

Computing the Number of Cubic Runs in Standard Sturmian Words

Marcin Piątkowski^{2*} and Wojciech Rytter^{1,2}

¹ Department of Mathematics, Computer Science and Mechanics,
University of Warsaw, Warsaw, Poland
rytter@mimuw.edu.pl

² Faculty of Mathematics and Informatics,
Nicolaus Copernicus University, Toruń, Poland
martinp@mat.umk.pl

Abstract. The *standard Sturmian* words are extensively studied in combinatorics of words. They are enough complicated to have many interesting properties and at the same time they are highly compressible. In this paper we present compact formulas for the number $\rho^{(3)}$ of cubic runs in any standard word. We show also that

$$\lim_{|w| \rightarrow \infty} \frac{\rho^{(3)}(w)}{|w|} = \frac{5\Phi + 3}{13\Phi + 9} \approx 0.36924841$$

and present the sequence of strictly growing standard words achieving this limit. The exact asymptotic ratio is here irrational, contrary to the situation of squares and runs in the same class of words. Furthermore we design an efficient algorithm computing the number of cubic runs in standard words in linear time with respect to the size of the compressed representation (recurrences) describing the word. The explicit size of the word can be exponential with respect to this representation. This is yet another example of a very fast computation on highly compressible texts.

Keywords: standard Sturmian words, repetitions, cubic runs, algorithms

1 Introduction

Repetitions in strings are important in combinatorics on words and many practical applications, see for instance [6], [11], [19] and [20]. The structure of repetitions is almost completely understood for the class of Fibonacci words, see [15], [17], [24], however it is not well understood for general words.

Runs are repetitions in which the period repeats at least twice. Highly repetitive segments, in which the repetitions ratio is at least 3, called the *cubic runs*, were introduced and studied in [10].

We say that a 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 will be 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 finite 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 $|u|$). *Cubic runs*

* The research supported by Ministry of Science and Higher Education of Poland, grant N N206 258035.

Figure 1. The structure of repetitions in the word $\text{Sw}(1, 2, 1, 3, 1)$. There are 19 runs and 4 cubic runs (marked in **bold**).

Example 1. Let $w = ababaababababababababababababababab$.

There are 5 runs with the period $|a|$:

$$\begin{aligned} w[5..6] &= a^2, & w[12..13] &= a^2, & w[19..20] &= a^2, \\ w[26..27] &= a^2, & w[31..32] &= a^2, \end{aligned}$$

5 runs with the period $|ab|$ (including 3 cubic runs):

$$\begin{aligned} w[1..5] &= (ab)^2a, & w[6..12] &= (ab)^3a, & w[13..19] &= (ab)^3a, \\ w[20..26] &= (ab)^3a, & w[27..31] &= (ab)^2a, \end{aligned}$$

4 runs with the period $|aba|$:

$$\begin{aligned} w[3..8] &= (aba)^2, & w[10..15] &= (aba)^2, \\ w[17..22] &= (aba)^2, & w[24..29] &= (aba)^2, \end{aligned}$$

4 runs with the period $|ababa|$:

$$\begin{aligned} w[1..10] &= (ababa)^2, & w[8..17] &= (ababa)^2, \\ w[15..24] &= (ababa)^2, & w[22..33] &= (ababa)^2ab, \end{aligned}$$

and 1 (cubic) run with the period $|ababaab|$:

$$w[1..31] = (ababaab)^4aba.$$

All together we have 19 runs and 4 cubic runs, see Figure 1 for comparison.

Denote by $\rho(w)$ and $\rho^{(3)}(w)$ the number of runs and cubic runs in the word w , and by $\rho(n)$ and $\rho^{(3)}(n)$ the maximal number of runs and cubic runs in the words of length n respectively. The most interesting and open conjecture about maximal repetitions is:

$$\rho(n) < n.$$

In 1999 Kolpakov and Kucherov (see [16]) showed that the number $\rho(w)$ of runs in a string w is $O(|w|)$, but the exact multiplicative constant coefficient is still unknown. The best known results related to the value of $\rho(n)$ are

$$0.944575712 n \leq \rho(n) \leq 1.029 n.$$

The upper bound is by [8], [9] and the lower bound is by [13], [14], [18], [27]. The best known results related to $\rho^{(3)}(n)$ are (due to [10]):

$$0.41 n \leq \rho^{(3)}(n) \leq 0.5 n.$$

For the class \mathcal{S} of standard Sturmian words there are known exact formulas for the number of runs and squares and their asymptotic behavior, see [2] and [22] for details. In this case we have

$$\lim_{n \rightarrow \infty} \frac{\rho(n)}{n} = 0.8.$$

This paper is devoted to the investigation of the structure of cubic runs in standard Sturmian words. We present the exact recurrence formulas for the number $\rho^{(3)}(w)$. Next we derive the algorithm computing $\rho^{(3)}(w)$ for any word $w \in \mathcal{S}$ in linear time with respect to the compressed representation of w , hence logarithmic time with respect to the length of the whole word w . We show also, that for any standard word w , we have

$$\rho^{(3)}(w_k) \leq 0.36924841 |w|,$$

and construct the sequence $\{w_k\}$ of strictly growing standard words, for which we have

$$\lim_{k \rightarrow \infty} \frac{\rho^{(3)}(w_k)}{|w_k|} = \frac{5\Phi + 3}{13\Phi + 9} \approx 0.36924841.$$

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], [4], [5], [7], [19], [25], [26], [28] and references therein. They have very compact representations in terms of sequences of integers, which has many algorithmic consequences.

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 . 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 grammar is proportional to the length of the directive sequence.

3 Morphic reduction of standard words

The recurrent definition of standard words leads to the simple characterization by the 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:

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

Lemma 4. *For $0 \leq i \leq n$ the morphism h_i transforms a standard word into another standard word, and we have:*

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

Proof. We will prove the above lemma by the induction on the length of the directive sequence. Recall that the standard word given by the empty directive sequence is a . For $|\gamma| = 1$ we have, by definition of standard words and the morphism h_n ,

$$\text{Sw}(\gamma_n) = a^{\gamma_n} b = h_n(a).$$

Assume now that $|\gamma| = k \geq 2$ and for directive sequences shorter than k the thesis holds. We have then:

$$\begin{aligned} \text{Sw}(\gamma_i, \dots, \gamma_n) &= [\text{Sw}(\gamma_i, \dots, \gamma_{n-1})]^{\gamma_n} \cdot \text{Sw}(\gamma_i, \dots, \gamma_{n-2}) \\ &\stackrel{\text{ind.}}{=} \left[h_i(\text{Sw}(\gamma_{i+1}, \dots, \gamma_{n-1})) \right]^{\gamma_n} \cdot h_i(\text{Sw}(\gamma_{i+1}, \dots, \gamma_{n-2})) \\ &= h_i([\text{Sw}(\gamma_{i+1}, \dots, \gamma_{n-1})]^{\gamma_n} \cdot \text{Sw}(\gamma_{i+1}, \dots, \gamma_{n-2})) \\ &= h_i(\text{Sw}(\gamma_{i+1}, \dots, \gamma_n)), \end{aligned}$$

which concludes the proof. □

Remark 5. As a direct conclusion from Lemma 4 we have that the standard word corresponding to the directive sequence $\gamma = (\gamma_0, \gamma_1, \dots, \gamma_n)$ is given as:

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

The inverse morphism h_i^{-1} can be seen as a reduction of the word $\text{Sw}(\gamma_i, \dots, \gamma_n)$ to the word $\text{Sw}(\gamma_{i+1}, \dots, \gamma_n)$ and allows us to reduce the computation of cubic runs in $\text{Sw}(\gamma_i, \dots, \gamma_n)$ to the same computation in $\text{Sw}(\gamma_{i+1}, \dots, \gamma_n)$.

Denote by $|w|_a$ the number of occurrences of the letters a in the word w . We define the function, which will be useful in the rest of this paper. For a directive sequence $\gamma = (\gamma_0, \dots, \gamma_n)$ and an integer $0 \leq k \leq n + 1$ we define

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

Moreover, for $k > n + 1$, we define $N_\gamma(k) = 0$.

Remark 6. As a direct conclusion from the above definition, the equation (1) and the equation (2) we have that for $0 \leq k \leq n$ the numbers $N_\gamma(k)$ satisfy:

$$N_\gamma(k) = \gamma_k N_\gamma(k + 1) + N_\gamma(k + 2). \tag{5}$$

Example 7. Let $\gamma = (1, 2, 1, 3, 1)$ be a directive sequence. We have then

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

Remark 8. In case of Fibonacci words the numbers $N_\gamma(k)$ are Fibonacci numbers:

$$N_\gamma(k) = |F_{n-k-1}| = f_{n-k-1}. \tag{6}$$

4 Formulas for the number of cubic runs

In this section we present and prove formulas for the number of cubic runs 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:

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

We begin with the characterization of possible periods of cubic runs in standard words. The following lemma is a consequence of the very special structure of subword graphs (especially their compacted versions) of those words. See [3] and [25] for more information.

Lemma 9. *The period of each cubic run in the standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ is of the form x_i , where x_i 's are as in equation (1).*

To prove the above lemma it is sufficient to show that no factor of the word $\text{Sw}(\gamma_0, \dots, \gamma_n)$, that does not satisfy the condition given there, could be the generator of a cubic run. We can use similar argumentation as in proof of Theorem 1 in [12]. The details are omitted in this version.

The main idea of the computation of cubic runs in a standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ is the partition of them into three separate categories depending on the length of their periods. We say that a cubic run is:

- short** – if it has the period of the form a or $a^k b$,
medium – if it has the period of the form x_2 ,
large – if it has the period of the form x_i , for $i > 2$.

Denote by $\rho_S^{(3)}(w)$, $\rho_M^{(3)}(w)$ and $\rho_L^{(3)}(w)$ the number of short, medium and large cubic runs in the word w , respectively. We will consider each type separately.

Example 10. Recall the word $w = \text{Sw}(1, 2, 1, 3, 1)$ from Example 1. We have:

- 3 short cubic runs (period ab),
- no medium cubic run,
- 1 large cubic run (period $ababaab$),

see Figure 1 for comparison.

4.1 Short runs

We start with the computation of the *short* cubic runs. These are the cubic runs with the periods of the form a or $a^k b$. Their number depends on the values of γ_0 and γ_1 .

Lemma 11. *The number $\rho_{S_1}^{(3)}$ of cubic runs with the period a in the standard word $w = \text{Sw}(\gamma_0, \gamma_1, \dots, \gamma_n)$ equals:*

$$\rho_{S_1}^{(3)}(w) = \begin{cases} 0 & \text{for } \gamma_0 = 1 \\ N_\gamma(2) - \text{odd}(n) & \text{for } \gamma_0 = 2 \\ N_\gamma(1) & \text{for } \gamma_0 > 2 \end{cases} \quad (7)$$

Proof. First assume that $\gamma_0 > 2$. Every cubic 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 cubic run in $\text{Sw}(\gamma_0, \gamma_1, \dots, \gamma_n)$ corresponds to the letter a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Hence in this case we have $N_\gamma(1)$ cubic runs with the period a .

Assume now that $\gamma_0 = 2$. In this case the word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ consists of the blocks of the two types: aab and $aaab$. Only the blocks of the second type include the cubic run with the period a . Due to Lemma 4 every such cubic 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 cubic runs equals the number of blocks ba in $\text{Sw}(\gamma_0, \dots, \gamma_n)$.

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 cubic 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 $N_\gamma(2)$. 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 cubic 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 the word $\text{Sw}(\gamma_1, \dots, \gamma_n)$, namely $N_\gamma(2) - 1$.

Finally assume that $\gamma_0 = 1$. In this case the word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ consists of the blocks of the two types: ab and aab . None of them includes a cubic run with the period a , and this completes the proof. \square

Lemma 12. *The number $\rho_{S_2}^{(3)}$ of cubic runs with the period $a^k b$ in the standard word $w = \text{Sw}(\gamma_0, \gamma_1, \dots, \gamma_n)$ equals:*

$$\rho_{S_2}^{(3)}(w) = \begin{cases} 0 & \text{for } \gamma_1 = 1 \\ N_\gamma(3) - \text{even}(n) & \text{for } \gamma_1 = 2 \\ N_\gamma(2) & \text{for } \gamma_1 > 2 \end{cases} \quad (8)$$

Proof. Notice that, due to equation (2) and Lemma 4, cubic runs with the periods of the form $a^k b$ in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ correspond to cubic runs with the period a in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Similar reasoning as above establishes the desired formula. \square

4.2 Medium runs

Recall that a cubic run is called *medium* if it has the period of the form x_2 . Observe that medium cubic runs appear in standard words generated by directive sequences of the length at least 3. We have to consider two cases: the directive sequences of the length 3 and the longer directive sequences. The values of γ_0 and γ_1 does not affect the number of medium cubic runs, hence to simplify the calculations we can assume in further proofs that $\gamma_0 = \gamma_1 = 1$.

We start with counting medium runs in standard words generated by directive sequences of the length greater than 3.

Lemma 13. *Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ be a standard word and $n \geq 3$. The number of medium cubic runs in w equals:*

$$\rho_M^{(3)}(w) = \begin{cases} N_\gamma(4) - 1 & \text{for } \gamma_2 = 1 \\ N_\gamma(3) & \text{for } \gamma_2 \geq 2 \end{cases} \quad (9)$$

Proof. We start with the assumption that $\gamma_2 > 2$. In this case every factor of the form $x_3 = x_2^{\gamma_2} x_1$ includes one cubic runs with the period x_2 . Hence the number of such cubic runs equals the number of factors x_3 in $\text{Sw}(\gamma_0, \dots, \gamma_n)$, namely $N_\gamma(3)$ (due to Lemma 4).

Assume now that $\gamma_2 = 2$. The word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ can be represented as a sequence of concatenated words x_3 and x_2 and has the form:

$$x_3^{\alpha_1} x_2 x_3^{\alpha_2} x_2 \cdots x_3^{\alpha_s} x_2 x_3 \quad \text{or} \quad x_3^{\beta_1} x_2 x_3^{\beta_2} x_2 \cdots x_3^{\beta_s} x_2.$$

Observe that $x_3 = x_2 x_2 x_1$ and every occurrence of x_3 in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ either follows some occurrence of x_2 or is followed by some occurrence of x_2 . In the first case we have $x_2 \cdot x_3 = x_2 \cdot x_2 x_2 x_1$ and there is a cubic run with period x_2 . In the second case we have $x_3 \cdot x_2 = x_2 x_2 x_1 \cdot x_2$, and there is also a cubic run with period x_2 , since x_1 is a prefix of x_2 . Therefore the number of medium cubic runs in this case equals the number of the factors x_3 in $\text{Sw}(\gamma_0, \dots, \gamma_n)$, namely $N_\gamma(3)$.

Finally assume that $\gamma_2 = 1$. The word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ can be represented as a sequence of concatenated words x_3 and x_4 and has the form:

$$x_4^{\alpha_1} x_3 x_4^{\alpha_2} x_3 \cdots x_4^{\alpha_s} x_3 x_4 \quad \text{or} \quad x_4^{\beta_1} x_3 x_4^{\beta_2} x_3 \cdots x_4^{\beta_s} x_3.$$

We have $x_3 = x_2 x_1$ and $x_4 = x_2 x_1 \cdots x_2 x_1 \cdot x_2$. Therefore only the last one occurrence of x_4 in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ does not correspond to a cubic run with the period x_2 and we have $N_\gamma(4) - 1$ such cubic runs in this case. This completes the proof. \square

Lemma 14. *The number of medium cubic runs in the word $w = \text{Sw}(\gamma_0, \gamma_1, \gamma_2)$ equals:*

$$\rho_M^{(3)}(w) = \begin{cases} 1 & \text{for } \gamma_2 > 2 \\ 0 & \text{for } \gamma_2 \leq 2 \end{cases}. \quad (10)$$

Proof. We have $\text{Sw}(\gamma_0, \gamma_1, \gamma_2) = x_2^{\gamma_2} x_1$. Hence there is only one medium run (with the period x_2) if $\gamma_2 > 2$ and no medium run otherwise. \square

4.3 Large runs

Recall that a cubic run is called *large* if it has the period of the form x_i for $i > 2$, where x_i are as in the equation (1). We reduce the problem of counting large cubic runs to the one for counting medium cubic runs, using morphic representation of standard words introduced in previous section.

Let h be a morphism, $v = a_1 a_2 \cdots a_k$ be a word of the length k and let $w = h(v)$. The morphism h defines the partition of w into segments $h(a_1), h(a_2), \dots, h(a_k)$. These segments are called the *h -blocks*. We say that a factor x of the word w is *synchronized* with the morphism h in w if and only if each occurrence of x in w starts at the beginning of some h -block and ends at the end of some h -block. Observe that every factor in w that is synchronized with h corresponds to some factor in v , hence the morphism h preserves the structure of the factors that are synchronized with it.

Example 15. Let $w = \text{Sw}(1, 2, 1, 3, 1)$ and $v = \text{Sw}(2, 1, 3, 1)$ be standard words and h_0 be a morphism defined as:

$$h_0 : \begin{cases} a \longrightarrow ab \\ b \longrightarrow a \end{cases}.$$

Recall that

$$\begin{aligned} \text{Sw}(1, 2, 1, 3, 1) &= h_0(\text{Sw}(2, 1, 3, 1)), \\ \text{Sw}(1, 2, 1, 3, 1) &= ababaabababaabababaabababaabababaab, \\ \text{Sw}(2, 1, 3, 1) &= aabaaabaaabaaabaaba. \end{aligned}$$

The factors $w[6..8] = aba$ and $w[13..17] = abaab$ are not synchronized with h_0 , because both of them end in the middle some h_0 -block. From the other hand, the factor $w[22..28] = ababaab$ and all its occurrences in w (namely $w[1..7]$, $w[8..14]$ and $w[15..21]$) start at the beginning of some h_0 -block and end at the end of some h_0 -block. Hence the factor $w[22..28]$ is synchronized with the morphism h_0 . Moreover it corresponds with the factor $v[13..16] = aaba$, see Figure 2 for comparison.

Lemma 16. *The periods of large cubic runs in the standard word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ are synchronized with the morphism h_0 .*

Proof. Let h_0 be the morphism defined as

$$h_0 : \begin{cases} a \longrightarrow a^{\gamma_0} b \\ b \longrightarrow a \end{cases}.$$

Figure 2. The factors aba and $ababa$ do not synchronize with the morphism h_0 in the word $\text{Sw}(1, 2, 1, 3, 1)$, while the factor $ababaab$ (in fact the period of the large cubic run) is synchronized with h_0 and corresponds to the factor $aaba$ in the word $\text{Sw}(2, 1, 3, 1)$.

Due to Lemma 4 we have $\text{Sw}(\gamma_0, \dots, \gamma_n) = h_0(\text{Sw}(\gamma_1, \dots, \gamma_n))$. Moreover, h_0 determines the partition of $\text{Sw}(\gamma_0, \dots, \gamma_n)$ into h_0 -blocks of the form $a^{\gamma_0}b$ and a , see Figure 2 for the partition of $\text{Sw}(1, 2, 1, 3, 1)$.

Recall that the period of each large cubic run in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ is of the form x_i , where $i \geq 3$. By definition of standard words the factor x_i starts with $a^{\gamma_0}b$, hence at the beginning of some h_0 -block.

For odd $i \geq 3$ the subword x_i ends with $x_1 = a^{\gamma_0}b$, hence at the end of some h_0 -block, and is obviously synchronized with h_0 .

For even $i \geq 3$ the factor x_i ends with

$$x_3 \cdot x_2 = x_2^{\gamma_2} x_1 \cdot x_1^{\gamma_1} x_0 = x_2^{\gamma_2} \cdot (a^{\gamma_0}b)^{\gamma_1+1} a.$$

First assume that x_i is followed by the block $a^{\gamma_0}b$. The single letter a at the end of x_i is then the whole h_0 -block and x_i is synchronized with the morphism h_0 .

Assume now that x_i ends with $(a^{\gamma_0}b)^{\gamma_1+1}a$ and is followed by $(a^{\gamma_0-1}b)$, namely it ends in the middle of some h_0 -block. In this case we have the occurrence of the factor $(a^{\gamma_0}b)^{\gamma_1+2}$ in $\text{Sw}(\gamma_0, \dots, \gamma_n)$, which is reduced by the morphism h_0^{-1} to the factor $a^{\gamma_1+2}b$ in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. By definition, the standard word $\text{Sw}(\gamma_1, \dots, \gamma_n)$ can include only the blocks of the two types: the short block $-a^{\gamma_1}b$ and the long block $-a^{\gamma_1+1}b$, hence we have the contradiction and the proof is complete. \square

The following lemma, which is a direct conclusion from Lemma 16, allows us to reduce the problem of counting large cubic runs in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ to counting large and medium cubic runs in $\text{Sw}(\gamma_1, \dots, \gamma_n)$.

Lemma 17. *Let $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ and $v = \text{Sw}(\gamma_1, \dots, \gamma_n)$ be standard words. The number of large cubic runs in w is given by the recurrence*

$$\rho_L^{(3)}(w) = \rho_L^{(3)}(v) + \rho_M^{(3)}(v).$$

Proof. Lemma 16 implies that the morphism defined in the equation (2) preserves the structure of long cubic runs in standard words. Recall that the word $\text{Sw}(\gamma_0, \dots, \gamma_n)$ is reduced by h_0^{-1} to the word $\text{Sw}(\gamma_1, \dots, \gamma_n)$. Therefore, every large cubic run α in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ corresponds to some cubic run β in $\text{Sw}(\gamma_1, \dots, \gamma_n)$.

Due to Lemma 9 the period of the cubic run α is of the form x_i , where $i \geq 3$. The corresponding cubic run β is either large (for $i > 3$) or medium (for $i = 3$). Hence to compute all large cubic runs in $\text{Sw}(\gamma_0, \dots, \gamma_n)$ it is sufficient to compute all large and medium cubic runs in $\text{Sw}(\gamma_1, \dots, \gamma_n)$. \square

The following theorem summarizes all the formulas developed above.

Theorem 18. *Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence, $w = \text{Sw}(\gamma_0, \dots, \gamma_n)$ and $w_i = \text{Sw}(\gamma_i, \dots, \gamma_n)$, for $0 \leq i \leq n$, be standard words. The number of cubic runs in w is given as:*

$$\rho^{(3)}(w) = \rho_{S_1}^{(3)}(w) + \rho_{S_1}^{(3)}(w) + \sum_{i=0}^{n-2} \rho_M^{(3)}(w_i). \quad (11)$$

Proof. The thesis of the theorem follows by combining the formulas (7), (8), the formula (9) repeated $n - 2$ times, and finally the formula (10). \square

Example 19. Consider a directive sequence $\gamma = (1, 2, 1, 3, 1)$. We compute the number of cubic runs in $\text{Sw}(1, 2, 1, 3, 1)$ using the formulas mentioned above. We have:

$$\begin{aligned} \text{short cubic runs with period } a: & \quad 0 \\ \text{short cubic runs with period } a^k b: & \quad |aaaba|_a - 1 = 3 \\ \text{medium cubic runs:} & \quad |ab|_a - 1 = 0 \\ \text{large cubic runs:} & \quad \rho_M^{(3)}(2, 1, 3, 1) + \rho_M^{(3)}(1, 3, 1) = |ab|_a + 0 = 1 \end{aligned}$$

Altogether there are 4 cubic runs, see Example 1 and Figure 1 for comparison.

4.4 Algorithm for computation of the number of cubic runs

The formulas investigated above allow us to develop an efficient algorithm computing the number of cubic runs in any standard Sturmian word.

Theorem 20. *Let $\gamma = (\gamma_0, \dots, \gamma_n)$ be a directive sequence and $\text{Sw}(\gamma)$ be a standard word. We can count the number of cubic runs in $\text{Sw}(\gamma)$ in linear time with respect to the length of the directive sequence $|\gamma|$ (logarithmic time with respect to the length of the whole word $|\text{Sw}(\gamma)|$).*

Proof. The formulas (7), (8), (9) and (10) for the number of cubic runs in a standard word $\text{Sw}(\gamma)$ depend directly on the components of the directive sequence γ and the numbers $N_\gamma(k)$. We can compute the numbers $N_\gamma(n), N_\gamma(n-1), \dots, N_\gamma(1)$ by consecutive iteration of the equation (5). In each step i of the computation we remember the number of cubic runs related to the value of the γ_i . The number of iterations performed by the algorithm correspond directly to the length of the directive sequence, hence it has the time complexity $O(|\gamma|)$. See Algorithm 1 for more details. \square

5 Asymptotic behaviour of the number of cubic runs

This section is devoted to the computation of the asymptotic limit

$$\lim_{|w| \rightarrow \infty} \frac{\rho^{(3)}(w)}{|w|} \quad (12)$$

for $w \in \mathcal{S}$.

Algorithm 1: Counting-Cubic-Runs(Sw(γ))

```

1  $(x, y, cr) \leftarrow (1, 0, 0)$ ;
2 if  $\gamma_n > 2$  then  $cr \leftarrow cr + 1$ ;
3 for  $i := n$  to 3 do
4    $(x, y) \leftarrow (\gamma_i \cdot x + y, x)$ ;
5   if  $\gamma_{i-1} \geq 2$  then  $cr \leftarrow cr + x$ ;
6   else  $cr \leftarrow cr + y - 1$ ;

7 if  $\gamma_1 = 2$  then
8    $cr \leftarrow cr + x$ ;
9   if  $n$  is even then  $cr \leftarrow cr - 1$ ;

10  $(x, y) \leftarrow (\gamma_2 \cdot x + y, x)$ ;
11 if  $\gamma_1 > 2$  then  $cr \leftarrow cr + x$ ;
12 if  $\gamma_0 = 2$  then
13    $cr \leftarrow cr + x$ ;
14   if  $n$  is odd then  $cr \leftarrow cr - 1$ ;

15 if  $\gamma_0 > 2$  then  $cr \leftarrow cr + \gamma_1 \cdot x + y$ ;
16 return  $cr$ ;

```

Theorem 21. Let $w = \text{Sw}(\gamma)$ be a standard word. Then we have

$$\rho^{(3)}(w) \leq 0.36924841 |w|.$$

Moreover there is infinite and strictly growing sequence of standard words achieving this asymptotic limit.

Proof. To prove the above theorem we will construct a directive sequence corresponding to a standard word for which the number of cubic runs will be maximal in relation to their length.

Let $\gamma = (\gamma_0, \gamma_1, \dots, \gamma_n)$ be a directive sequence and $w = \text{Sw}(\gamma)$ be a standard word. The number of cubic runs with the period of the form a in w corresponds directly to the value γ_0 , see equation (7). The word w consists of blocks of the two types: $a^{\gamma_0}b$ and $a^{\gamma_0+1}b$. For $\gamma_0 \geq 3$ every such block contains a desired cubic run, for $\gamma_0 = 2$ only the second type of blocks contains a short cubic run, and for $\gamma_0 < 2$ there is no short cubic run in w . Moreover, for $\gamma_0 > 3$ the number of considered cubic runs does not change while the length of the word increases significantly.

For $\gamma_0 = 2$ we have, by the equations (5) and (7):

$$|w| = (\gamma_1 + 1) N_\gamma(2) + 3 \gamma_1 N_\gamma(3) \quad \text{and} \quad \rho_{S_1}^{(3)}(w) = N_\gamma(2) \pm 1,$$

and for $\gamma_0 = 3$ we have:

$$|w| = (4 \gamma_1 + 1) N_\gamma(3) + 4 N_\gamma(3) \quad \text{and} \quad \rho_{S_1}^{(3)}(w) = \gamma_1 N_\gamma(2) + N_\gamma(3).$$

Therefore for the change of the value γ_0 from 2 to 3 the increase of $\rho_{S_1}^{(3)}$ (namely: $(\gamma_1 - 1) N_\gamma(2) + N_\gamma(3)$) is significant in relation to the increase of the length of the whole word (namely: $\gamma_1 N_\gamma(2) + \gamma_1 N_\gamma(3)$). Hence $\gamma_0 = 3$ is the optimal value. It does

not affect the number of cubic runs with longer periods, hence we assume in further discussion its optimal value.

The number of cubic runs with the period of the form $a^k b$ in w depends on the value of γ_1 , see equation (8). Similar argumentation as above shows that γ_1 must be greater than 1 and no more than 3. For $\gamma_1 = 2$ we have, by the equations (5) and (8):

$$|w| = (9\gamma_2 + 4)N_\gamma(3) + 9N_\gamma(4) \quad \text{and} \quad \rho_{S_2}^{(3)}(w) = N_\gamma(3) \pm 1,$$

and for $\gamma_1 = 3$ we have:

$$|w| = (13\gamma_2 + 4)N_\gamma(3) + 13N_\gamma(4) \quad \text{and} \quad \rho_{S_2}^{(3)}(w) = \gamma_2 N_\gamma(3) + N_\gamma(4).$$

Therefore the change of the value of γ_1 from 2 to 3 increases the number of cubic runs by: $(\gamma_2 - 1)N_\gamma(3) + N_\gamma(4) \pm 1$ and at the same time increases the length of the word by: $4\gamma_2 N_\gamma(3) + 4 N_\gamma(4)$. Hence we conclude that $\gamma_1 = 2$ is the optimal value.

The number of medium cubic runs in the word w corresponds to the value of γ_2 . It is easy to see that γ_2 must be at most 2, otherwise the length of the word increases significantly and the value $\rho_M^{(3)}(w)$ does not change.

For $\gamma_2 = 1$ we have, by the equations (5) and (9):

$$|w| = (13\gamma_3 + 9)N_\gamma(4) + 13N_\gamma(5) \quad \text{and} \quad \rho_M^{(3)}(w) = N_\gamma(4) - 1.$$

and for $\gamma_2 = 2$ we have:

$$|w| = (22\gamma_3 + 9)N_\gamma(4) + 22N_\gamma(5) \quad \text{and} \quad \rho_M^{(3)}(w) = \gamma_3 N_\gamma(4) + N_\gamma(5).$$

Therefore the change of the value of γ_2 from 1 to 2 increases the number of cubic runs by: $(\gamma_3 - 1) N_\gamma(4) + N_\gamma(5) + 1$ and at the same time increases the length of the word by: $9\gamma_3 N_\gamma(4) + 9 N_\gamma(5)$. Hence we conclude that $\gamma_2 = 1$ is the optimal value.

We compute large cubic runs by reduction of them to medium runs, see Lemma 17. Applying $n - 2$ times the above argumentation for the medium cubic runs we conclude that optimal value of $\gamma_3, \gamma_4, \dots, \gamma_{n-1}$ is also 1. Similarly for $\gamma_n > 1$ there is one additional long run whereas the length of the word increases more than two times.

We have shown above, that the maximal value of the quotient of the number of cubic runs to the length of the word is achieved by the standard words generated by directive sequences of the form $\gamma = (3, 2, 1, 1, \dots, 1)$. Now we are ready to compute the value of the asymptotic limit from the equation (12).

Let us consider a sequence of standard words:

$$w_k = (3, 2, \underbrace{1, 1, 1, \dots, 1}_k).$$

We have by definition of standard words and Remark 8:

$$|w_k| = 13 N_\gamma(3) + 9 N_\gamma(4) = 13 f_{k-1} + 9 f_{k-2},$$

and by Theorem 18 and Remark 8:

$$\rho^{(3)}(w_k) = 5 N_\gamma(3) + 3 N_\gamma(4) - k \pm 1 = 5 f_{k-1} + 3 f_{k-2}.$$

We have also that:

$$\lim_{n \rightarrow \infty} \frac{f_n}{f_{n-1}} = \Phi \approx 1.61803390,$$

hence

$$\lim_{k \rightarrow \infty} \frac{\rho^{(3)}(w_k)}{|w_k|} \approx \frac{5\Phi + 3}{13\Phi + 9} \approx 0.36924841,$$

and this completes the proof. □

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 Sturmian words*, in Proceedings of the 13th international conference on Implementation and Applications of Automata, vol. 5148 of Lecture Notes in Computer Science, Springer, 2008, pp. 252–261.
3. P. BATURO, M. PIĄTKOWSKI, AND W. RYTTER: *Usefulness of directed acyclic subword graphs in problems related to standard Sturmian words*. International Journal of Foundations of Computer Science, 20(6) 2009, pp. 1005–1023.
4. P. BATURO AND W. RYTTER: *Compressed string-matching in standard Sturmian words*. Theoretical Computer Science, 410(30–32) 2009, pp. 2804–2810.
5. 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.
6. J. BERSTEL AND J. KARHUMAKI: *Combinatorics on words: a tutorial*. Bulletin of the EATCS, 79 2003, pp. 178–228.
7. 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.
8. 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.
9. 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.
10. M. CROCHEMORE, C. S. ILIOPOULOS, M. KUBICA, J. RADOSZEWSKI, W. RYTTER, AND T. WALEN: *On the maximal number of cubic runs in a string*, in Proceedings of the International Conference on Implementation and Applications of Automata, 2010, pp. 227–238.
11. M. CROCHEMORE AND W. RYTTER: *Jewels of Stringology: text algorithms*, World Scientific, 2003.
12. D. DAMANIK AND D. LENZ: *Powers in Sturmian sequences*. European Journal of Combinatorics, 24(4) 2003, pp. 377–390.
13. 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.
14. 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.
15. C. S. ILIOPOULOS, D. MOORE, AND W. F. SMYTH: *A characterization of the squares in a Fibonacci string*. Theoretical Computer Science, 172(1–2) 1997, pp. 281–291.
16. R. KOLPAKOV AND G. KUCHEROV: *Finding maximal repetitions in a word in linear time*, in Proceedings of the 40th Annual Symposium on Foundations of Computer Science, IEEE Computer Society, 1999, pp. 596–604.
17. R. M. KOLPAKOV AND G. KUCHEROV: *On maximal repetitions in words*, in Proceedings of 12th International Symposium on Fundamentals of Computation Theory, vol. 1684 of Lecture Notes in Computer Science, Springer, 1999, pp. 374–385.
18. 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.
19. M. LOTHAIRE: *Algebraic Combinatorics on Words*, vol. 90 of Encyclopedia of mathematics and its application, Cambridge University Press, 2002.
20. M. LOTHAIRE: *Applied Combinatorics on Words*, vol. 105 of Encyclopedia of Mathematics and its Application, Cambridge University Press, 2005.
21. M. PIĄTKOWSKI: *Stringological applets* .
<http://www.mat.umk.pl/~martinp/stringology/applets>.
22. M. PIĄTKOWSKI AND W. RYTTER: *Asymptotic behaviour of the maximal number of squares in standard Sturmian words*, in Proceedings of the 14-th Prague Stringology Conference, Czech Technical University, 2009, pp. 237–248, accepted to International Journal of Foundations of Computer Science.

23. N. PYTHEAS FOGG: *Substitutions in Dynamics, Arithmetics and Combinatorics*, vol. 1794 of Lecture Notes in Mathematics, Springer, 2002.
24. W. RYTTER: *The structure of subword graphs and suffix trees of Fibonacci words*. Theoretical Computer Science, 363(2) 2006, pp. 211–223.
25. 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.
26. J. SHALLIT: *Characteristic words as fixed points of homomorphisms*, Tech. Rep. CS-91-72, University of Waterloo, Department of Computer Science, 1991.
27. J. SIMPSON: *Modified Padovan words and the maximum number of runs in a word*. Australian Journal of Combinatorics, 46 2010, pp. 129–145.
28. H. USCKA-WEHLOU: *Digital lines, Sturmian words, and continued fractions*, PhD thesis, Department of Mathematics, Uppsala University, 2009.