. *Discrete Math.*, 94(3):209–219, 1991. doi:10.1016/0012-365X(91)90026-X.
- [35] Hamoon Mousavi. Automatic theorem proving in Walnut, 2016. doi:10.48550/ARXIV.1603.06017.
- [36] Peter Ochsenschläger. Binomialkoeffizienten und Shuffle-Zahlen. Technischer Bericht, Fachbereich Informatik, T.H. Darmstadt, 1981.
- [37] Jarkko Peltomäki. Introducing privileged words: Privileged complexity of Sturmian words. *Theor. Comput. Sci.*, 500:57–67, 2013. doi:10.1016/j.tcs.2013.05.028.
- [38] Jean-Éric Pin and Pedro V. Silva. A noncommutative extension of Mahler’s theorem on interpolation series. *European J. Combin.*, 36:564–578, 2014. doi:10.1016/j.ejc.2013.09.009.
- [39] Gwénaél Richomme, Kalle Saari, and Luca Q. Zamboni. Balance and abelian complexity of the tribonacci word. *Adv. Appl. Math.*, 45(2):212–231, 2010. doi:10.1016/j.aam.2010.01.006.
- [40] Gwénaél Richomme, Kalle Saari, and Luca Q. Zamboni. Abelian complexity of minimal subshifts. *J. Lond. Math. Soc.*, 83(1):79–95, 2011. doi:10.1112/jlms/jdq063.
- [41] Gwénaél Richomme and Patrice Séébold. On factorially balanced sets of words. *Theoret. Comput. Sci.*, 412(39):5492–5497, 2011. doi:10.1016/j.tcs.2011.06.027.
- [42] Michel Rigo. Relations on words. *Indag. Math., New Ser.*, 28(1):183–204, 2017. doi:10.1016/j.indag.2016.11.018.
- [43] Michel Rigo and Pavel Salimov. Another generalization of abelian equivalence: binomial complexity of infinite words. *Theor. Comput. Sci.*, 601:47–57, 2015. doi:10.1016/j.tcs.2015.07.025.
- [44] Michel Rigo, Manon Stipulanti, and Markus A. Whiteland. Binomial complexities and Parikh-collinear morphisms. In Volker Diekert and Mikhail V. Volkov, editors, *Developments in Language Theory - 26th International Conference, DLT 2022, Tampa, FL, USA, May 9-13, 2022, Proceedings*, volume 13257 of *Lecture Notes in Computer Science*, pages 251–262. Springer, 2022. doi:10.1007/978-3-031-05578-2\_20.

- [45] Michel Rigo, Manon Stipulanti, and Markus A. Whiteland. Automaticity and Parikh-collinear morphisms. In *Combinatorics on words*, volume 13899 of *Lecture Notes in Comput. Sci.*, pages 247–260. Springer, Cham, [2023] ©2023. doi : 10.1007/978-3-031-33180-0\_19.
- [46] Arto Salomaa. Counting (scattered) subwords. *Bull. Eur. Assoc. Theor. Comput. Sci. EATCS*, 81:165–179, 2003.
- [47] Jeffrey Shallit. *The Logical Approach to Automatic Sequences: Exploring Combinatorics on Words with Walnut*. London Mathematical Society Lecture Note Series. Cambridge University Press, 2022. doi : 10.1017/9781108775267.
- [48] Markus A. Whiteland. Equations over the k-binomial monoids. In *WORDS 2021*, volume 12847 of *LNCS*, pages 185–197. Springer, Cham., 2021. doi : 10.1007/978-3-030-85088-3\_16.