

Aperiodic pseudorandom number generators based on infinite words

Lubomíra Balková^a, Michelangelo Bucci^b, Alessandro De Luca^c, Jiří Hladký^a, Svetlana Puzynina^{d,e}

^aDepartment of Mathematics, FNSPE, Czech Technical University in Prague, Trojanova 13, 120 00 Praha 2, Czech Republic

^bDepartment of Mathematics, University of Liège, Grande traverse 12 (B37), B-4000 Liège, Belgium

^cDIETI, Università degli Studi di Napoli Federico II, via Claudio 21, 80125 Napoli, Italy

^dLIP, ENS Lyon, 46 Allée d'Italie, Lyon 69364, France

^eSobolev Institute of Mathematics, Russia

Abstract

In this paper we study how certain families of aperiodic infinite words can be used to produce aperiodic pseudorandom number generators (PRNGs) with good statistical behavior. We introduce the *well distributed occurrences* (WELLDOC) combinatorial property for infinite words, which guarantees absence of the lattice structure defect in related pseudorandom number generators. An infinite word u on a d -ary alphabet has the WELLDOC property if, for each factor w of u , positive integer m , and vector $\mathbf{v} \in \mathbb{Z}_m^d$, there is an occurrence of w such that the Parikh vector of the prefix of u preceding such occurrence is congruent to \mathbf{v} modulo m . (The Parikh vector of a finite word v over an alphabet \mathcal{A} has its i -th component equal to the number of occurrences of the i -th letter of \mathcal{A} in v .) We prove that Sturmian words, and more generally Arnoux-Rauzy words and some morphic images of them, have the WELLDOC property. Using the TestU01 [11] and PractRand [5] statistical tests, we moreover show that not only the lattice structure is absent, but also other important properties of PRNGs are improved when linear congruential generators are combined using infinite words having the WELLDOC property.

Keywords: Pseudorandom number generator, Well distributed occurrences, Sturmian word, Arnoux-Rauzy word, Linear congruential generator

2010 MSC: 68R15, 65C10

Introduction

Pseudorandom number generators aim to produce random numbers using a deterministic process. No wonder they suffer from many defects. The most usual ones – linear congruential generators – are known to produce periodic sequences with a defect called the lattice structure. Guimond et al. [14] proved that when two linear congruential generators are combined using infinite words coding certain classes of quasicrystals or, equivalently, of cut-and-project sets, the resulting sequence is aperiodic and has no lattice structure. For some other related results concerning aperiodic pseudorandom generators we refer to [12, 13]. We mention that although the lattice structure is considered as a defect of a random number generator, it can be useful in some applications for approximation of the uniform distribution [9].

We have found a combinatorial condition – *well distributed occurrences*, or WELLDOC for short – that also guarantees absence of the lattice structure in related pseudorandom generators. The WELLDOC

Email addresses: lubomira.balkova@gmail.com (Lubomíra Balková), michelangelo.bucci@gmail.com (Michelangelo Bucci), alessandro.deluca@unina.it (Alessandro De Luca), hladky.jiri@gmail.com (Jiří Hladký), s.puzynina@gmail.com (Svetlana Puzynina)

property for an infinite word u over an alphabet \mathcal{A} means that for any integer m and any factor w of u , the set of Parikh vectors modulo m of prefixes of u preceding the occurrences of w coincides with $\mathbb{Z}_m^{|\mathcal{A}|}$ (see Definition 2.1). In other words, among Parikh vectors modulo m of such prefixes one has all possible vectors. Besides giving generators without lattice structure, the WELLDOC property is an interesting combinatorial property of infinite words itself. We prove that the WELLDOC property holds for the family of Sturmian words, and more generally for Arnoux-Rauzy words.

Sturmian words constitute a well studied family of infinite aperiodic words. Let u be an infinite word, i. e., an infinite sequence of elements from a finite set called an alphabet. The (*factor*) *complexity* function counts the number of distinct factors of u of length n . A fundamental result of Morse and Hedlund [18] states that a word u is eventually periodic if and only if for some n its complexity is less than or equal to n . Infinite words of complexity $n + 1$ for all n are called *Sturmian words*, and hence they are aperiodic words of the smallest complexity. The most studied Sturmian word is the so-called Fibonacci word

01001010010010100101001001010010...

fixed by the morphism $0 \mapsto 01$ and $1 \mapsto 0$. (See Section 2 for formal definitions.) The first systematic study of Sturmian words was given by Morse and Hedlund in [19]. Such sequences arise naturally in many contexts, and admit various types of characterizations of geometric and combinatorial nature (see, e.g., [16]).

Arnoux-Rauzy words were introduced in [1] as natural extensions of Sturmian words to multiliteral alphabets (see Definition 4.4). Despite the fact that they were introduced as generalizations of Sturmian words, Arnoux-Rauzy words display a much more complex behavior. In particular, we have two different proofs of the WELLDOC property for Sturmian words, and only one of them can be generalized to Arnoux-Rauzy words. In the sequel we provide both of them.

An infinite word with the WELLDOC property is then used to combine two linear congruential generators and form an infinite aperiodic sequence with good statistical behavior. Using the TestU01 [11] and PractRand [5] statistical tests, we have moreover shown that not only the lattice structure is absent, but also other important properties of PRNGs are improved when linear congruential generators are combined using infinite words having the WELLDOC property.

The paper is organized as follows. In the next section, we give some background on pseudorandom number generation. Next, in Section 2, we give the basic combinatorial definitions needed for our main results, including the WELLDOC property, and we prove that the WELLDOC property of u guarantees absence of the lattice structure of the PRNG based on u . In Sections 3 and 4, we prove that the property holds for Sturmian and Arnoux-Rauzy words. Finally, in Section 5, we present results of empirical tests of PRNGs based on words having the WELLDOC property.

A preliminary version of this paper [2], using the acronym *WDO* instead of WELLDOC, was presented at the WORDS 2013 conference.

1. Pseudorandom Number Generators and Lattice Structure

For the sake of our discussion, any infinite sequence of integers can be understood as a *pseudorandom number generator (PRNG)*; see also [14]. The generators the most widely used in the past – linear congruential generators – are known to suffer from a defect called the lattice structure (they possess it already from dimension 2 as shown in [17]).

Let $Z = (Z_n)_{n \in \mathbb{N}}$ be a PRNG whose output is a finite set $M \subset \mathbb{N}$. We say that Z has the *lattice structure* if there exists $t \in \mathbb{N}$ such that the set

$$\{(Z_i, Z_{i+1}, \dots, Z_{i+t-1}) \mid i \in \mathbb{N}\}$$

is covered by a family of parallel equidistant hyperplanes and at the same time, this family does not cover the whole lattice

$$M^t = \{(A_1, A_2, \dots, A_t) \mid A_i \in M \text{ for all } i \in \{1, \dots, t\}\}.$$

Recall that a *linear congruential generator* (LCG) $(Z_n)_{n \in \mathbb{N}}$ is given by parameters $a, m, c \in \mathbb{N}$ and defined by the recurrence relation $Z_{n+1} = aZ_n + c \pmod m$. Let us mention a famous example of a LCG whose lattice structure is striking. For $t = 3$, the set of triples of RANDU, i.e., $\{(Z_i, Z_{i+1}, Z_{i+2}) \mid i \in \mathbb{N}\}$ is covered by only 15 parallel equidistant hyperplanes, see Figure 1.

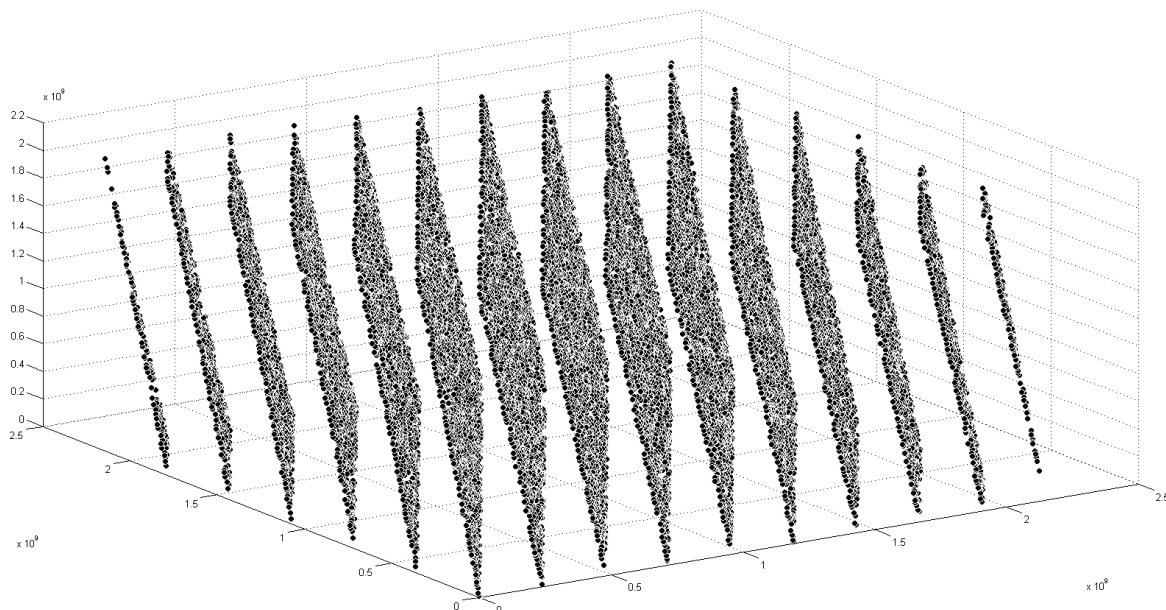


Figure 1: The triples of RANDU – the LCG with $a = (2^{16} + 3)$, $m = 2^{31}$, $c = 0$ – are covered by as few as 15 parallel equidistant planes.

In the paper of Guimond et al. [14], a restricted version of the following sufficient condition for the absence of the lattice structure is formulated.

Proposition 1.1. *Let Z be a PRNG whose output is a finite set $M \subset \mathbb{N}$ containing at least two elements. Assume there exists for any $A, B \in M$ and for any $\ell \in \mathbb{N}$ an ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ such that both $(A_1, A_2, \dots, A_\ell, A)$ and $(A_1, A_2, \dots, A_\ell, B)$ are $(\ell + 1)$ -tuples of the generator Z . Then Z does not have the lattice structure.*

Remark 1.2. Proposition 1.1 can be reformulated in terms of combinatorics on words (see Section 2) as follows: Let Z be a PRNG whose output is a finite set $M \subset \mathbb{N}$ containing at least two elements. If for any $A, B \in M$ and any length ℓ Z has a right special factor of length ℓ with right extensions A and B , then Z does not have the lattice structure.

Since Proposition 1.1 is formulated for a restricted class of generators in [14] (see Lemma 2.3 ibidem), we will provide its proof. However, we point out that all ideas of the proof are taken from [14]. We start with an auxiliary lemma.

Let us set $\lambda = \gcd\{A - B \mid A, B \in M\}$.

Lemma 1.3. *Let Z be a PRNG satisfying all assumptions of Proposition 1.1. Let $\bar{\mathbf{n}}$ be the unit normal vector of a family of parallel equidistant hyperplanes covering all t -tuples of Z . Assume $\bar{\mathbf{e}}_i$ (the i -th vector of the canonical basis of the Euclidean space \mathbb{R}^t) is not orthogonal to $\bar{\mathbf{n}}$. Then the distance d_i of adjacent hyperplanes in the family along $\bar{\mathbf{e}}_i$ is of the form λ/k for some $k \in \mathbb{N}$.*

Remark 1.4. The distance d_i of adjacent hyperplanes W_0, W_1 along $\bar{\mathbf{e}}_i$ means $|x_i - y_i|$ for any $\bar{\mathbf{x}} \in W_0$ and $\bar{\mathbf{y}} \in W_1$, where the j -th components of $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ satisfy $x_j = y_j$ for all $j \in \{1, \dots, t\}, j \neq i$. This is a well defined term because the hyperplanes in the family are of the form $W_j \equiv \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} = \alpha + jd, j \in \mathbb{Z}$, where d is the distance of adjacent hyperplanes in the family and \cdot denotes the standard scalar product. Thus, without loss of generality, consider the adjacent hyperplanes

$$W_0 \equiv \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} = \alpha \quad \text{and} \quad W_1 \equiv \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} = \alpha + d.$$

Then for any $\bar{\mathbf{x}} \in W_0$ and $\bar{\mathbf{y}} = \bar{\mathbf{x}} + s\bar{\mathbf{e}}_i$ from W_1 , we have

$$\begin{aligned} \bar{\mathbf{y}} \cdot \bar{\mathbf{n}} &= \alpha + d = \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} + d, \\ \bar{\mathbf{y}} \cdot \bar{\mathbf{n}} &= \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} + s\bar{\mathbf{e}}_i \cdot \bar{\mathbf{n}} = \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} + sn_i, \end{aligned}$$

where n_i is the i -th component of $\bar{\mathbf{n}}$. Consequently, $d_i = |s| = \left| \frac{d}{n_i} \right|$ and is the same for any choice of $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ which differ only in their i -th component and belong to adjacent hyperplanes.

Proof of Lemma 1.3. Let us start with a useful observation. Let $\bar{\mathbf{z}}$ belong to a hyperplane W of the family in question.

1. If $\bar{\mathbf{e}}_j$ is orthogonal to $\bar{\mathbf{n}}$, then we may change the j -th component of $\bar{\mathbf{z}}$ in an arbitrary way and the resulting vector will belong to the same hyperplane, i.e., if $W \equiv \bar{\mathbf{x}} \cdot \bar{\mathbf{n}} = \alpha$, then clearly $(\bar{\mathbf{z}} + \beta\bar{\mathbf{e}}_j) \cdot \bar{\mathbf{n}} = \bar{\mathbf{z}} \cdot \bar{\mathbf{n}} = \alpha$ for any $\beta \in \mathbb{R}$, thus $\bar{\mathbf{z}} + \beta\bar{\mathbf{e}}_j$ belongs to W .
2. If $\bar{\mathbf{e}}_j$ is not orthogonal to $\bar{\mathbf{n}}$ and the distance d_j of adjacent hyperplanes along $\bar{\mathbf{e}}_j$ in the family is of the form λ/k for some $k \in \mathbb{N}$, then $\bar{\mathbf{z}} + r\lambda\bar{\mathbf{e}}_j$ belongs to the family for any $r \in \mathbb{Z}$. This follows from a repeated application of the fact that if $\bar{\mathbf{z}}$ belongs to a hyperplane W , then $\bar{\mathbf{z}} + \frac{\lambda}{k}\bar{\mathbf{e}}_j$ belongs to an adjacent hyperplane of W .

Let us proceed by contradiction, i.e., we assume that there exists $i \in \{1, \dots, t\}$ such that $\bar{\mathbf{e}}_i$ is not orthogonal to $\bar{\mathbf{n}}$ and the distance along $\bar{\mathbf{e}}_i$ of adjacent hyperplanes of the family in question is not of the form λ/k , $k \in \mathbb{N}$. Let ℓ denote the largest of such indices. Choose $A, B \in M$ arbitrarily. According to assumptions, there exists an $(\ell - 1)$ -tuple $(A_1, A_2, \dots, A_{\ell-1})$ such that both $(A_1, A_2, \dots, A_{\ell-1}, A)$ and $(A_1, A_2, \dots, A_{\ell-1}, B)$ are ℓ -tuples of Z . It is therefore possible to find two t -tuples of Z such that the first one is of the form $(A_1, A_2, \dots, A_{\ell-1}, A, A_{\ell+1}, \dots, A_t)$ and the second one of the form $(A_1, A_2, \dots, A_{\ell-1}, B, \hat{A}_{\ell+1}, \dots, \hat{A}_t)$. These two t -tuples – considered as vectors in \mathbb{R}^t – belong by the assumption of Lemma 1.3 to some hyperplanes in the family. Since all vectors $\bar{\mathbf{e}}_j$, $j \in \{\ell + 1, \dots, t\}$ are either orthogonal to $\bar{\mathbf{n}}$ or the distance of adjacent hyperplanes along $\bar{\mathbf{e}}_j$ is of the form λ/k for some $k \in \mathbb{N}$, we can change the last $t - \ell$ coordinates $\hat{A}_{\ell+1}, \dots, \hat{A}_t$ of the second vector to arbitrary values from M (we transform them into $A_{\ell+1}, \dots, A_t$) and it will still belong to a hyperplane in the family. This is a consequence of the observation at the beginning of this proof. Hence, both vectors $(A_1, A_2, \dots, A_{\ell-1}, A, A_{\ell+1}, \dots, A_t)$ and $(A_1, A_2, \dots, A_{\ell-1}, B, A_{\ell+1}, \dots, A_t)$ belong to some hyperplanes of the family. Their distance along $\bar{\mathbf{e}}_\ell$ equals $|A - B|$, i.e., d_ℓ divides $A - B$. Since A, B have been chosen arbitrarily, it follows that d_ℓ divides λ , i.e., $\lambda = kd_\ell$ for some $k \in \mathbb{N}$, which is a contradiction with the choice of $\bar{\mathbf{e}}_\ell$. \square

Proof of Proposition 1.1. Let $\bar{\mathbf{n}}$ be the unit normal vector of a family of parallel equidistant hyperplanes covering all t -tuples of Z . Suppose without loss of generality that $\bar{\mathbf{e}}_1, \dots, \bar{\mathbf{e}}_\ell$ are not orthogonal to $\bar{\mathbf{n}}$ and $\bar{\mathbf{e}}_{\ell+1}, \dots, \bar{\mathbf{e}}_t$ are orthogonal to $\bar{\mathbf{n}}$. Let $\bar{\mathbf{z}} = (Z_n, Z_{n+1}, \dots, Z_{n+t-1})$ be a t -tuple of Z , thus $\bar{\mathbf{z}}$ belongs to one of the hyperplanes. Take any vector $\bar{\mathbf{y}} \in M^t$ and let us show that it belongs to a hyperplane in the family.

1. Any vector from M^t which differs from $\bar{\mathbf{z}}$ only in the first ℓ components belongs to a hyperplane of the family. This comes from Lemma 1.3 because when we change for $i \in \{1, \dots, \ell\}$ the i -th component of $\bar{\mathbf{z}}$ by $d_i = \frac{\lambda}{k}$, then we jump on the adjacent parallel hyperplane. So, any transformation of the i -th component of $\bar{\mathbf{z}}$ into another value from M means a finite number of jumps from one hyperplane onto another. Hence, we may transform $\bar{\mathbf{z}}$ so that it has the first ℓ components equal to $\bar{\mathbf{y}}$ and the obtained vector $\bar{\mathbf{x}}$ belongs to a hyperplane in the family.
2. Any vector from M^t which differs from $\bar{\mathbf{x}}$ only in the last $t - \ell$ components belongs to the same hyperplane as $\bar{\mathbf{x}}$. This comes from the orthogonality $\bar{\mathbf{e}}_i \perp \bar{\mathbf{n}}$ for $i > \ell$ (the argument is the same as in the proof of Lemma 1.3). Since $\bar{\mathbf{y}}$ differs from $\bar{\mathbf{x}}$ only in the last $t - \ell$ components, $\bar{\mathbf{y}}$ belongs to a hyperplane in the family. \square

2. Combinatorics on Words and the WELLD OC Property

2.1. Backgrounds on Combinatorics on Words

In the following, \mathcal{A} denotes a finite set of symbols called *letters*; the set \mathcal{A} is therefore called an *alphabet*. A *finite word* is a finite string $w = w_1 w_2 \dots w_n$ of letters from \mathcal{A} ; its length is denoted by $|w| = n$ and $|w|_a$ denotes the number of occurrences of $a \in \mathcal{A}$ in w . The empty word, a neutral element for concatenation of finite words, is denoted ε and it is of zero length. The set of all finite words over the alphabet \mathcal{A} is denoted by \mathcal{A}^* .

Under an *infinite word* we understand an infinite sequence $u = u_0 u_1 u_2 \dots$ of letters from \mathcal{A} . A finite word w is a *factor* of a word v (finite or infinite) if there exist words p and s such that $v = pws$. If $p = \varepsilon$, then w is said to be a *prefix* of v ; if $s = \varepsilon$, then w is a *suffix* of v . The set of factors and prefixes of v are denoted by $\text{Fact}(v)$ and $\text{Pref}(v)$, respectively. If $v = ps$ for finite words v, p, s , then we write $p = vs^{-1}$ and $s = p^{-1}v$.

An infinite word u over the alphabet \mathcal{A} is called *eventually periodic* if it is of the form $u = v w^\omega$, where v, w are finite words over \mathcal{A} and ω denotes an infinite repetition. An infinite word is called *aperiodic* if it is not eventually periodic.

For any factor w of an infinite word u , every index i such that w is a prefix of the infinite word $u_i u_{i+1} u_{i+2} \dots$ is called an *occurrence* of w in u . An infinite word u is *recurrent* if each of its factors has infinitely many occurrences in u .

The *factor complexity* of an infinite word u is a map $C_u : \mathbb{N} \mapsto \mathbb{N}$ such that $C_u(n)$ is the number of factors of length n contained in u . The factor complexity of eventually periodic words is bounded, while the factor complexity of an aperiodic word u satisfies $C_u(n) \geq n + 1$ for all $n \in \mathbb{N}$. A *right extension* of a factor w of u over the alphabet \mathcal{A} is any letter $a \in \mathcal{A}$ such that wa is a factor of u . Of course, any factor of u has at least one right extension. A factor w is called *right special* if w has at least two right extensions. Similarly, one can define a *left extension* and a *left special* factor. A factor is *bispecial* if it is both right and left special. An aperiodic word contains right special factors of any length.

The *Parikh vector* of a finite word w over an alphabet $\{0, 1, \dots, d-1\}$ is defined as $(|w|_0, |w|_1, \dots, |w|_{d-1})$. For a finite or infinite word $u = u_0 u_1 u_2 \dots$, $\text{Pref}_n u$ will denote the prefix of length n of u , i.e., $\text{Pref}_n u = u_0 u_1 \dots u_{n-1}$.

In some of the examples we consider are morphic words. A *morphism* is a function $\varphi : \mathcal{A}^* \rightarrow \mathcal{B}^*$ such that $\varphi(\varepsilon) = \varepsilon$ and $\varphi(wv) = \varphi(w)\varphi(v)$, for all $w, v \in \mathcal{A}^*$. Clearly, a morphism is completely defined by the images of the letters in the domain. A morphism is *prolongable* on $a \in \mathcal{A}$, if $|\varphi(a)| \geq 2$ and a is a prefix of $\varphi(a)$. If φ is prolongable on a , then $\varphi^n(a)$ is a proper prefix of $\varphi^{n+1}(a)$, for all $n \in \mathbb{N}$. Therefore, the sequence $(\varphi^n(a))_{n \geq 0}$ of words defines an infinite word u that is a fixed point of φ . Such a word u is a (pure) *morphic* word.

Let us introduce a combinatorial condition on infinite words that – as we will see later – guarantees no lattice structure for the associated PRNGs.

Definition 2.1 (The WELLDOC property). We say that an aperiodic infinite word u over the alphabet $\{0, 1, \dots, d-1\}$ has *well distributed occurrences* (or has *the WELLDOC property*) if for any $m \in \mathbb{N}$ and any factor w of u the word u satisfies the following condition. If i_0, i_1, \dots denote the positions of the occurrences of w in u , then

$$\{(|\text{Pref}_{i_j} u|_0, \dots, |\text{Pref}_{i_j} u|_{d-1}) \bmod m \mid j \in \mathbb{N}\} = \mathbb{Z}_m^d;$$

that is, the Parikh vectors of $\text{Pref}_{i_j} u$ for $j \in \mathbb{N}$, when reduced modulo m , give the whole set \mathbb{Z}_m^d .

We define the WELLDOC property for aperiodic words since it clearly never holds for periodic ones. It is easy to see that if a recurrent infinite word u has the WELLDOC property, then for every vector $\mathbf{v} \in \mathbb{Z}_m^d$ there are infinitely many values of j such that the Parikh vector of $\text{Pref}_{i_j} u$ is congruent to \mathbf{v} modulo m .

Example 2.2. The Thue-Morse word

$$u = 01101001100101101001011001101001 \dots,$$

which is a fixed point of the morphism $0 \mapsto 01, 1 \mapsto 10$, does not satisfy the WELLDOC property. Indeed, take $m = 2$ and $w = 00$, then w occurs only in odd positions i_j so that $(|\text{Pref}_{i_j} u|_0 + |\text{Pref}_{i_j} u|_1) = i_j$ is odd. Thus, e.g.,

$$(|\text{Pref}_{i_j} u|_0, |\text{Pref}_{i_j} u|_1) \bmod 2 \neq (0, 0),$$

and hence

$$\{(|\text{Pref}_{i_j} u|_0, |\text{Pref}_{i_j} u|_1) \bmod 2 \mid j \in \mathbb{N}\} \neq \mathbb{Z}_2^2.$$

Example 2.3. We say that an infinite word u over an alphabet \mathcal{A} , $|\mathcal{A}| = d$, is *universal* if it contains all finite words over \mathcal{A} as its factors. It is easy to see that any universal word satisfies the WELLDOC property. Indeed, for any word $w \in \mathcal{A}^*$ and any m there exists a finite word v such that if i_0, i_1, \dots, i_k denote the occurrences of w in v , then

$$\{(|\text{Pref}_{i_j} v|_0, \dots, |\text{Pref}_{i_j} v|_{d-1}) \bmod m \mid j \in \{0, 1, \dots, k\}\} = \mathbb{Z}_m^d.$$

Since u is universal, v is a factor of u . Denoting by i an occurrence of v in u , one gets that the positions $i + i_j$ are occurrences of w in u . Hence

$$\begin{aligned} & \{(|\text{Pref}_{i+i_j} u|_0, \dots, |\text{Pref}_{i+i_j} u|_{d-1}) \bmod m \mid j \in \{0, 1, \dots, k\}\} \\ &= \{(|\text{Pref}_i u|_0, \dots, |\text{Pref}_i u|_{d-1}) + (|\text{Pref}_{i_j} v|_0, \dots, |\text{Pref}_{i_j} v|_{d-1}) \bmod m \mid j \in \{0, 1, \dots, k\}\} = \mathbb{Z}_m^d. \end{aligned}$$

Therefore, u satisfies the WELLDOC property.

2.2. Combination of PRNGs

In order to eliminate the lattice structure, it helps to combine PRNGs in a smart way. Such a method was introduced in [13]. Let $X = (X_n)_{n \in \mathbb{N}}$ and $Y = (Y_n)_{n \in \mathbb{N}}$ be two PRNGs with the same output $M \subset \mathbb{N}$ and the same period $m \in \mathbb{N}$, and let $u = u_0u_1u_2 \dots$ be a binary infinite word over the alphabet $\{0, 1\}$.

The PRNG $Z = (Z_n)_{n \in \mathbb{N}}$ based on u is obtained by the following algorithm:

1. Read step by step the letters of u .
2. When you read 0 for the i -th time, copy the i -th symbol from X to the end of the constructed sequence Z .
3. When you read 1 for the i -th time, copy the i -th symbol from Y to the end of the constructed sequence Z .

This construction can be generalized for non-binary alphabets: Using infinite words over a multiliteral alphabet, one can combine more than two PRNGs. We remark that following terminology from [3], the sequence Z is obtained as a *shuffle* of the sequences X and Y with the steering word u .

In order to distinguish between generators and infinite words used for their combination, we always use capital letters X, Y, Z, \dots for generators and lower-case letters u, v, w for words (the same convention is applied for their outputs: A, B, \dots for output values of generators (elements of M), a, b, \dots for letters of words). Finite sequences of successive elements $\bar{x} = (X_i, X_{i+1}, \dots, X_{i+t-1})$ of a PRNG X are called t -tuples, or vectors, while in the case of an infinite word u , we call $u_iu_{i+1} \dots u_{i+t-1}$ a factor of length t .

2.3. The WELLDOC Property and Absence of the Lattice Structure

Guimond et al. in [14] have shown that PRNGs based on infinite words coding a certain class of cut-and-project sets have no lattice structure. In the sequel, we will generalize their result and find larger classes of words guaranteeing no lattice structure for associated generators. We focus on the binary alphabet, although everything works for multiliteral words as well (and for combination of more generators therefore), since the proofs become more technical in non-binary case.

Theorem 2.4. *Let $X = (X_n)_{n \in \mathbb{N}}$ and $Y = (Y_n)_{n \in \mathbb{N}}$ be PRNGs with output $M \subseteq \mathbb{N}$ and period $m > 1$. Let Z be the PRNG obtained from X and Y , based on a binary infinite word u with the WELLDOC property. Then Z has no lattice structure.*

Proof. According to Proposition 1.1, it suffices to check that its assumptions are met. Let $A, B \in M$ and $\ell \in \mathbb{N}$. Assume $A = X_i$ and $B = Y_j$ for some i, j . Consider a right special factor w of u of length ℓ , i.e., both words $w0$ and $w1$ are factors of u (such a factor w exists since u is an aperiodic word because of the WELLDOC property). By Definition 2.1, it is possible to find an occurrence i_k of $w0$ in u such that

$$|\text{Pref}_{i_k} u|_0 \equiv i - |w|_0 - 1 \pmod{m}, \quad |\text{Pref}_{i_k} u|_1 \equiv j - |w|_1 - 1 \pmod{m}.$$

Reading the word $w0$ at the occurrence i_k , the corresponding ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ of the generator Z consists of symbols

$$X_{(i-|w|_0) \bmod m}, \dots, X_{(i-1) \bmod m} \text{ and } Y_{(j-|w|_1) \bmod m}, \dots, Y_{(j-1) \bmod m}.$$

When reading 0 after w , the symbol $X_i = A$ from the first generator follows $(A_1, A_2, \dots, A_\ell)$.

Again, by Definition 2.1, there exists an occurrence i_s of $w1$ in u such that

$$|\text{Pref}_{i_s} u|_0 \equiv i - |w|_0 - 1 \pmod{m}, \quad |\text{Pref}_{i_s} u|_1 \equiv j - |w|_1 - 1 \pmod{m}.$$

When reading the word w at the occurrence i_s , the same ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ of Z as previously occurs. This time, however, $(A_1, A_2, \dots, A_\ell)$ is followed by B because we read $w1$ and $Y_j = B$. Thus, we have found an ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ of Z followed in Z by both A and B . \square

Remark 2.5. The WELLDOC property is sufficient, but not necessary for absence of the lattice structure. For example, consider a modified Fibonacci word \hat{u} where the letter 2 is inserted after each letter, i.e., $\hat{u} = 0212020212021202 \dots$. It is easy to verify that \hat{u} does not have well distributed occurrences. However, we will show the following: Let Z be the PRNG combining three generators $X = (X_n)_{n \in \mathbb{N}}$, $Y = (Y_n)_{n \in \mathbb{N}}$ and $V = (V_n)_{n \in \mathbb{N}}$ with the same output $M \subset \mathbb{N}$ and the same period $m \in \mathbb{N}$ according to the modified Fibonacci word \hat{u} . Then Z has no lattice structure.

It suffices to verify assumptions of Proposition 1.1. Let $A, B \in M$ and $\ell \in \mathbb{N}$, ℓ an even number (the proof is analogous for odd ℓ). Assume $A = X_i$ and $B = Y_j$. Consider a right special factor w of the Fibonacci word u of length $\ell/2$. Since u has the WELLDOC property, there exists an occurrence i_k of $w0$ in u such that

$$|\text{Pref}_{i_k} u|_0 \equiv i - |w|_0 - 1 \pmod{m}, \quad |\text{Pref}_{i_k} u|_1 \equiv j - |w|_1 - 1 \pmod{m}.$$

Then if we insert the letter 2 after each letter of w , we obtain a right special factor \hat{w} of the modified Fibonacci word \hat{u} of length ℓ . It holds then that

$$\begin{aligned} |\text{Pref}_{2i_k} \hat{u}|_0 &\equiv i - |w|_0 - 1 \equiv i - |\hat{w}|_0 - 1 \pmod{m}, \\ |\text{Pref}_{2i_k} \hat{u}|_1 &\equiv j - |w|_1 - 1 \equiv j - |\hat{w}|_1 - 1 \pmod{m}, \\ |\text{Pref}_{2i_k} \hat{u}|_2 &\equiv i - |w|_0 - 1 + j - |w|_1 - 1 \equiv i + j - |\hat{w}|_2 - 2 \pmod{m}. \end{aligned}$$

When reading the word $\hat{w}0$ at the occurrence $2i_k$, the corresponding ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ of the generator Z is followed by the symbol $X_i = A$ from the first generator.

Again, by the WELLDOC property of u , there exists an occurrence i_s of $w1$ in u such that

$$|\text{Pref}_{i_s} u|_0 \equiv i - |w|_0 - 1 \pmod{m}, \quad |\text{Pref}_{i_s} u|_1 \equiv j - |w|_1 - 1 \pmod{m}.$$

It holds then that

$$\begin{aligned} |\text{Pref}_{2i_s} \hat{u}|_0 &\equiv i - |w|_0 - 1 \equiv i - |\hat{w}|_0 - 1 \pmod{m}, \\ |\text{Pref}_{2i_s} \hat{u}|_1 &\equiv j - |w|_1 - 1 \equiv j - |\hat{w}|_1 - 1 \pmod{m}, \\ |\text{Pref}_{2i_s} \hat{u}|_2 &\equiv i - |w|_0 - 1 + j - |w|_1 - 1 \equiv i + j - |\hat{w}|_2 - 2 \pmod{m}. \end{aligned}$$

When reading the word \hat{w} at the occurrence $2i_s$, the same ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ of Z as previously occurs. This time, however, $(A_1, A_2, \dots, A_\ell)$ is followed by B because we read $\hat{w}1$ and $Y_j = B$. Thus, we have found an ℓ -tuple $(A_1, A_2, \dots, A_\ell)$ of Z followed in Z by both A and B . Therefore Z has no lattice structure.

Remark 2.6. In the proof of Theorem 2.4, the modulus m from the WELLDOC property is set to be equal to the period of the combined generators. Therefore, if we require absence of the lattice structure for a PRNG obtained when combining PRNGs with a fixed period \hat{m} , then it is sufficient to use an infinite word u that satisfies the WELLDOC property for the modulus $m = \hat{m}$. This means for instance that the Thue-Morse word is not completely out of the game, but it cannot be used to combine periodic PRNGs with the period being a power of 2.

We have formulated a combinatorial condition – well distributed occurrences – guaranteeing no lattice structure of the associated generator. It is now important to find classes of words satisfying such a condition.

3. Sturmian Words

In this section we show that Sturmian words have well distributed occurrences.

Definition 3.1. An aperiodic infinite word u is called *Sturmian* if its factor complexity satisfies $C_u(n) = n + 1$ for all $n \in \mathbb{N}$.

So, Sturmian words are by definition binary and they have the lowest possible factor complexity among aperiodic infinite words. Sturmian words admit various types of characterizations of geometric and combinatorial nature. One of such characterizations is via irrational rotations on the unit circle. In [19] Hedlund and Morse showed that each Sturmian word may be realized measure-theoretically by an irrational rotation on the circle. That is, every Sturmian word is obtained by coding the symbolic orbit of a point on the circle of circumference one under a rotation R_α by an irrational angle¹ α , $0 < \alpha < 1$, where the circle is partitioned into two complementary intervals, one of length α and the other of length $1 - \alpha$. Conversely, each such coding gives rise to a Sturmian word.

Definition 3.2. The *rotation* by angle α is the mapping R_α from $[0, 1)$ (identified with the unit circle) to itself defined by $R_\alpha(x) = \{x + \alpha\}$, where $\{x\} = x - \lfloor x \rfloor$ is the fractional part of x . Considering a partition of $[0, 1)$ into $I_0 = [0, 1 - \alpha)$, $I_1 = [1 - \alpha, 1)$, define a word

$$s_{\alpha,\rho}(n) = \begin{cases} 0 & \text{if } R_\alpha^n(\rho) = \{\rho + n\alpha\} \in I_0, \\ 1 & \text{if } R_\alpha^n(\rho) = \{\rho + n\alpha\} \in I_1. \end{cases}$$

One can also define $I'_0 = (0, 1 - \alpha]$, $I'_1 = (1 - \alpha, 1]$, the corresponding word is denoted by $s'_{\alpha,\rho}$.

Remark that some but not all Sturmian words are morphic. In fact, it is known that a characteristic Sturmian word (i.e., $\rho = \alpha$) is morphic if and only if the continuous fraction expansion of α is periodic. For more information on Sturmian words we refer to [16, Chapter 2].

Theorem 3.3. *Let u be a Sturmian word. Then u has the WELLDOC property.*

Proof. In the proof we use the definition of Sturmian word via rotation. The main idea is controlling the number of 1's modulo m by taking circle of length m , and controlling the length taking the rotation by $m\alpha$.

For the proof we will use an equivalent reformulation of the theorem:

Let u be a Sturmian word on $\{0, 1\}$, for any natural number m and any factor w of u let us denote i_0, i_1, \dots the occurrences of w in u . Then

$$\{(i_j, |\text{Pref}_{i_j} u|_1) \bmod m \mid j \in \mathbb{N}\} = \mathbb{Z}_m^2.$$

That is, we control the number of 1's and the length instead of the number of 0's.

Since a Sturmian word can be defined via rotations by an irrational angle on a unit circle, without loss of generality we may assume that $u = s_{\alpha,\rho}$ for some $0 < \alpha < 1$, $0 \leq \rho < 1$, α irrational (see Definition 3.2). Equivalently, we can consider m copies of the circle connected into one circle of length m with m intervals I_1^i of length α corresponding to 1. The Sturmian word is obtained by rotation by α on this circle of length m (see Fig. 2).

¹Measured by arc length (thus equivalent to $2\pi\alpha$ radians).

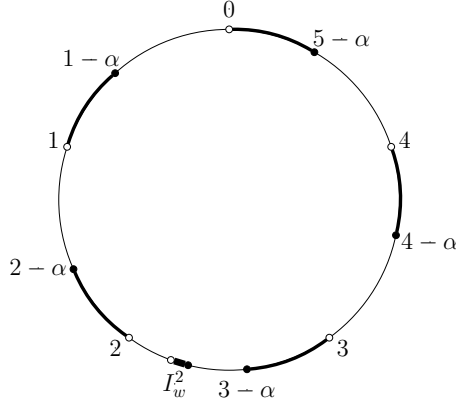


Figure 2: Illustration to the proof of Theorem 3.3: the example for $m = 5$.

Namely, we define the rotation $R_{\alpha,m}$ as the mapping from $[0, m)$ (identified with the circle of length m) to itself defined by $R_{\alpha,m}(x) = \{x + \alpha\}_m$, where $\{x\}_m = x - \lfloor x/m \rfloor m$ and for $m = 1$ coincides with the fractional part of x . A partition of $[0, m)$ into $2m$ intervals $I_0^i = [i, i + 1 - \alpha)$, $I_1^i = [i + 1 - \alpha, i + 1)$, $i = 0, \dots, m - 1$ defines the Sturmian word $u = s_{\alpha,\rho}$:

$$s_{\alpha,\rho}(n) = \begin{cases} 0 & \text{if } R_{\alpha,m}^n(\rho) = \{\rho + n\alpha\}_m \in I_0^i \text{ for some } i = 0, \dots, m - 1, \\ 1 & \text{if } R_{\alpha,m}^n(\rho) = \{\rho + n\alpha\}_m \in I_1^i \text{ for some } i = 0, \dots, m - 1. \end{cases}$$

It is well known that any factor $w = w_0 \cdots w_{k-1}$ of u corresponds to an interval I_w in $[0, 1)$, so that whenever you start rotating from the interval I_w , you obtain w . Namely, $x \in I_w$ if and only if $x \in I_{w_0}, R_\alpha(x) \in I_{w_1}, \dots, R_\alpha^{|w|-1}(x) \in I_{w_{|w|-1}}$.

Similarly, we can define m intervals corresponding to w in $[0, m)$ (circle of length m), so that if $I_w = [x_1, x_2)$, then $I_w^i = [x_1 + i, x_2 + i)$, $i = 0, \dots, m - 1$.

Fix a factor w of u , take arbitrary $(j, i) \in \mathbb{Z}_m^2$. Now let us organize (j, i) among the occurrences of w , i.e., find l such that $u_l \dots u_{l+|w|-1} = w$, $l \bmod m = j$ and $|\text{Pref}_l u|_1 \bmod m = i$:

Consider rotation $R_{m\alpha,m}(x)$ by $m\alpha$ instead of rotation by α , and start m -rotating from $j\alpha + \rho$. Formally, $R_{m\alpha,m}(x) = \{x + m\alpha\}_m$, where, as above, $\{x\}_m = x - \lfloor x/m \rfloor m$. This rotation will put us to positions $mk + j$, $k \in \mathbb{N}$, in the Sturmian word: for $a \in \{0, 1\}$ one has $s_{\alpha,\rho}(mk + j) = a$ if $R_{m\alpha,m}^k(j\alpha + \rho) = \{j\alpha + \rho + km\alpha\}_m \in I_a^i$ for some $i = 0, \dots, m - 1$.

Remark that the points in the orbit of an m -rotation of a point on the m -circle are dense, and hence the rotation comes infinitely often to each interval. So pick k when $j\alpha + mk\alpha + \rho \in I_w^i \subset [i, i + 1)$ (and actually there exist infinitely many such k). Then the length l of the corresponding prefix is equal to $km + j$, and the number of 1's in it is $i + mp$, where p is the number of complete circles you made, i.e., $p = \lfloor (j\alpha + mk\alpha + \rho)/m \rfloor$. \square

4. Arnoux-Rauzy Words

In this section we show that Arnoux-Rauzy words [1], which are natural extensions of Sturmian words to larger alphabets, also satisfy the WELLDOC property. Note that the proof for Sturmian words cannot be generalized to Arnoux-Rauzy words, because it is based on the geometric interpretation of Sturmian words via rotations, while this interpretation does not extend to Arnoux-Rauzy words.

4.1. Basic Definitions

The definitions and results we remind in this subsection are well-known and mostly taken from [1, 8] and generalize the ones given for binary words in [4].

Definition 4.1. Let \mathcal{A} be a finite alphabet. The *reversal operator* is the operator $\sim: \mathcal{A}^* \mapsto \mathcal{A}^*$ defined by recurrence in the following way:

$$\tilde{\varepsilon} = \varepsilon, \quad \widetilde{va} = a\tilde{v}$$

for all $v \in \mathcal{A}^*$ and $a \in \mathcal{A}$. The fixed points of the reversal operator are called *palindromes*.

Definition 4.2. Let $v \in \mathcal{A}^*$ be a finite word over the alphabet \mathcal{A} . The *right palindromic closure* of v , denoted by $v^{(+)}$, is the shortest palindrome that has v as a prefix. It is readily verified that if p is the longest palindromic suffix of $v = wp$, then $v^{(+)} = wp\tilde{w}$.

Definition 4.3. We call the *iterated (right) palindromic closure operator* the operator ψ recurrently defined by the following rules:

$$\psi(\varepsilon) = \varepsilon, \quad \psi(va) = (\psi(v)a)^{+}$$

for all $v \in \mathcal{A}^*$ and $a \in \mathcal{A}$. The definition of ψ may be extended to infinite words u over \mathcal{A} as $\psi(u) = \lim_n \psi(\text{Pref}_n u)$, i.e., $\psi(u)$ is the infinite word having $\psi(\text{Pref}_n u)$ as its prefix for every $n \in \mathbb{N}$.

Definition 4.4. Let Δ be an infinite word on the alphabet \mathcal{A} such that every letter occurs infinitely often in Δ . The word $c = \psi(\Delta)$ is then called a *characteristic (or standard) Arnoux-Rauzy word* and Δ is called the *directive sequence* of c . An infinite word u is called an *Arnoux-Rauzy word* if it has the same set of factors as a (unique) characteristic Arnoux-Rauzy word, which is called the *characteristic word* of u . The *directive sequence* of an Arnoux-Rauzy word is the directive sequence of its characteristic word.

Let us also recall the following well-known characterization (see e.g. [8]):

Theorem 4.5. *Let u be an aperiodic infinite word over the alphabet \mathcal{A} . Then u is a standard Arnoux-Rauzy word if and only if the following hold:*

1. *Fact(u) is closed under reversal (that is, if v is a factor of u so is \tilde{v}).*
2. *Every left special factor of u is also a prefix.*
3. *If v is a right special factor of u then va is a factor of u for every $a \in \mathcal{A}$.*

From the preceding theorem, it can be easily verified that the bispecial factors of a standard Arnoux-Rauzy correspond to its palindromic prefixes (including the empty word), and hence to the iterated palindromic closure of the prefixes of its directive sequence. That is, if

$$\varepsilon = b_0, b_1, b_2, \dots$$

is the sequence, ordered by length, of bispecial factors of the standard Arnoux-Rauzy word u , $\Delta = \Delta_0 \Delta_1 \dots$ its directive sequence (with $\Delta_i \in \mathcal{A}$ for every i), we have $b_{i+1} = (b_i \Delta_i)^{+}$.

A direct consequence of this, together with the preceding definitions, is the following statement, which will be used in the sequel.

Lemma 4.6. *Let u be a characteristic Arnoux-Rauzy word and let Δ and $(b_i)_{i \geq 0}$ be defined as above. If Δ_i does not occur in b_i , then $b_{i+1} = b_i \Delta_i b_i$. Otherwise let $j < i$ be the largest integer such that $\Delta_j = \Delta_i$. Then $b_{i+1} = b_i b_j^{-1} b_i$.*

4.2. Parikh Vectors and Arnoux-Rauzy Factors

Where no confusion arises, given an Arnoux-Rauzy word u , we will let

$$\varepsilon = b_0, b_1, \dots, b_n, \dots$$

denote the sequence of bispecial factors of u ordered by length, and for any $i \in \mathbb{N}$ we will let $\bar{\mathbf{b}}_i$ denote the Parikh vector of b_i .

Remark 4.7. By the pigeonhole principle, it is clear that for every $m \in \mathbb{N}$ there exists an integer $N \in \mathbb{N}$ such that, for every $i \geq N$, the set $\{j > i \mid \bar{\mathbf{b}}_j \equiv_m \bar{\mathbf{b}}_i\}$ is infinite. Where no confusion arises and with a slight abuse of notation, fixed m , we will always write N for the smallest of such integers.

Lemma 4.8. *Let u be a characteristic Arnoux-Rauzy word and let $m \in \mathbb{N}$. Let*

$$\alpha_1 \bar{\mathbf{b}}_{j_1} + \dots + \alpha_k \bar{\mathbf{b}}_{j_k} \equiv_m \bar{\mathbf{v}} \in \mathbb{Z}_m^d$$

be a linear combination of Parikh vectors such that $\sum_{i=1}^k \alpha_i = 0$, with $j_i \geq N$ and $\alpha_i \in \mathbb{Z}$ for all $i \in \{1, \dots, k\}$. Then, for any $\ell \in \mathbb{N}$, there exists a prefix v of u such that the Parikh vector of v is congruent to $\bar{\mathbf{v}}$ modulo m and vb_ℓ is also a prefix of u .

Proof. Without loss of generality, we can assume $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_k$, hence there exists k' such that

$$\alpha_1 \geq \alpha_{k'} \geq 0 \geq \alpha_{k'+1} \geq \alpha_k.$$

We will prove the result by induction on $\beta = \sum_{j=1}^{k'} \alpha_j$. If $\beta = 0$, trivially, we can take $v = \varepsilon$ and the statement is clearly verified. Let us assume the statement true for all $0 \leq \beta < n$ and let us prove it for $\beta = n$. By the remark preceding this lemma, for every ℓ we can choose $i' > j' > \ell$ such that $\bar{\mathbf{b}}_{j_1} \equiv_m \bar{\mathbf{b}}_{i'}$ and $\bar{\mathbf{b}}_{j_k} \equiv_m \bar{\mathbf{b}}_{j'}$. Since every bispecial factor is a prefix and suffix of all the bigger ones, in particular we have that $b_{j'}$ is a suffix of $b_{i'}$, and b_ℓ is a prefix of $b_{j'}$; this implies that $b_{i'} b_{j'}^{-1} b_\ell$ is actually a prefix of $b_{i'}$. By assumption, the Parikh vector of $b_{i'} b_{j'}^{-1} b_\ell$ is clearly $\bar{\mathbf{b}}_{i'} - \bar{\mathbf{b}}_{j'} \equiv_m \bar{\mathbf{b}}_{j_1} - \bar{\mathbf{b}}_{j_k}$. Since $\alpha_1 \geq 1$ implies $\alpha_k \leq -1$, we have, by induction hypothesis, that there exists a prefix w of u such that the Parikh vector of w is congruent modulo m to

$$(\alpha_1 - 1) \bar{\mathbf{b}}_{j_1} + \dots + (\alpha_k + 1) \bar{\mathbf{b}}_{j_k}$$

and $w b_{i'}$ is a prefix of u . Hence $w b_{i'} b_{j'}^{-1} b_\ell$ is also a prefix of u and, by simple computation, the Parikh vector of $v = w b_{i'} b_{j'}^{-1} b_\ell$ is congruent modulo m to $\bar{\mathbf{v}} = \alpha_1 \bar{\mathbf{b}}_{j_1} + \dots + \alpha_k \bar{\mathbf{b}}_{j_k}$. \square

Definition 4.9. Let $n \in \mathbb{Z}$. We will say that an integer linear combination of integer vectors is a n -combination if the sum of all the coefficients equals n .

Lemma 4.10. *Let u be a characteristic Arnoux-Rauzy word and let $n \in \mathbb{N}$. Every n -combination of Parikh vectors of bispecial factors can be expressed as an n -combination of Parikh vectors of arbitrarily large bispecials. In particular, for every $K, L \in \mathbb{N}$, it is possible to find a finite number of integers $\alpha_1, \dots, \alpha_k$ such that $\bar{\mathbf{b}}_K = \alpha_1 \bar{\mathbf{b}}_{j_1} + \dots + \alpha_k \bar{\mathbf{b}}_{j_k}$ with $j_i > L$ for every i and $\alpha_1 + \dots + \alpha_k = 1$.*

Proof. A direct consequence of Lemma 4.6 is that for every i such that Δ_i appears in b_i , we have $\bar{\mathbf{b}}_{i+1} = 2\bar{\mathbf{b}}_i - \bar{\mathbf{b}}_j$, where $j < i$ is the largest such that $\Delta_j = \Delta_i$. This in turn (since every letter in Δ appears infinitely many times from the definition of Arnoux-Rauzy word) implies that for every non-negative integer j , there exists a positive k such that $\bar{\mathbf{b}}_j = 2\bar{\mathbf{b}}_{j+k} - \bar{\mathbf{b}}_{j+k+1}$, that is, we can substitute each Parikh vector of a bispecial with a 1-combination of Parikh vectors of strictly larger bispecials. Simply iterating the process, we obtain the statement. \square

In the following we will assume the set \mathcal{A} to be a finite alphabet of cardinality d . For every set $X \subseteq \mathcal{A}^*$ of finite words, we will let $\text{PV}(X) \subseteq \mathbb{Z}^d$ denote the set of Parikh vectors of elements of X and for every $m \in \mathbb{N}$ we will let $\text{PV}_m(X) \subseteq \mathbb{Z}_m^d$ denote the set of elements of $\text{PV}(X)$ reduced modulo m .

For an infinite word u over \mathcal{A} , and a factor v of u , let $S_v(u)$ denote the set of all prefixes of u followed by an occurrence of v . In other words,

$$S_v(u) = \{p \in \text{Pref}(u) \mid pv \in \text{Pref}(u)\}.$$

Definition 4.11. For any set of finite words $X \subseteq \mathcal{A}^*$, we will say that u has the property \mathcal{P}_X (or, for short, that u has \mathcal{P}_X) if, for every $m \in \mathbb{N}$ and for every $v \in X$ we have that

$$\text{PV}_m(S_v(u)) = \mathbb{Z}_m^d.$$

That is to say, for every vector $\bar{w} \in \mathbb{Z}_m^d$ there exists a word $w \in S_v(u)$ such that the Parikh vector of w is congruent to \bar{w} modulo m .

With this notation, an infinite word u has the WELLDOC property if and only if it has the property $\mathcal{P}_{\text{Fact}(u)}$.

Proposition 4.12. *Let u be a characteristic Arnoux-Rauzy word over the d -letter alphabet \mathcal{A} . Then u has the property $\mathcal{P}_{\text{Pref}(u)}$.*

Proof. Let us fix an arbitrary $m \in \mathbb{N}$. We want to show that, for every $v \in \text{Pref}(u)$, $\text{PV}_m(S_v(u)) = \mathbb{Z}_m^d$. Let then $\bar{v} \in \mathbb{Z}^d$ and ℓ be the smallest number such that v is a prefix of b_ℓ . Let $i_1 < i_2 < \dots < i_d$ be such that Δ_{i_j} does not appear in b_{i_j} , where Δ is the directive word of u . Without loss of generality, we can rearrange the letters so that each Δ_{i_j} is lexicographically smaller than $\Delta_{i_{j+1}}$. With this assumption if, for every j , we set $\bar{v}_j = \bar{\mathbf{b}}_{i_{j+1}}$, i.e., equal to the Parikh vector of $b_{i_{j+1}}$, which, by the first part of Lemma 4.6, equals $b_{i_j} \Delta_{i_j} b_{i_j}$, we can find $j-1$ positive integers μ_1, \dots, μ_{j-1} such that $\bar{v}_j = (\mu_1, \mu_2, \dots, \mu_{j-1}, 1, 0, \dots, 0)$. It is easy to show, then, that the set $V = \{\bar{v}_1, \dots, \bar{v}_d\}$ generates \mathbb{Z}^d , hence there exists an integer n such that \bar{v} can be expressed as an n -combination of elements of V (which are Parikh vectors of bispecial factors of u). Trivially, then, $\bar{v} = \bar{v} - n\bar{\mathbf{0}} = \bar{v} - n\bar{\mathbf{b}}_0$; thus, it is possible to express \bar{v} as a 0-combination of Parikh vectors of (by the previous Lemma 4.10) arbitrarily large bispecial factors of u . By Lemma 4.8, then there exists a prefix p of u whose Parikh vector $\bar{\mathbf{p}}$ satisfies $\bar{\mathbf{p}} \equiv_m \bar{v}$ and pb_ℓ is a prefix of u . Since we picked ℓ such that v is a prefix of b_ℓ , we have that $p \in S_v(u)$. From the arbitrariness of v , \bar{v} and m , we obtain the statement. \square

As a corollary of Proposition 4.12, we obtain the main result of this section.

Theorem 4.13. *Let u be an Arnoux-Rauzy word over the d -letter alphabet \mathcal{A} . Then u has the property $\mathcal{P}_{\text{Fact}(u)}$, or equivalently, u has the WELLDOC property.*

Proof. Let m be a positive integer and let c be the characteristic word of u . Let v be a factor of u and xvy be the shortest bispecial containing v . By Proposition 4.12, we have that $\text{PV}_m(S_{xv}(c)) = \mathbb{Z}_m^d$ and, since the set is finite, we can find a prefix p of c such that $\text{PV}_m(S_{xv}(p)) = \mathbb{Z}_m^d$. Let w be a prefix of u such that wp is a prefix of u . If \bar{x} and \bar{w} are the Parikh vectors of, respectively, x and w , it is easy to see that

$$\bar{w} + \bar{x} + \text{PV}(S_{xv}(p)) \subseteq \bar{w} + \text{PV}(S_v(p)) \subseteq \text{PV}(S_v(u))$$

Since we have chosen p such that $\text{PV}_m(S_{xv}(p)) = \mathbb{Z}_m^d$, we clearly obtain that $\text{PV}_m(S_v(u)) = \mathbb{Z}_m^d$ and hence, by the arbitrariness of v and m , the statement. \square

Remark 4.14. Now we introduce a simple method of obtaining words satisfying the WELLDOC property. Take a word u with the WELLDOC property over an alphabet $\{0, 1, \dots, d-1\}$, $d > 2$, apply a morphism $\varphi : d-1 \mapsto 0, i \mapsto i$ for $i = 0, \dots, d-2$, i.e., φ joins two letters into one. It is straightforward that $\varphi(u)$ has the WELLDOC property. So, taking Arnoux-Rauzy words and joining some letters, we obtain other words than Sturmian and Arnoux-Rauzy satisfying the WELLDOC property.

Remark 4.15. Now we introduce another class of morphisms preserving the WELLDOC property. Recall that the *adjacency matrix* Φ of a morphism $\varphi : \mathcal{A} \rightarrow \mathcal{A}$, with $\mathcal{A} = \{0, 1, \dots, d-1\}$, is defined by $\Phi_{i,j} = |\varphi(j-1)|_{i-1}$ for $1 \leq i, j \leq d$. By definition, it follows that if \bar{v} is the Parikh vector of $v \in \mathcal{A}^*$, then $\Phi\bar{v}$ is the Parikh vector of $\varphi(v)$.

Let us show that if $\det \Phi = \pm 1$ and u has the WELLDOC property, then so does $\varphi(u)$. Indeed, let w be any factor of $\varphi(u)$, and suppose $xwy = \varphi(v)$ for some $v \in \text{Fact}(u)$ and $x, y \in \mathcal{A}^*$. We then have $S_w(\varphi(u)) \supseteq \varphi(S_v(u))x$, so that, writing \bar{x} for the Parikh vector of x , we have for any $m > 0$

$$\text{PV}_m(S_w(\varphi(u))) \supseteq \Phi \cdot \text{PV}_m(S_v(u)) + \bar{x} \pmod{m}.$$

Since u has the WELLDOC property, $\text{PV}_m(S_v(u)) = \mathbb{Z}_m^d$. As $\det \Phi = \pm 1$, Φ is invertible (even modulo m), so that $\Phi \cdot \mathbb{Z}_m^d + \bar{x} \pmod{m} = \mathbb{Z}_m^d$. Hence $\text{PV}_m(S_w(\varphi(u))) = \mathbb{Z}_m^d$, showing that $\varphi(u)$ has the WELLDOC property by the arbitrariness of w and m .

5. Statistical Tests of PRNGs

In the previous part, we showed that PRNGs based on infinite words with well distributed occurrences have no lattice structure. In this section we perform empirical statistical tests. We will show how mixing based on aperiodic infinite words will cope with known weaknesses of LCGs: the statistical test we perform show significant improvements.

5.1. Computer Generation of Morphic Words

From practical point of view it is important to find algorithms for generating infinite words (more precisely, their prefixes) that are efficient both in memory footprint and CPU time. In [20] an efficient algorithm for generating the Fibonacci word was introduced: The prefix of length n is generated in $O(\log(n))$ space and $O(n)$ time. We generalize this method for any Sturmian and Arnoux-Rauzy word being a fixed point of a morphism φ . The main ingredient is that we consider φ^n instead of φ ; we precompute and store in the memory $\varphi^n(a)$ for any $a \in \mathcal{A}$. The runtime and memory consumption to generate 10^{10} letters of the Fibonacci and the Tribonacci word is summarized in Table 1. There are the following observations we would like to point out:

1. There is no need to store the first n letters in memory to generate the $(n+1)$ -th letter. Letters are generated on the fly and only nodes of the traversal tree are kept in the memory. The algorithm also supports leap frogging, generation can be started at any position in the word. The consequence is that the algorithm can be easily parallelized to produce multiple streams [10].
2. Using the method from [20] together with our improvement for generation of Sturmian and Arnoux-Rauzy words, the speed of generation of their prefixes is much higher than the speed of generation of LCGs output values. For example, generation of 10^{10} 32-bit values using a LCG modulo 2^{64} takes 14.3 seconds on our machine. Compare it to 0.5 seconds for generation of 10^{10} letters of a fixed point of a morphism with the same hardware. Thus, using a fixed point to combine LCGs causes only a negligible runtime penalty.

- The speed of generation can be further improved by using a higher initial memory footprint and CPU that can effectively copy such larger chunks of memory (size of L1 data cache is a limiting factor). Thus the new method scales nicely and can benefit from the future CPUs with higher L1 caches. The only requirement is to precompute $\varphi^n(a)$, $a \in \mathcal{A}$, for larger n . Our program does this automatically based on the limit on the initial memory consumption provided by the user.

Word	Fibonacci	Tribonacci
φ morphism rule	115s / 336 Bytes	107s / 256 Bytes
φ^n morphism rule	0.41s / 32 Bytes	0.36s / 32 Bytes

Table 1: The comparison of time in seconds and memory consumption to hold the traversal tree state needed to generate the first 10^{10} letters of the Fibonacci and the Tribonacci word using the original [20] (1st line) and the new algorithm (2nd line). The iteration n in the φ^n rule was chosen so that the length of $\varphi^n(a)$ does not exceed 4096 bytes for any $a \in \mathcal{A}$. The measurement was done on Intel Core i7-3520M CPU running at 2.90GHz.

5.2. Testing PRNGs Based on Sturmian and Arnoux-Rauzy words

We will present results for PRNGs based on:

- the Fibonacci word (as an example of a Sturmian word), i.e., the fixed point of the morphism $0 \mapsto 01, 1 \mapsto 0$,
- the modified Fibonacci word – Fibonacci2 – with the letter 2 inserted after each letter (see Remark 2.5),
- the Tribonacci word (as the simplest example of a ternary Arnoux-Rauzy word), i.e., the fixed point of $0 \mapsto 01, 1 \mapsto 02, 2 \mapsto 0$.
- the class of pure morphic Arnoux-Rauzy words where the morphism is a combination of morphisms $\sigma_0 = \{0 \mapsto 0, 1 \mapsto 10, 2 \mapsto 20\}, \sigma_1 = \{0 \mapsto 01, 1 \mapsto 1, 2 \mapsto 21\}, \sigma_2 = \{0 \mapsto 02, 1 \mapsto 12, 2 \mapsto 2\}$ and each morphism $\sigma_{0..2}$ is used at least once. We developed an algorithm to generate any word from this class. We tested several different generators based on Arnoux-Rauzy words from this class but since the results were very similar we will present here only results for the word $AR(\sigma_0, \sigma_1, \sigma_2, \sigma_1)$. The resulting morphism is $\{0 \mapsto 0102010102010, 1 \mapsto 102010, 2 \mapsto 2010102010\}$. We chose this particular word because its definition is still relatively simple (it is a combination of only 4 σ rules) but at the same time it differs significantly from Tribonacci word (for example, the frequencies of letters are different).

All programs are available online, together with a description [15]. We included the modified Fibonacci word that does not have the WELLDOC property, but guarantees no lattice structure for the arising generator. However, this word gives worse results in testing than the Fibonacci word (which is expectable since the resulting PRNG has an LCG on even positions).

5.2.1. Combining LCGs

Instead of combining plain LCGs, we will apply some modifications before their combination. Those modifications turn out to be useful according to the known weaknesses of LCGs.

We chose LCGs with the period m in range from $2^{47} - 115$ to 2^{64} , but we use only their upper 32 bits as the output since the statistical tests require 32-bit sequences as the input. Their output is thus in all cases $M = \{0, 1, \dots, 2^{32} - 1\}$.

We use two batteries of random tests – TestU01 BigCrush and PractRand. They operate differently. The first one includes 160 statistical tests, many of them tailored to the specific classes of PRNGs. It is a reputable test, however its drawback is that it works with a fixed amount of data and discards the least significant bit (for some tests even two bits) of the 32-bit numbers being tested. The second battery consists of three different tests where one is adapted on short range correlations, one reveals long range violations, and the last one is a variation on the classical Gap test. Details can be found in [6, 7]. Moreover, the PractRand battery applies automatically various filters on the input data. For our purpose the lowbit filter is interesting – it is passing various number of the least significant bits to the statistical tests. Lower bits of LCGs output with power-of-2 modulo have a much shorter period than the LCG itself. Therefore, the lowbit filter is useful to check how proposed mixing scheme can cope with this weakness. The PractRand tests are also able to treat very long input sequences, up to a few exabytes. To control the runtime we have limited the length of input sequences to 16TB.

The first column of Table 2 shows the list of tested LCGs. The BigCrush column shows how many tests of the TestU01 BigCrush battery failed. The PractRand column gives the \log_2 of sample datasize in Bytes for which the results of the PractRand tests started to be “very suspicious” (p -values smaller than 10^{-5}). One LCG did not show any failures in the PractRand tests which is denoted as > 44 – the meaning is that the PractRand test has passed successfully 16TB of input data and the test was stopped there. The last column provides time in seconds to generate the first 10^{10} 32-bit sequences of output on Intel i7-3520M CPU running at 2.90GHz.

Generator	Legend	BigCrush	PractRand	Time 10^{10}
LCG($2^{47} - 115, 71971110957370, 0$)	L47-115	14	40	281
LCG($2^{63} - 25, 2307085864, 0$)	L63-25	2	>44	277
LCG($2^{59}, 13^{13}, 0$)	L59	19	27	14.1
LCG($2^{63}, 5^{19}, 1$)	L63	19	33	14.4
LCG($2^{64}, 2862933555777941757, 1$)	L64_28	18	35	14.0
LCG($2^{64}, 3202034522624059733, 1$)	L64_32	14	34	14.1
LCG($2^{64}, 3935559000370003845, 1$)	L64_39	13	33	14.0

Table 2: List of the used LCGs with parameters LCG(m, a, c). Results in the BigCrush (number of failed tests) and in the PractRand (\log_2 of sample size for which the test started to fail) battery of statistical tests. Time in seconds to generate the first 10^{10} 32-bit words of output on Intel i7-3520M CPU running at 2.90GHz.

From Table 2 it can be seen that the LCGs with $m \in \{2^{47} - 115, 2^{63} - 25\}$ have the best statistical properties from the chosen LCGs. At the same time, these LCGs are 20 times slower than the other LCGs used. This is because we have used 128-bit integer arithmetic to compute their internal state and because explicit modulo operation cannot be avoided. As the CPU used does not have the 128-bit integer arithmetic, it has to be implemented in software (in this case via GCC’s `__int128` type) which is much slower than the 64-bit arithmetic wired on CPU.

5.2.2. Results in Statistical Tests

We will present results for the PRNGs based on the Fibonacci, Fibonacci2, Tribonacci and $AR(\sigma_0, \sigma_1, \sigma_2, \sigma_1)$ words using the different combinations of LCGs from Table 2. It includes also the situations where the instances of the same LCG are used. Each instance has its own state. The LCGs were seeded with the value 1. The PRNGs were warmed up by generating 10^9 values before statistical tests started. Since the relative frequency of the letters in the aperiodic words differ a lot (for example for the Fibonacci word the ratio of zeroes to ones is given by $\tau = \frac{1+\sqrt{5}}{2}$), the warming procedure will guarantee that the state of instances of LCGs will differ even when the same LCGs are used. Even more importantly, the distance between the LCGs is growing as the new output of PRNGs is generated.

Summary of results is in Table 3. The BigCrush column is using the following notation: the first number indicates how many tests from the BigCrush battery have clearly failed and the optional second number in parenthesis denotes how many tests have suspiciously low p -value in the range from 10^{-6} to 10^{-4} . The PractRand column gives the \log_2 of sample datasize in Bytes for which the results of the PractRand tests started to be “very suspicious” (p -values smaller than 10^{-5}). The maximum sample data size used was $16\text{TB} \doteq 2^{44}\text{B}$. The Time column gives runtime in seconds to generate the first 10^{10} 32-bit words of output on Intel i7-3520M CPU running at 2.90GHz. The source code of the testing programs is in [15].

Word	Group	0	1	2	BigCrush	PractRand	Time 10^{10}
Fib	A	L64_28	L64_28		0	41	30.2
		L64_32	L64_28		0(1)	41	29.3
		L64_39	L64_28		0 (2)	41	31
		L64_28	L64_32		0	41	30.2
		L64_32	L64_32		0	41	30.1
		L64_39	L64_32		0	41	30.1
		L64_28	L64_39		0	42	30.2
		L64_32	L64_39		0	40	30.5
		L64_39	L64_39		0	42	30.1
	B	L47-115	L47-115		1(1)	>44	302
		L63-25	L63-25		0(1)	>44	299
		L59	L59		0(1)	34	28.7
		L63	L63		0	40	29.8
	C	L63-25	L59		0	38	198
		L59	L63-25		0(1)	35	134
		L63-25	L64_39		0	>44	199
		L64_39	L63-25		0	41	135
		L59	L64_39		0	35	30.4
		L64_39	L59		0	37	31.3
Fib2	A	L64_28	L64_28	L64_28	0	40	28.4
		L64_39	L64_28	L64_28	0(2)	40	27.9
		L64_39	L64_32	L64_28	0	39	27.5
		L64_28	L64_39	L64_28	0	40	27.3
		L64_32	L64_39	L64_28	0	40	27.5
		L64_39	L64_39	L64_28	0	40	27.4
		L64_39	L64_28	L64_32	0	40	27.3
		L64_28	L64_39	L64_32	0	40	27.9

Continued on the next page

Table 3 – Continued from the previous page

Word	Group	0	1	2	BigCrush	PractRand	Time 10 ¹⁰	
		L64_28	L64_28	L64_39	0(1)	40	27.4	
		L64_32	L64_28	L64_39	0	39	27.7	
		L64_39	L64_28	L64_39	0	40	27.3	
		L64_28	L64_32	L64_39	0	40	27.3	
		L64_28	L64_39	L64_39	0	40	27.3	
		L64_39	L64_39	L64_39	0	40	27.4	
	B	L47-115	L47-115	L47-115	0(2)	>44	297.0	
		L63-25	L63-25	L63-25	0(2)	>44	293.0	
		L59	L59	L59	0(1)	32	27.4	
		L63	L63	L63	0	38	27.3	
	C	L63-25	L59	L64_39	0(1)	39	113.0	
		L63-25	L64_39	L59	0	32	113.0	
		L59	L63-25	L64_39	0	38	81.1	
		L59	L64_39	L63-25	0	39	158.3	
		L64_39	L63-25	L59	0	31	81.0	
		L64_39	L59	L63-25	0	42	159.0	
	Trib	A	L64_28	L64_28	L64_28	0(2)	42	27.2
			L64_39	L64_28	L64_28	0	43	27.1
L64_39			L64_32	L64_28	0(1)	42	28.0	
L64_28			L64_39	L64_28	0(1)	42	28.1	
L64_32			L64_39	L64_28	0	42	27.1	
L64_39			L64_39	L64_28	0(1)	42	27.2	
L64_39			L64_28	L64_32	0	43	27.1	
L64_28			L64_39	L64_32	0(1)	42	27.1	
L64_28			L64_28	L64_39	0	42	28.0	
L64_32			L64_28	L64_39	0	42	27.2	
L64_39			L64_28	L64_39	0(1)	43	27.1	
L64_28			L64_32	L64_39	0	43	27.1	
L64_28			L64_39	L64_39	0(2)	42	27.3	
L64_39			L64_39	L64_39	0	43	27.1	
B		L47-115	L47-115	L47-115	1	>44	299.0	
		L63-25	L63-25	L63-25	0(1)	>44	298.0	
		L59	L59	L59	0	35	27.2	
		L63	L63	L63	0(1)	41	27.2	
C		L63-25	L59	L64_39	0(1)	39	172.0	
		L63-25	L64_39	L59	0(1)	41	173.0	
		L59	L63-25	L64_39	0	35	106.0	
		L59	L64_39	L63-25	0	34	70.5	
		L64_39	L63-25	L59	0	41	107.0	
		L64_39	L59	L63-25	0(1)	40	74.3	
AR	A	L64_28	L64_28	L64_28	0	43	27.3	
		L64_39	L64_28	L64_28	0	43	27.0	

Continued on the next page

Table 3 – Continued from the previous page

Word	Group	0	1	2	BigCrush	PractRand	Time 10^{10}
		L64_39	L64_32	L64_28	0	43	27.4
		L64_28	L64_39	L64_28	0(1)	43	27.8
		L64_32	L64_39	L64_28	0	42	27.1
		L64_39	L64_39	L64_28	0	43	27.5
		L64_39	L64_28	L64_32	0(1)	43	27.2
		L64_28	L64_39	L64_32	0	43	27.1
		L64_28	L64_28	L64_39	0	43	27.3
		L64_32	L64_28	L64_39	0	42	27.1
		L64_39	L64_28	L64_39	0(1)	43	27.4
		L64_28	L64_32	L64_39	0	42	27.9
		L64_28	L64_39	L64_39	0	43	27.8
		L64_39	L64_39	L64_39	0(1)	43	27.4
	B	L47-115	L47-115	L47-115	1	>44	299
		L63-25	L63-25	L63-25	0	>44	297
		L59	L59	L59	0	36	27.3
		L63	L63	L63	0	42	27.2
	C	L63-25	L59	L64_39	0	39	173
		L63-25	L64_39	L59	0	43	174
		L59	L63-25	L64_39	0	36	104
		L59	L64_39	L63-25	0	36	70.1
		L64_39	L63-25	L59	0	42	107
		L64_39	L59	L63-25	0	39	74.7

Table 3: Summary of results of statistical tests for PRNGs based on the Fibonacci, Fibonacci2, Tribonacci and $AR(\sigma_0, \sigma_1, \sigma_2, \sigma_1)$ words and different combinations of LCGs from Table 2.

We can make the following observations based on the results in statistical tests:

1. The quality of LCGs improved substantially when we combined them according to infinite words with the WELLDLOC property. This can be seen in the TestU01 BigCrush results. While for the plain LCGs 13 to 19 tests have clearly failed (the only exception is the generator L63-25 with two failures – see Table 2), almost all of the BigCrush tests passed for the combined generators. The worst result was to have one BigCrush test failed for the combination of two (respective three) instances of L47-115 based on Fibonacci, Tribonacci and $AR(\sigma_0, \sigma_1, \sigma_2, \sigma_1)$ words. The likely reason is that the generator L47-115 has the shortest period of all tested LCGs.
2. The results of the PractRand battery confirm the above findings. For instance, in the case of LCGs with modulo 2^{64} , the test started to find irregularities in the distribution of the least significant bit of tested PRNGs output at around 2TB sample size. Compare it with the sample size of 8GB to 32GB when fast plain LCGs started to fail the test. The PractRand battery applies different filters on the input stream and all failures appeared for Low1/32 filter where only the least significant bit of the PRNG output is used. It corresponds to a known weakness of power-of-2 modulo LCGs: lower bits of the output have significantly smaller period than the LCG itself. The quality of the PRNGs can be therefore further improved by combining LCGs that do not show flaws for the least significant bits or by using for example just 16 upper bits of the LCGs output.

3. The quality of the PRNG is linked to the quality of the underlying LCG. When looking at the group B in Table 3, we observe that the PractRand results of the arising PRNGs are closely related to the success of LCGs from Table 2 in the PractRand tests. We also tested LCGs with $m = 2^{31} - 1$. It has revealed that if the underlying generators have poor statistical properties, then the PRNG will not be able to mask it. In particular, you cannot expect that PRNGs – despite their infinite aperiodic nature – will fix the short period problem. Once the period of the underlying LCG is exhausted, statistical tests will find irregularities in the output of the PRNG.
4. Using the instances of the same LCG (with only sufficiently distinct seeds) produces as good results as combination of different LCGs (multipliers and shifts are different, but the modulus is the same). It is just important to make sure that starting states of the LCGs are far apart enough. Refer to the group A in Table 3.
5. The lower quality LCG dictates the quality of resulting PRNG. When mixing LCGs with different quality, use better ones as replacement for more frequent letters in the aperiodic word. We refer to the group C in Table 3. For example for the Fibonacci word compare first two rows in the group C: the order of LCGs is merely swapped but the difference in the sample size for which PractRand starts to fail is $8\times$. This is even more significant for the Tribonacci based generators where the difference between the worst and best PractRand results when reordering the underlying LCGs is given by factor $128\times$.
6. On the other hand, results from the group A in Table 3 demonstrate that when using generators of similar quality (same modulus, similar deficiencies), the order in which generators are used to substitute the letters of the infinite word does not influence the quality of the resulting generator.
7. We can also see that the modified Fibonacci word (see Remark 2.5) does not produce better results than the Fibonacci word. Clearly, a regular structure of 2's on every other position does not help to produce a better random sequence even if we mix now three LCGs instead of two as in the case of the Fibonacci word.
8. Results for the Tribonacci and $AR(\sigma_0, \sigma_1, \sigma_2, \sigma_1)$ words are better than for the Fibonacci word. (We observed this fact for all ternary Arnoux-Rauzy words in comparison to Sturmian words.) It seems therefore that mixing three LCGs is better than using just two LCGs, assuming that an infinite word with the WELLDOC property is used for mixing. We expect naturally that the better chosen LCGs (or even some other modern fast linear PRNGs, e.g. *mt19937* or nonlinear PRNGs based on the AES cipher) we combine according to an infinite word with the WELLDOC property, the better their results in statistical tests will be.
9. The results for all the tested Arnoux-Rauzy words from the class described in Section 5.2 were very similar. This class enables to define an infinite number of different aperiodic words. For example, one can define following scheme to generate aperiodic word based on user supplied seed. The seed can be converted to ternary numeral system where each ternary digit will represent one σ rule. To satisfy the condition that each of $\sigma_0, \sigma_1, \sigma_2$ rules has to be used at least once we can prepend the seed with prefix, for example $\sigma_0, \sigma_1, \sigma_2$. The resulting Arnoux-Rauzy word will be defined according to rule $\sigma_0, \sigma_1, \sigma_2, \sigma_{s_1}, \dots, \sigma_{s_n}$ where s_1, \dots, s_n are ternary digits of the seed.
10. Clearly, not all Sturmian and Arnoux-Rauzy words are equally good from a statistical point of view. For instance, an undesirable property for the steering word is the occurrence of long powers with a short root, that would lead to long stretches where few LCGs alternate periodically. This is easily avoidable, for example, by choosing a Sturmian word having a slope whose continued fraction expansion does not contain large partial quotients. Please note that this condition also guarantees that the dominating letter frequency is not too close to 1.

In conclusion, we summarize the main results from the user point of view:

- Using different instances of the same LCG to form a new generator based on the infinite word with the WELLDOC property gives a generator with improved statistical properties.
- The introduced method of generation of morphic words is very fast and supports parallel processing.
- The period of underlying generators has to be large enough – much larger than the number of needed values.
- When using different types of the underlying LCGs to form a PRNG, the generator with the worst properties should be used to replace the least frequent letter of the aperiodic word. Moreover, statistical properties of the resulting PRNG are ruled by the deficiencies of the worst used generator.
- In the construction of PRNGs based on aperiodic infinite words, one can use other random number generators instead of LCGs. We have done testing with two instances (respectively three for the Tribonacci and Arnoux-Rauzy words) of Mersenne twister 19937 as the underlying generator. The newly constructed generator has passed all the empirical tests on randomness we have executed (in contrary to Mersenne twister 19937 itself which is failing two tests from TestU01’s BigCrush battery).

6. Open problems and future research

Concerning the combinatorial part of our paper, one of the interesting open questions there is finding large families of infinite words satisfying the WELLDOC property. For example, which morphic words have the WELLDOC property? Also, from a practical point of view it seems to be meaningful to study a weaker WELLDOC property where in Definition 2.1 instead of every $m \in \mathbb{N}$ we consider only a particular m . For instance, one can search for words satisfying such a modified WELLDOC condition for $m = 2$, $m = 2^\ell$ etc. Another question to be asked is how to construct words with the WELLDOC property over larger alphabets using words with such a property over smaller alphabets.

Acknowledgements

The first author was supported by the Czech Science Foundation grant GAČR 13-03538S, and thanks L’Oréal Czech Republic for the Fellowship Women in Science. The third author was partially supported by the Italian Ministry of Education (MIUR), under the PRIN 2010–11 project “Automi e Linguaggi Formali: Aspetti Matematici e Applicativi”. The fifth author was supported by the LABEX MILYON (ANR-10-LABX-0070) of Université de Lyon, within the program “Investissements d’Avenir” (ANR-11-IDEX-0007) operated by the French National Research Agency (ANR).

References

- [1] P. Arnoux, G. Rauzy, Représentation géométrique de suites de complexité $2n + 1$, *Bulletin de la Société Mathématique de France* 119 (1991) 199–215.
- [2] v. Balková, M. Bucci, A. De Luca, S. Puzynina, Infinite words with well distributed occurrences, in: J. Karhumäki, A. Lepistö, L. Zamboni (Eds.), *Combinatorics on Words*, volume 8079 of *Lecture Notes in Computer Science*, Springer, 2013, pp. 46–57.
- [3] É. Charlier, T. Kamae, S. Puzynina, L. Q. Zamboni, Infinite self-shuffling words, *J. Combin. Theory Ser. A* 128 (2014) 1–40.
- [4] A. de Luca, Sturmian words: structure, combinatorics, and their arithmetics, *Theoret. Comput. Sci.* 183 (1997) 45–82.

- [5] C. Doty-Humphrey, 2010, Practically Random: C++ library of statistical tests for RNGs, URL: <https://sourceforge.net/projects/pracrand>.
- [6] C. Doty-Humphrey, 2014, Specific tests in PractRand, URL: http://pracrand.sourceforge.net/Tests_engines.txt.
- [7] C. Doty-Humphrey, J. Hladký, 2013, Practically Random: Discussion of testing results, URL: <http://sourceforge.net/p/pracrand/discussion/366935/thread/a2eaaad12>.
- [8] X. Droubay, J. Justin, G. Pirillo, Episturmian words and some constructions of de Luca and Rauzy, *Theoret. Comput. Sci.* 255 (2001) 539–553.
- [9] P. L’Ecuyer, Random number generation, in: *Handbook of computational statistics*, Springer, Heidelberg, 2012, pp. 35–71. doi:10.1007/978-3-642-21551-3_3.
- [10] P. L’Ecuyer, D. Munger, B. Oreshkin, R. Simard, Random numbers for parallel computers: Requirements and methods, with emphasis on GPUs, 2015. <http://www.iro.umontreal.ca/~lecuyer/myftp/papers/parallel-rng-imacs.pdf>, preprint.
- [11] P. L’Ecuyer, R. Simard, TestU01: a C library for empirical testing of random number generators, *ACM Trans. Math. Software* 33 (2007) Art. 22, 40.
- [12] L.-S. Guimond, J. Patera, Proving the deterministic period breaking of linear congruential generators using two tile quasicrystals, *Math. Comp.* 71(237) (2002) 319–332.
- [13] L.-S. Guimond, J. Patera, J. Patera, Combining random number generators using cut and project sequences, *Czechoslovak J. Phys.* 51 (2001) 305–311. DI-CRM Workshop on Mathematical Physics (Prague, 2000).
- [14] L.-S. Guimond, J. Patera, J. Patera, Statistical properties and implementation of aperiodic pseudorandom number generators, *Appl. Numer. Math.* 46(3–4) (2003) 295–318. *Applied numerical computing: grid generation and solution methods for advanced simulations*.
- [15] J. Hladký, 2013, Random number generators based on the aperiodic infinite words, source programs for statistical tests. URL: <https://github.com/jirka-h/aprng>.
- [16] M. Lothaire, *Algebraic Combinatorics on Words*, volume 90 of *Encyclopedia of Mathematics and its Applications*, Cambridge University Press, 2002.
- [17] G. Marsaglia, Random numbers fall mainly in the planes, *Proc. Nat. Acad. Sci. U.S.A.* 61(1) (1968) 25–28.
- [18] M. Morse, G. A. Hedlund, Symbolic dynamics, *Amer. J. Math.* 60 (1938) 815–866.
- [19] M. Morse, G. A. Hedlund, Symbolic dynamics II. Sturmian trajectories, *Amer. J. Math.* 62(1) (1940) 1–42.
- [20] J. Patera, Generating the Fibonacci chain in $O(\log n)$ space and $O(n)$ time, *Phys. Particles Nuclei* 33(7) (2002) 225–234.