

The Sum of Exponents of Maximal Repetitions in Standard Sturmian Words

Marcin Piątkowski

Faculty of Mathematics and Computer Science,
Nicolaus Copernicus University, Toruń, Poland
`marcin.piatkowski@mat.umk.pl`

Abstract. A maximal repetition is a non-extendable (with the same period) periodic segment in a string, in which the period repeats at least twice. In this paper we study problems related to the structure of maximal repetitions in standard Sturmian words and present the formulas for the sum of their exponents. Moreover, we show how to compute the sum of exponents of maximal repetitions in any standard Sturmian word in linear time with respect to the (total) size of its compressed representation. The presented formulas and algorithm can be easily modified to obtain the total run length of the word.

Keywords: Sturmian words, repetitions, runs, algorithm

1 Introduction

Problems related to repetitions are fundamental in combinatorics on words and many practical applications: data compression, computational biology, pattern-matching and so on, see for instance [6], [7], [10], [11], [14] and references therein. The most important type of repetitions are maximal repetitions, i.e. non-extendable (with the same period) periodic segments in a string, in which the period repeats at least twice. This paper complements the work [2], where the exact formula for the number of runs in standard Sturmian words was presented. We investigate here the structure of runs in standard Sturmian words in more details to obtain a formula for the sum of their exponents. We show also an algorithm, derived from our formula, which computes the sum of exponents of maximal repetitions in any standard word in linear time with respect to the (total) size of its compressed representation (i.e. the directive sequence).

Throughout the paper we use the standard notions of combinatorics on words. In particular, words are finite sequences over a finite set Σ of letters, called the alphabet. For a word $w = w_1w_2 \cdots w_n$, by w_i we denote its i -th letter, by $w[i..j]$ the subword $w_iw_{i+1} \cdots w_j$, by $|w|$ its length and by $|w|_a$ the number of letters a occurring in w . The number i is a period of the word w if $w_j = w_{i+j}$ for all i with $i + j \leq |w|$. The minimal period of w is denoted by $period(w)$. We say that a word w is periodic if $period(w) \leq \frac{|w|}{2}$. A word w is said to be *primitive* if w is not of the form z^k , where z is a nonempty word and $k \geq 2$ is a natural number.

A *maximal repetition* (a *run*, in short) in a word w is an interval $\alpha = [i..j]$, such that $w[i..j] = u^k v$ ($k \geq 2$) is a nonempty periodic subword of w , where u is of the minimal length and v is a proper prefix (possibly empty) of u , that can not be extended (neither $w[i-1..j]$ nor $w[i..j+1]$ is a run with the period $|u|$). The factor v is called the remainder of α and the number $k + \frac{|v|}{|u|}$ is called the exponent of α . The

In 1999 Kolpakov and Kucherov showed that the number of runs in a word is linear with respect to its length (see [13]). The stronger property of runs is that the sum of their exponents is also linear with respect to the length of the word. Kolpakov and Kucherov conjectured that for all w we have $\sigma(w) \leq 2 \cdot |w|$. In 2012 Crochemore with coauthors contradicted this conjecture and showed that the upper bound for $\sigma(w)$ is $2.035 \cdot |w| \leq \sigma(w) \leq 4.1 \cdot |w|$. In this paper we investigate this problem in very special class of strings – the standard Sturmian words. We present compact formulas for the sum of runs exponents and an algorithm for its efficient computation.

Recently a new measure of a string periodicity was proposed by Glen and Simpson (see [12]). The *total run length* (TRL) of a word w is the sum of the lengths of all runs in w . Since this notion is similar to the sum of exponents of maximal repetitions, our formulas and algorithm could be easily adopted to compute also the total run length of any standard Sturmian word.

The paper is organized as follows. In section 2 we introduce the definition of standard Sturmian words and some of their basic properties. Next, in section 3 we study the structure of repetitions in standard Sturmian words and present a few facts necessary in further investigation. Finally, we show and prove the formulas for the sum of exponents of maximal repetitions together with an algorithm for its fast computation. Some useful applets related to problems considered in this paper can be found on the web site:

<http://www.mat.umk.pl/~martinp/stringology/applets/>

2 Standard Sturmian words

Standard Sturmian words (standard words in short) are one of the most investigated class of strings in combinatorics on words, see for instance [1], [3], [4], [5], [15], [17], [18] and references therein. They have very compact representations in terms of sequences of integers, which has many algorithmic consequences.

The *directive sequence* is the integer sequence: $\gamma = (\gamma_0, \gamma_1, \dots, \gamma_n)$, where $\gamma_0 \geq 0$ and $\gamma_i > 0$ for $i = 1, 2, \dots, n$. The standard word corresponding to γ , denoted by $\text{Sw}(\gamma)$, is described by the recurrences of the form:

$$x_{-1} = b, \quad x_0 = a, \quad \dots, \quad x_n = (x_{n-1})^{\gamma_{n-1}} x_{n-2}, \quad x_{n+1} = (x_n)^{\gamma_n} x_{n-1}, \quad (1)$$

where $\text{Sw}(\gamma) = x_{n+1}$. For simplicity we denote $q_i = |x_i|$.

Example 2.

Consider the directive sequence $\gamma = (1, 2, 1, 3, 1)$. We have $\text{Sw}(\gamma) = x_5$, where:

$x_{-1} = b$	$q_{-1} = 1$
$x_0 = a$	$q_0 = 1$
$x_1 = (x_0)^1 \cdot x_{-1} = a \cdot b$	$q_1 = 2$
$x_2 = (x_1)^2 \cdot x_0 = ab \cdot ab \cdot a$	$q_2 = 5$
$x_3 = (x_2)^1 \cdot x_1 = ababa \cdot ab$	$q_3 = 7$
$x_4 = (x_3)^3 \cdot x_2 = ababaab \cdot ababaab \cdot ababaab \cdot ababa$	$q_4 = 26$
$x_5 = (x_4)^1 \cdot x_3 = ababaabababaabababaabababa \cdot ababaab$	$q_5 = 33$

The sequence of words $\{x_i\}_{i=0}^{n+1}$ is called the standard sequence. Every word occurring in a standard sequence is a standard word, and every standard word occurs in some standard sequence. We assume that the standard word given by the empty directive sequence is a and $\text{Sw}(0) = b$.

Observe that for even $n > 0$ the standard word x_n has the suffix ba , and for odd $n > 0$ it has the suffix ab . Moreover, for $\gamma_0 > 0$ we have standard words starting with the letter a and for $\gamma_0 = 0$ we have standard words starting with the letter b . In fact the word $\text{Sw}(0, \gamma_1, \dots, \gamma_n)$ can be obtained from $\text{Sw}(\gamma_1, \dots, \gamma_n)$ by switching the letters a and b . Without loss of generality we consider here standard words starting with the letter a , therefore we assume $\gamma_0 > 0$. Words starting with the letter b can be considered similarly.

Remark 3.

The special kind of standard words are well known Fibonacci words. They are formed by repeated concatenation in the same way that the Fibonacci numbers are formed by repeated addition. By definition Fibonacci words are standard words given by directive sequences of the form $\gamma = (1, 1, \dots, 1)$ (n -th Fibonacci word F_n corresponds to a sequence of n ones).

The number $N = |\text{Sw}(\gamma)|$ is the (real) size of the word, while $(n + 1) = |\gamma|$ can be thought as its compressed size. Observe that, by the definition of standard words, N is exponential with respect to n . Moreover, each directive sequence corresponds to a *grammar-based compression*, which consists in describing a given word by a context-free grammar G generating this (single) word. The size of the grammar G is the total length of all productions of G . In our case the size of the considered grammar is proportional to the length of the directive sequence.

2.1 Morphic reduction of standard words

The recurrent definition of standard words from equation (1) leads to their simple characterization by a composition of morphisms. Let $\gamma = (\gamma_0, \gamma_1, \dots, \gamma_n)$ be a directive sequence. We associate with γ a sequence of morphisms $\{h_i\}_{i=0}^n$, defined as:

$$h_i : \begin{cases} a \longrightarrow a^{\gamma_i} b \\ b \longrightarrow a \end{cases} \quad \text{for } 0 \leq i \leq n. \quad (2)$$

The following fact describes another simple method of standard word generation. It can be proven by a simple induction, see [2] for more details.

Lemma 4 (see [2]).

For $0 \leq i \leq n$ the morphism h_i transforms a standard word into another standard word, and we have:

$$\begin{aligned} \text{Sw}(\gamma_n) &= h_n(a), \\ \text{Sw}(\gamma_i, \gamma_{i+1}, \dots, \gamma_n) &= h_i(\text{Sw}(\gamma_{i+1}, \gamma_{i+2}, \dots, \gamma_n)). \end{aligned}$$

As a direct corollary to Lemma 4 we have that for $\gamma = (\gamma_0, \gamma_1, \dots, \gamma_n)$:

$$\text{Sw}(\gamma_0, \gamma_1, \dots, \gamma_n) = h_0 \circ h_1 \circ \dots \circ h_n(a). \quad (3)$$

Moreover, the inverse morphism h_i^{-1} can be seen as a reduction of a standard word $w^{(i)} = \text{Sw}(\gamma_i, \dots, \gamma_n)$ to $w^{(i+1)} = \text{Sw}(\gamma_{i+1}, \dots, \gamma_n)$.

Recall that $|w|_a$ denotes the number of occurrences of the letter a in the word w . In the rest of this paper, for $\gamma = (\gamma_0, \dots, \gamma_n)$ and $0 \leq k \leq n$, we use the following notation:

$$N_\gamma(k) = |\text{Sw}(\gamma_k, \gamma_{k+1}, \dots, \gamma_n)|_a, \quad (4)$$

which enables us to simplify the formulas for the sum of runs exponents. Observe that equations (2) and (4) imply:

$$N_\gamma(k) = \gamma_k \cdot N_\gamma(k+1) + N_\gamma(k+2). \quad (5)$$

Example 5.

Consider a directive sequence $\gamma = (1, 2, 1, 3, 1)$. We have (compare with Example 2):

$$\begin{aligned} \text{Sw}(1, 2, 1, 3, 1) &= ababaabababaabababaabababaab & N_\gamma(0) &= 19, \\ \text{Sw}(2, 1, 3, 1) &= aabaaabaaabaaabaaba & N_\gamma(1) &= 14, \\ \text{Sw}(1, 3, 1) &= abababaab & N_\gamma(2) &= 5, \\ \text{Sw}(3, 1) &= aaaba & N_\gamma(3) &= 4, \\ \text{Sw}(1) &= ab & N_\gamma(4) &= 1, \\ \text{Sw}(\varepsilon) &= a & N_\gamma(5) &= 1. \end{aligned}$$

As a straightforward corollary to equations (2), (4) and (5) we have:

Corollary 6.

The number of letters b in a word $\text{Sw}(\gamma_i, \dots, \gamma_n)$ equals $N_\gamma(i+1)$.

2.2 The m -partition of a standard word

The concept of the m -partition of a standard word is crucial in the maximal repetitions structure investigation. It allows us to divide the set of all runs in a standard word to disjoint sets depending on the length of their periods and simplify the considered problems. The following fact is a direct consequence of the recurrent definition of standard words.

Proposition 7.

Every standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ can be represented as a sequence of concatenated words x_m and x_{m-1} , and has the form:

$$(i) \quad x_m^{\alpha_1} x_{m-1} x_m^{\alpha_2} x_{m-1} \cdots x_m^{\alpha_s} x_{m-1} x_m \quad \text{or} \quad (ii) \quad x_m^{\beta_1} x_{m-1} x_m^{\beta_2} x_{m-1} \cdots x_m^{\beta_s} x_{m-1},$$

where $\alpha_k, \beta_k \in \{\gamma_m, \gamma_m + 1\}$, $0 \leq m \leq n$, and x_m are as in equation (1).

Such a decomposition of a standard word w is called the m -partition of w . The block x_m is called the *repeating block* and x_{m-1} – the *single block*. Recall that for $m > 0$ the last two letters of x_m are ab for an odd m and ba for an even m . Therefore the m -partition of $x_{n+1} = \text{Sw}(\gamma_0, \dots, \gamma_n)$ is of the form (i) if m has the same parity as $(n+1)$, and of the form (ii) otherwise (see Example 9 and Figure 2).

Note that the 0-partition of a standard word is its decomposition into letters. Moreover, Proposition 7, Lemma 4 and equation (3) imply the following fact.

To prove the above lemma it is sufficient to show that no factor of a standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ that does not satisfy the condition given there could be the generator of some repetition, see the proof of Theorem 1 in [9] for more details.

Let us denote by \widehat{w} the word w with two last letters removed and by \widetilde{w} the word w with two last letters exchanged. The following fact can be proven by a simple induction, see for instance [15].

Lemma 13.

Let x_i be as in equation (1) and $i > 1$. Then:

1. We have $x_{i-1} \cdot x_i = x_i \cdot \widetilde{x_{i-1}}$,
2. The longest prefix of $x_{i-1} \cdot x_i$ with the period of the length q_i is of the form $x_i \cdot \widehat{x_{i-1}}$.

Example 14.

Recall the word $\text{Sw}(1, 2, 1, 3, 1)$ from Example 2, where $x_3 = ababaab$, $x_2 = ababa$. Then we have $\widetilde{x_2} = abaab$, $\widehat{x_2} = aba$ and

$$x_2 \cdot x_3 = ababa \cdot ababaab = ababaab \cdot abaab = x_3 \cdot \widetilde{x_2}.$$

Moreover, the longest prefix of $x_2 \cdot x_3$ with the period of the length q_2 is of the form:

$$\underbrace{\overbrace{abab}^{x_2} \overbrace{aaba}^{x_3}}_{x_3 \cdot \widehat{x_2}}$$

Observe that by equation (1) we have

$$\text{Sw}(\gamma_0, \dots, \gamma_n, 1) = (x_n)^{\gamma_n} \cdot x_{n-1} \cdot x_n \quad \text{and} \quad \text{Sw}(\gamma_0, \dots, \gamma_n + 1) = (x_n)^{\gamma_n} \cdot x_n \cdot x_{n-1}.$$

Therefore, as a straightforward corollary to the first point of Lemma 13 we get:

Corollary 15.

Standard words $\text{Sw}(\gamma_0, \dots, \gamma_n, 1)$ and $\text{Sw}(\gamma_0, \dots, \gamma_n + 1)$ differ only in the order of the last two letters.

See Figure 3 for an illustration of this fact. To properly count the exponents of runs in standard words we need also the following fact.

Proposition 16.

Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word and $2 \leq i \leq n - 2$. If x_{i-1} is the last block of the i -partition of w , then it is preceded by $(x_i)^{\gamma_i + 1}$.

Proof.

Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word and $2 \leq i \leq n - 2$. By equation (1) we have $x_i = \text{Sw}(\gamma_0, \dots, \gamma_i)$ and $x_{i-1} = \text{Sw}(\gamma_0, \dots, \gamma_{i-1})$. Recall that x_i ends with ba for even $i > 0$ (i.e. for the odd length of a directive sequence) and with ab for odd $i > 0$ (i.e. for the even length of a directive sequence). Consider that w has the suffix $(x_i)^\alpha x_{i-1}$. Then n and i have the same parity and the number $n - m + 1$ is odd, hence the word $w^{(m)} = \text{Sw}(\gamma_m, \dots, \gamma_n)$ ends with ab . More precisely, due to Proposition 8, $w^{(m)}$ ends with $a^\alpha b$. By Lemma 4, the suffix $a^{\gamma_i} b$ of $w^{(m)}$ corresponds to the last letter a of $w^{(m+1)} = \text{Sw}(\gamma_{m+1}, \dots, \gamma_n)$. Since $n - m + 2$ is even and $w^{(m+1)}$ ends with ba , due to Lemma 4 the suffix $a^{\gamma_i} b$ of $w^{(m)}$ have to be preceded by a single occurrence of a . Therefore, we have $\alpha = \gamma_i + 1$ and this completes the proof. \square

4 The sum of exponents of maximal repetitions

In this section we present and prove formulas for the the sum of exponents of maximal repetitions in any standard word, that depend only on its compressed representation – the directive sequence. The following zero-one functions for testing the parity of a nonnegative integer i will be useful to simplify those formulas:

$$\text{even}(i) = \begin{cases} 1 & \text{for even } i \\ 0 & \text{for odd } i \end{cases} \quad \text{and} \quad \text{odd}(i) = \begin{cases} 1 & \text{for odd } i \\ 0 & \text{for even } i \end{cases}.$$

Moreover, we define an auxiliary function $\Delta_n : \mathbf{N} \rightarrow \mathbf{N}$:

$$\Delta_n(i) = |n - i + 1| \bmod 2.$$

In other words, $\Delta_n(i) = 1$ if and only if the numbers n and i have the same parity, and $\Delta_n(i) = 0$ otherwise. Recall also that for simplicity we denote $|x_i| = q_i$.

The main idea of the computation of the sum of runs exponents in a standard word w is the partition of the set of all maximal repetitions in w into separate categories depending on the length of their periods. Runs in w with the period of the form x_i and $(x_i)^k x_{i-1}$ (for $1 < k < \gamma_i$), where x_i are as in equation (1), are called the runs of *type* i . We study runs of each type separately.

Let $\sigma_i(\gamma)$ denotes the sum of exponents of type i runs. Then the sum of exponents of all runs in $\text{Sw}(\gamma)$ can be computed using the following theorem.

Theorem 17.

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence. The sum of exponents of runs in $\text{Sw}(\gamma)$ is given as:

$$\sigma(\gamma) = \sum_{i=1}^n \sigma_i(\gamma).$$

The detailed computation of $\sigma_i(\gamma)$ for each $0 \leq i \leq n$ is provided below.

4.1 The general case

We start with an investigation of a general case, i.e. maximal repetitions of the type i for $2 \leq i \leq n - 1$. First, we consider runs with the period of the form x_i .

Lemma 18.

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word. For $2 \leq i \leq n - 1$ the sum of exponents of runs with the period x_i in w equals:

$$\sigma'_i(\gamma) = N_\gamma(i + 1) \cdot \left(\gamma_i + 1 + \frac{q_{i-1} - 2}{q_i} \right) + \left(N_\gamma(i + 2) - 1 \right) + \Delta_n(i) \frac{2}{q_i}. \quad (6)$$

Proof.

Let us denote

$$w = \text{Sw}(\gamma_0, \dots, \gamma_n), \quad w^{(i)} = \text{Sw}(\gamma_i, \dots, \gamma_n) \quad \text{and} \quad w^{(i+1)} = \text{Sw}(\gamma_{i+1}, \dots, \gamma_n).$$

Due to Lemma 11, each maximal repetition with the period x_i in w is aligned to the i -partition of w , hence it corresponds to a block $(x_i)^\alpha x_{i-1}$, where $\alpha \in \{\gamma_i, \gamma_i + 1\}$.

Each internal block of this form is followed by a subsequent x_i . Due to Lemma 13, the longest prefix of $x_{i-1}x_i$ with the period x_i equals $x_i \cdot \widehat{x_{i-1}}$. Therefore, the period of the considered run repeats $\alpha + 1$ times and its fractional part has the length $q_{i-1} - 2$.

Consider the i -partition of w . By Proposition 8 occurrences of x_i correspond to occurrences of a in $w^{(i)}$ and occurrences of x_{i-1} correspond to occurrences of b in $w^{(i)}$. Therefore, a block $(x_i)^\alpha x_{i-1}$ correspond to the block $a^\alpha b$ in $w^{(i)}$. Moreover, due to Lemma 4, each block of the form $a^{\gamma_i+1}b$ in $w^{(i)}$ corresponds to the letter a preceded by the letter b in $w^{(i+1)}$ and each block of the form $a^{\gamma_i}b$ in $w^{(i)}$ corresponds to the letter a not preceded by the letter b in $w^{(i+1)}$.

The rightmost occurrence of $(x_i)^\alpha x_{i-1}$ have to be considered separately. Due to Proposition 7, if i and n have different parity the i -partition of w ends with $(x_i)^\alpha x_{i-1}x_i$. In this case the period of the considered repeats $\alpha + 1$ times and its fractional part has the length $q_{i-1} - 2$. On the other hand, if i and n have the same parity, the i -partition of w ends with $(x_i)^\alpha x_{i-1}$. Due to Proposition 16, $\alpha = \gamma_i + 1$. Moreover, since x_{i-1} is a prefix of x_i , the fractional part of considered run consists of the whole word x_i and has the length q_{i-1} .

Summing up, in the computation of the sum of runs exponents, we count $\gamma_i + 1 + \frac{q_{i-1}-2}{q_i}$ for each occurrence of a in $w^{(i+1)}$, namely $N_\gamma(i + 1)$ times, and an additional 1 for each b in $w^{(i+1)}$ (except the rightmost one), namely $N_\gamma(i + 2) - 1$ times. Finally, we must take care of the remainder of the rightmost run with period x_i and we obtain the statement of the lemma. See Figure 3 for the illustration of type-2 runs structure in example words and two possible remainders of the rightmost run. \square

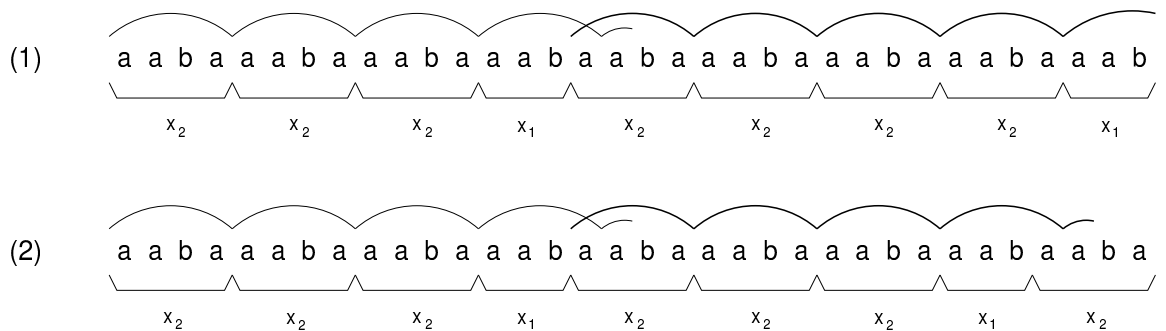


Figure 3. The structure of runs with the period x_2 in a standard word $\text{Sw}(2, 1, 3, 1, 1)$ (1) compared to $\text{Sw}(2, 1, 3, 2)$ (2).

Observe that the maximal repetitions with the period of the form $(x_i)^k x_{i-1}$, where $1 \leq k < \gamma_i$, appear only for $\gamma_i > 1$. The sum of exponents of such runs is given by the following fact.

Lemma 19.

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word. For $1 \leq i \leq n - 1$ the sum of exponents of runs in w with the period $(x_i)^k x_{i-1}$, where $1 \leq k < \gamma_i$, equals:

$$\sigma''_i(\gamma) = \left(N_\gamma(i + 1) - 1 \right) \cdot \sum_{k=1}^{\gamma_i-1} \left(2 + \frac{q_i - 2}{k \cdot q_i + q_{i-1}} \right). \tag{7}$$

Proof.

Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ and $u = (x_i)^k x_{i-1}$, where $1 \leq k < \gamma_i$. Due to Lemma 11, each occurrence of u is aligned to the i -partition of w . Consider a repetition of the form u^m in w and denote it as $u^{(1)}u^{(2)} \dots u^{(m)}$. Observe that each $u^{(2)}, \dots, u^{(m)}$ have to be preceded by the suffix of u , namely x_{i-1} . Since each two consecutive occurrences of x_{i-1} in the i -partition of w are separated by at least γ_i occurrences of x_i and $k < \gamma_i$, the factor u cannot have more than two consecutive occurrences. Therefore, the considered run with the period u has the form $u^{(1)}u^{(2)} \cdot v$, where v is a prefix of u .

The suffix x_{i-1} of $u^{(2)}$ starts at the beginning of an x_i block followed by x_{i-1} , which appears either as block of the i -partition of w or as a prefix of a subsequent block x_i . Due to Lemma 13, the considered factor has the form $x_i \cdot x_{i-1} = x_{i-1} \cdot \widetilde{x}_i$. Therefore, the fractional part of the considered run has the length $q_i - 2$.

Observe, that occurrences of $u^{(1)}$ in w are aligned with occurrences of x_{i-1} in the i -partition of w . Therefore, each such occurrence of x_{i-1} (except the rightmost one) corresponds to $\gamma_i - 1$ runs with a period $(x_i)^k x_{i-1}$, for $1 \leq k < \gamma_i$. Due to Proposition 8, each occurrence of x_{i-1} in the i -partition of w corresponds to an occurrence of b in $\text{Sw}(\gamma_i, \dots, \gamma_n)$. Summing up exponents of all $\gamma_i - 1$ runs for each b in $\text{Sw}(\gamma_i, \dots, \gamma_n)$ (except the rightmost one), namely $N_\gamma(i+1) - 1$ occurrences, we obtain the statement of the lemma. See Figure 4 for an illustration of the structure of runs of this type. \square

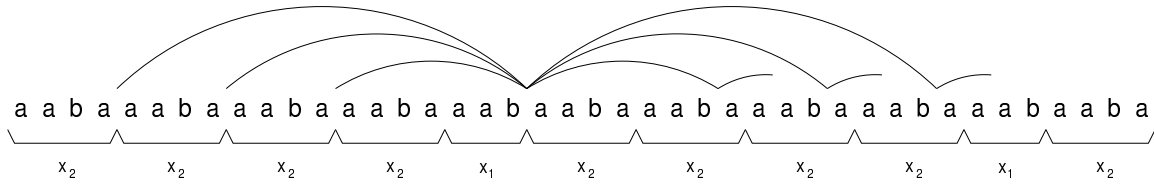


Figure 4. The structure of runs with the period $(x_2)^k x_1$ ($1 \leq k \leq 3$) in $\text{Sw}(2, 1, 4, 2)$.

The complete formula for the sum of exponents of all type- i runs can be obtained by combining the formulas from Lemma 18 and Lemma 19.

Lemma 20.

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word. For $2 \leq i \leq n - 1$ the sum of exponents of type i runs in w equals:

$$\begin{aligned} \sigma_i(\gamma) = & N_\gamma(i+1) \cdot \left(\gamma_i + 1 + \frac{q_{i-1} - 2}{q_i} \right) + \left(N_\gamma(i+2) - 1 \right) + \Delta_n(i) \frac{2}{q_i} \\ & + \left(N_\gamma(i+1) - 1 \right) \cdot \sum_{k=1}^{\gamma_i-1} \left(2 + \frac{q_i}{k \cdot q_i + q_{i-1}} \right). \end{aligned} \tag{8}$$

4.2 Boundary cases

For a standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ runs of types 0, 1 and n have to be investigated differently. We start with the analyze of runs of type 0, i.e. the runs with the period of the form a .

Lemma 21 (Type 0).

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard

word. The sum of exponents of type 0 runs in w equals:

$$\sigma_0(\gamma) = \begin{cases} 2(N_\gamma(2) - \text{odd}(n)) & \text{for } \gamma_0 = 1 \\ \gamma_0 N_\gamma(1) + N_\gamma(2) - \text{odd}(n) & \text{for } \gamma_0 > 1 \end{cases}. \quad (9)$$

Proof.

Each standard word consists of blocks of repeated occurrences of the letter a separated by single occurrences of the letter b . The length of the blocks of the form $a \cdots a$ depends on the value of γ_0 .

First assume that $\gamma_0 = 1$. In this case the word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ consists of the blocks of two types: ab or aab and only the blocks of the second type include the runs with the period a and exponent 2. Due to Lemma 4, every such run in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ corresponds to the letter b followed by the letter a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Hence, the number of such runs equals the number of blocks ba in $\text{Sw}(\gamma_1, \dots, \gamma_n)$.

Recall that for an even length of the directive sequence $|(\gamma_1, \dots, \gamma_n)|$ (n is even) the word $\text{Sw}(\gamma_1, \dots, \gamma_n)$ ends with ba and in this case the number of runs with the period a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$ equals the number of the letters b in $\text{Sw}(\gamma_1, \dots, \gamma_n)$, namely $N_\gamma(2)$. On the other hand, for an odd length of the directive sequence $|(\gamma_1, \dots, \gamma_n)|$ (n is odd) the word $\text{Sw}(\gamma_1, \dots, \gamma_n)$ ends with ab and the last letter b does not correspond to a run in $\text{Sw}(\gamma_0, \dots, \gamma_n)$. In this case, the number of runs with the period a in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ is one less than the number of the letters b in $\text{Sw}(\gamma_1, \dots, \gamma_n)$, namely $N_\gamma(2) - 1$. Hence, in this case the sum of type-0 runs exponents equals

$$\sigma_0(\gamma) = 2(N_\gamma(2) - \text{odd}(n)).$$

Assume now that $\gamma_0 > 1$. Every run with the period a in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ equals a^{γ_0} or a^{γ_0+1} and is followed by the single letter b . Due to Lemma 4, every such run in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ corresponds to the letter a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Hence in this case we have $N_\gamma(1)$ runs with the period a .

By Lemma 4 each occurrence of a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$ preceded by b produces a run a^{γ_0+1} in $\text{Sw}(\gamma_0, \dots, \gamma_n)$, and each occurrence of a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$ not preceded by b produces a run a^{γ_0} in $\text{Sw}(\gamma_0, \dots, \gamma_n)$. Therefore, in computation of the sum of runs exponents, we count γ_0 for each a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$ and an additional 1 for each b . As in the previous case, for odd n , the rightmost b does not correspond to a run in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Therefore, in this case the sum of type-0 runs exponents equals

$$\sigma_0(\gamma) = \gamma_0 N_\gamma(1) + N_\gamma(2) - \text{odd}(n).$$

□

The next boundary case, strongly related to the case considered above, is the sum of exponents of runs with the period of the form x_1 .

Lemma 22 (Type 1).

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word. The sum of exponents of runs with the period x_1 in w equals:

$$\sigma'_1(\gamma) = \begin{cases} \left((N_\gamma(3) - 1) \cdot \left(2 + \frac{\gamma_0}{\gamma_0 + 1} \right) + \text{odd}(n) \cdot \left(2 + \frac{1}{\gamma_0 + 1} \right) \right) & \text{for } \gamma_1 = 1 \\ N_\gamma(2) \cdot \left(\gamma_1 + \frac{\gamma_0}{\gamma_0 + 1} \right) + (N_\gamma(3) - 1) + \text{odd}(n) \cdot \frac{\gamma_0 - 1}{\gamma_0 + 1} & \text{for } \gamma_1 > 1 \end{cases}. \quad (10)$$

Proof.

Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$. By definition we have $x_1 = a^{\gamma_0}b$. Therefore, the remainder of each internal run with the period x_1 has the length γ_0 .

Consider the 1-partition of w . By Lemma 4 occurrences of blocks of the form $a^{\gamma_0}b$ correspond to occurrences of letters a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$ and occurrences of blocks of the form a to occurrences of letters b in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Therefore, following similar argumentation as in proof of Lemma 21, we obtain the formula for the sum of exponents of internal runs with the period x_1 in w .

Let us now consider the rightmost run with the period x_1 in w . If n is even, w ends with $a \cdot a^{\gamma_0}b$ and this occurrence of x_1 does not correspond to a run in w . On the other hand, if n is odd, due to Proposition 16 w ends with $(a^{\gamma_0}b)^{\gamma_1+1}a$. Such a suffix corresponds to a run with the total part of exponent equal $\gamma_1 + 1$ and the remainder a , and we should include it in our formula. \square

The sum of exponents of runs with the period $(x_1)^k x_0$ for $1 \leq k < \gamma_1$ follows from Lemma 19. As a final step of investigation we count the sum of exponents of type- n runs.

Lemma 23 (Type n).

Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word. The sum of exponents of runs of type n in w is given by the formula:

$$\sigma_n(\gamma) = \begin{cases} 0 & \text{for } \gamma_n = 1 \\ \gamma_n + \frac{q_{n-1}}{q_n} & \text{for } \gamma_n > 1 \end{cases} \quad (11)$$

Proof.

We have $w = (x_n)^{\gamma_n} x_{n-1}$. Therefore, for $\gamma_n = 1$ there is no run of type n in w . On the other hand, for $\gamma_n > 1$, w contains only one run of type n . Its generator $-x_n -$ repeats undivided γ_n times. Moreover, since x_{n-1} is a prefix of x_n , the total exponent of a equals $\gamma_n + \frac{q_{n-1}}{q_n}$. \square

Now we can combine the formulas from equations (6), (7), (9), (10) and (11) and obtain the formula from Theorem 17.

4.3 Algorithm

The formulas from equations (6), (7), (9), (10) and (11) lead to simple and efficient algorithm for computation of the sum of runs exponents in any standard word. Its time complexity depends only on the coefficients of the directive sequence, which is the compressed representation of a considered word.

Theorem 24.

Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma)$ be a standard word. The sum of exponents of maximal repetitions in w can be computed in time $O(\|\gamma\|)$, where $\|\gamma\| = \gamma_0 + \gamma_1 + \dots + \gamma_n$.

Proof.

Observe that, by equations (6), (7), (9), (10) and (11), the value of each formula $\sigma_i(\gamma)$ depends only on coefficients of γ and the values of $N_\gamma(i + 1)$, $N_\gamma(i + 2)$, q_i and q_{i-1} .

Therefore, we can iterate through all types of runs from 0 to n computing the value of $\sigma_i(\gamma)$ and simultaneously updating the values of $N_\gamma(i+1)$, $N_\gamma(i+2)$, q_i and q_{i-1} . See Algorithm 1 for details.

The main loop of presented algorithm (lines 8-13) performs $n+1$ iterations. The most time consuming part of each iteration is the computation of the sum of exponents of maximal repetitions with the period $(x_i)^k x_{i-1}$ (line 10), namely the component

$$\sum_{k=1}^{\gamma_i-1} \left(2 + \frac{q_i - 2}{k \cdot q_i + q_{i-1}} \right).$$

It can be done in $O(\gamma_i)$ time. Hence, the time complexity of the whole algorithm is $O(\|\gamma\|)$, where $\|\gamma\| = \gamma_0 + \gamma_1 + \dots + \gamma_n$. \square

Algorithm 1: Sum-Of-Exponents($\text{Sw}(\gamma)$)

Input: $\gamma = (\gamma_0, \dots, \gamma_n)$
Output: $\sigma(\gamma)$

```

1 result  $\leftarrow$  0;
2  $N_\gamma(n+1) \leftarrow$  1;
3  $N_\gamma(n+2) \leftarrow$  0;
4  $q_0 \leftarrow$  1;
5  $q_{-1} \leftarrow$  1;

6 for  $i := 1$  to  $n$  do
7    $(q_{i+1}, q_i) \leftarrow (\gamma_i \cdot q_i, q_i)$ ;

8 for  $i := n$  downto 0 do
9   compute  $\sigma'_i(\gamma)$ ; // runs with period  $x_i$ ;
10  compute  $\sigma''_i(\gamma)$ ; // runs with period  $(x_i)^k x_{i-1}$ ;
11  result  $\leftarrow$  result +  $\sigma'_i(\gamma) + \sigma''_i(\gamma)$ ;
12   $(q_i, q_{i-1}) \leftarrow (q_{i-1}, q_i - \gamma_{i-1} \cdot q_{i-1})$ ;
13   $(N_\gamma(i), N_\gamma(i+1)) \leftarrow (\gamma_i \cdot N_\gamma(i) + N_\gamma(i+1), N_\gamma(i))$ ;

14 return result;
```

Final remarks

The aim of this paper was to study problems related to repetitions in standard Sturmian words – one of the most thoroughly investigated class of strings in combinatorics of words. We presented the formulas for the sum of exponents of maximal repetitions in any standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ that depend only on its compressed representation (the directive sequence). We proposed also an algorithm based on those formulas that computes the sum of runs exponents in any standard word in linear time with respect to the (total) size of the directive sequence, i.e. in time $O(\|\gamma\|)$, where $\|\gamma\| = \gamma_0 + \gamma_1 + \dots + \gamma_n$.

The notion of total run length (TRL) proposed in [12] can be considered similarly. To obtain the formulas for the total run length of a standard word we can use modified formulas for the sum of runs exponents. We only needed to multiply the total part of each exponent by the length of related period (either q_i or $k \cdot q_i + q_{i-1}$) and remove

the denominator from its fractional part. The described change could be also taken into account in the presented algorithm.

In the case of the total run length computation, the component

$$\sum_{k=1}^{\gamma_i-1} \left(2 + \frac{q_i - 2}{k \cdot q_i + q_{i-1}} \right)$$

of equation 7 has the form

$$\sum_{k=1}^{\gamma_i-1} \left((k+1)q_i + q_{i-1} - 2 \right).$$

The above formula is a sum of an arithmetic progression, hence it can be simplified as

$$(\gamma_i - 1) \frac{(\gamma_i + 2)q_i + 2q_{i-1} - 4}{2}.$$

Therefore, in each iteration of the main loop of the modified algorithm, we have to compute the value of a single arithmetic formula and update the values of $N_\gamma(i+1)$, $N_\gamma(i+2)$, q_i and q_{i-1} . This way we obtain the algorithm computing the total run length of any standard word $\text{Sw}(\gamma)$ in time $O(|\gamma|)$, where $|\gamma|$ denotes the length of the directive sequence.

References

1. J. ALLOUCHE AND J. SHALLIT: *Automatic Sequences. Theory, Applications, Generalizations.*, Cambridge University Press, 2003.
2. P. BATURO, M. PIĄTKOWSKI, AND W. RYTTER: *The Number of Runs in Standard Sturmian Words*. *Electronic Journal of Combinatorics*, 20(1) 2013.
3. P. BATURO AND W. RYTTER: *Compressed string-matching in standard Sturmian words*. *Theoretical Computer Science*, 410(30–32) 2009, pp. 2804–2810.
4. J. BERSTEL: *Sturmian and Episturmian words: a survey of some recent results*, in *Proceedings of the 2nd international conference on Algebraic informatics*, vol. 4728 of *Lecture Notes in Computer Science*, Springer, 2007, pp. 23–47.
5. J. BERSTEL, A. LAUVE, C. REUTENAUER, AND F. SALIOLA: *Combinatorics on Words: Christoffel Words and Repetitions in Words*, CRM monograph series, Providence, R.I: American Mathematical Society, 2009.
6. M. CROCHEMORE AND L. ILIE: *Analysis of maximal repetitions in strings*, in *Proceedings of the 32nd International Conference on Mathematical Foundations of Computer Science*, vol. 4708 of *Lecture Notes in Computer Science*, Springer, 2007, pp. 465–476.
7. M. CROCHEMORE, L. ILIE, AND L. TINTA: *Towards a solution of the "runs" conjecture*, in *Proceedings of the 19th annual symposium on Combinatorial Pattern Matching*, vol. 5029 of *Lecture Notes in Computer Science*, Springer, 2008, pp. 290–302.
8. D. DAMANIK AND D. LENZ: *The index of sturmian sequences*. *European Journal of Combinatorics*, 23(1) 2002, pp. 23–29.
9. D. DAMANIK AND D. LENZ: *Powers in Sturmian sequences*. *European Journal of Combinatorics*, 24(4) 2003, pp. 377–390.
10. F. FRANEK, R. J. SIMPSON, AND W. F. SMYTH: *The maximum number of runs in a string*, in *Proceedings of 14th Australian Workshop on Combinatorial Algorithms*, 2003, pp. 26–35.
11. F. FRANEK AND Q. YANG: *An asymptotic lower bound for the maximal number of runs in a string*. *International Journal of Foundations of Computer Science*, 19(1) 2008, pp. 195–203.
12. A. GLEN AND J. SIMPSON: *The total run length of a word*. arXiv:1301.6568, 2013.
13. R. KOLPAKOV AND G. KUCHEROV: *On the sum of exponents of maximal repetitions in a word*, Tech. Rep. 99-R-034, LORIA, 1999.

14. K. KUSANO, W. MATSUBARA, A. ISHINO, H. BANNAI, AND A. SHINOBARA: *New lower bound for the maximum number of runs in a string*. Computing Research Repository, abs/0804.1214 2008.
15. M. LOTHAIRE: *Algebraic Combinatorics on Words*, vol. 90 of Encyclopedia of mathematics and its application, Cambridge University Press, 2002.
16. M. PIĄTKOWSKI: *Stringology applets* .
<http://www.mat.umk.pl/~martinp/stringology/applets>.
17. M. SCIORTINO AND L. ZAMBONI: *Suffix automata and standard Sturmian words*, in Proceedings of the 11th International Conference on Developments in Language Theory, vol. 4588 of Lecture Notes in Computer Science, Springer, 2007, pp. 382–398.
18. J. SHALLIT: *Characteristic words as fixed points of homomorphisms*, Tech. Rep. CS-91-72, University of Waterloo, Department of Computer Science, 1991.