

Fewest repetitions in infinite binary words

Golnaz Badkobeh, Maxime Crochemore

► **To cite this version:**

Golnaz Badkobeh, Maxime Crochemore. Fewest repetitions in infinite binary words. *RAIRO - Theoretical Informatics and Applications (RAIRO: ITA)*, EDP Sciences, 2012, 46 (1), pp.17-31. 10.1051/ita/2011109 . hal-00742086

HAL Id: hal-00742086

<https://hal-upec-upem.archives-ouvertes.fr/hal-00742086>

Submitted on 13 Feb 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Fewest repetitions in infinite binary words

Golnaz Badkobeh

King's College London, London, UK

Maxime Crochemore

King's College London, London, UK
and Université Paris-Est, France

July 25, 2012

Abstract

A square is the concatenation of a nonempty word with itself. A word has period p if its letters at distance p match. The exponent of a nonempty word is the quotient of its length over its smallest period.

In this article we give a proof of the fact that there exists an infinite binary word which contains finitely many squares and simultaneously avoids words of exponent larger than $7/3$.

Our infinite word contains 12 squares, which is the smallest possible number of squares to get the property, and 2 factors of exponent $7/3$. These are the only factors of exponent larger than 2.

The value $7/3$ introduces what we call the finite-repetition threshold of the binary alphabet. We conjecture it is $7/4$ for the ternary alphabet, like its repetitive threshold.

Keywords: combinatorics on words, repetitions, word morphisms.

MSC: 68R15 Combinatorics on words.

1 Introduction

Repetitions in words is a basic question in Theoretical Informatics, certainly because it is related to many applications although it has first been studied by Thue at the beginning of the twentieth century [11] with a pure theoretical objective. Related results apply to the design of efficient string pattern matching algorithm, to text compression methods and entropy analysis, as well as to the study of repetitions in biological molecular sequences among others.

The knowledge of the strongest constraints an infinite word can tolerate help for the design and analysis of efficient algorithms. The optimal bound on the maximal exponent of factors of the word has been studied by Thue and many other authors after him. One of the first discoveries was that an infinite binary word can avoid factors with an exponent larger than 2, called 2^+ -powers. This has been extended by Dejean [3] to the ternary alphabet and her famous

conjecture on the repetitive threshold for larger alphabets has eventually been proved recently after a series of partial results by different authors (see [9, 2] and references therein).

Another constraint is considered by Fraenkel and Simpson [4]: their parameter to the complexity of binary infinite words is the number of squares occurring in them without any restriction on the number of occurrences. It is fairly straightforward to check that no infinite binary word can contain less than three squares and they proved that some of them contain exactly three. Two of these squares appear in the cubes 000 and 111 so that the maximum exponent is 3 in their word. In this article we produce an infinite word with few distinct squares and a smaller maximal exponent.

Fraenkel and Simpson's proof uses a pair of morphisms, one to get an infinite word by iteration, the other to produce the final translation on the binary alphabet. Their result has been proved with different pairs of morphisms by Rampersad et al. [8] (the first morphism is uniform), by Harju and Nowotka [5] (the second morphism accepts any infinite square-free word), and by Badkobeh and Crochemore [1] (the simplest morphisms).

In this article we show that we can combine the two types of constraints for the binary alphabet: producing an infinite word whose maximal exponent of its factor is the smallest possible while containing the smallest number of squares. The maximal exponent is $7/3$ and the number of squares is 12 to which can be added two words of exponent $7/3$.

It is known from Karhumäki and Shallit [6] that if an infinite binary word avoids $7/3$ -powers it contains an infinite number of squares. Proving that it contains more than 12 squares is indeed a matter of simple computation.

Shallit [10] has built an infinite binary word avoiding $7/3^+$ -powers and all squares of period at least 7. His word contains 18 squares.

Our infinite binary word avoids the same powers but contains only 12 squares, the largest having period 8. As before the proof relies on a pair of morphisms satisfying suitable properties. Both morphisms are almost uniform (up to one unit). The first morphism is weakly square-free on a 6-letter alphabet, and the second does not even correspond to a uniquely-decipherable code but admits a unique decoding on the words produced by the first. To get the morphisms, we first examined carefully the structure of long words satisfying the conditions and obtained by backtracking computation. Then, we inferred the morphisms from the regularities found in the words.

After introducing the definitions and main results in the next section, we provide a weakly square-free morphism and the infinite square-free word on 6 letters it generates in Section 3. Section 4 shows how this word is translated into an infinite binary word satisfying the constraints. In the conclusion we define the new notion of finite-repetition threshold and state a conjecture on its value for the 3-letter alphabet.

2 Repetitions in binary words

A word is a sequence of letters drawn from a finite alphabet. We consider the binary alphabet $B = \{0, 1\}$, the ternary alphabet $A_3 = \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$, and the 6-letter alphabet $A_6 = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}, \mathbf{e}, \mathbf{f}\}$.

A square is a word of the form uu where u is a nonempty (finite) word. A word has period p if its letters at distance p are equal. The exponent of a nonempty word is the quotient of its length over its smallest period. Thus, a square is any word with an even integer exponent.

In this article we consider infinite binary words in which a small number of squares occur.

The maximal length of a binary word containing less than three square is finite. It can be checked that it is 18, e.g. 010011000111001101 contains only 00 and 11. But, as recalled above, this length is infinite if 3 squares are allowed to appear in the word. A simple proof of it relies on two morphisms f and h_0 defined as follows. The morphism f is defined from A_3^* to itself by

$$\begin{cases} f(\mathbf{a}) = \mathbf{abc}, \\ f(\mathbf{b}) = \mathbf{ac}, \\ f(\mathbf{c}) = \mathbf{b}. \end{cases}$$

It is known that the infinite word $\mathbf{f} = f(\mathbf{a})^\infty$ it generates is square-free (see [7, Chapter 2]). The morphism h_0 is from A_3^* to B^* and defined by

$$\begin{cases} h(\mathbf{a}) = 01001110001101, \\ h(\mathbf{b}) = 0011, \\ h(\mathbf{c}) = 000111. \end{cases}$$

Then the result is a consequence of the next statement.

Theorem 1 ([1]) *The infinite word $\mathbf{h}_0 = h_0(f(\mathbf{a})^\infty)$ contains the 3 squares 00, 11 and 1010 only. The cubes 000 and 111 are the only factors occurring in \mathbf{h} and of exponent larger than 2.*

It is impossible to avoid 2^+ -powers and keep a bounded number of squares. As proved by Karhumäki and Shallit [6], the exponent has to go up to $7/3$ to allow the property.

In the two following sections we define two morphisms and derive the properties that we need to prove the next statement.

Theorem 2 *There exists an infinite binary word whose factors have an exponent at most $7/3$ and that contains 12 squares, the fewest possible.*

Our infinite binary word contain the 12 squares 0^2 , 1^2 , $(01)^2$, $(10)^2$, $(001)^2$, $(010)^2$, $(011)^2$, $(100)^2$, $(101)^2$, $(110)^2$, $(01101001)^2$, $(10010110)^2$, and the two words 0110110 and 1001001 of exponent $7/3$.

Proving that it is impossible to have less than 12 squares in the previous statement results from the next table. It has been obtained by pruned backtracking sequential computation that avoids exhaustive search. It shows the

maximal length of binary words whose factors have an exponent at most $7/3$, for each number s of squares, $0 \leq s \leq 11$.

| | | | | | | | | | | | | |
|-----------|-----|---|---|----|----|----|----|----|----|----|----|-----|
| s | = 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $\ell(s)$ | = 3 | 5 | 8 | 12 | 14 | 18 | 24 | 30 | 37 | 43 | 83 | 116 |

3 A weakly square-free morphism on six letters

In this section we consider a specific morphism used for the proof of Theorem 2. It is called g and defined from A_6^* to itself by:

$$\left\{ \begin{array}{l} g(\mathbf{a}) = \mathbf{abac}, \\ g(\mathbf{b}) = \mathbf{babd}, \\ g(\mathbf{c}) = \mathbf{eabdf}, \\ g(\mathbf{d}) = \mathbf{fbace}, \\ g(\mathbf{e}) = \mathbf{bace}, \\ g(\mathbf{f}) = \mathbf{abdf}. \end{array} \right.$$

We prove below that the morphism is weakly square-free in the sense that $\mathbf{g} = g^\infty(\mathbf{a})$ is an infinite square-free word, that is, all its finite factors have an exponent smaller than 2. Note that however it is not square-free since for example $g(\mathbf{cf}) = \mathbf{eabdfabdf}$ contains the square $(\mathbf{abdf})^2$. This prevents from using characterisation of square-freeness of the morphism, or equivalently of the fixed points of the morphism. As far as we know only an ad hoc proof is possible.

The set of codewords $g(a)$'s ($a \in A_6$) is a prefix code and therefore a uniquely-decipherable code. Note also that any occurrence of \mathbf{abac} in $g(w)$, for $w \in A_6^*$, uniquely corresponds to an occurrence of \mathbf{a} in w . The proof below relies on the fact that not all doublets and triplets (words of length 2 and 3 respectively) occur in \mathbf{g} , as the next statements show.

Lemma 1 *The set of doublets occurring in \mathbf{g} is*

$$D = \{\mathbf{ab}, \mathbf{ac}, \mathbf{ba}, \mathbf{bd}, \mathbf{cb}, \mathbf{ce}, \mathbf{da}, \mathbf{df}, \mathbf{ea}, \mathbf{fb}\}.$$

Proof. Note that all letters of A_6 appear in \mathbf{g} . Then doublets $\mathbf{ab}, \mathbf{ac}, \mathbf{ba}, \mathbf{bd}, \mathbf{ce}, \mathbf{df}, \mathbf{ea}, \mathbf{fb}$ appear in \mathbf{g} because they appear in the images of one letter. The images of these doublets generate two more doublets, \mathbf{cb} and \mathbf{da} , whose images do not create new doublets. ■

Lemma 2

The set of triplets in \mathbf{g} is

$$T = \{\mathbf{aba}, \mathbf{abd}, \mathbf{acb}, \mathbf{ace}, \mathbf{bab}, \mathbf{bac}, \mathbf{bda}, \mathbf{bdf}, \mathbf{cba}, \mathbf{cea}, \mathbf{dab}, \mathbf{dfb}, \mathbf{eab}, \mathbf{fba}\}.$$

Proof. Triplets appear in the images of a letter or of a doublet. Triplets found in images of one letter are: $\mathbf{aba}, \mathbf{abd}, \mathbf{ace}, \mathbf{bab}, \mathbf{bac}, \mathbf{bdf}, \mathbf{eab}, \mathbf{fba}$. The images of doublets occurring in \mathbf{g} , in set D of Lemma 1, contain the extra triplets: $\mathbf{acb}, \mathbf{bda}, \mathbf{cba}, \mathbf{cea}, \mathbf{dab}, \mathbf{dfb}$. ■

Table 1: Gaps of **abac**: words between consecutive occurrences of **abac** in **g**. They are images of gaps between consecutive occurrences of **a**.

| | | | |
|--------------------|---|--------------------------|----|
| $g(\mathbf{b})$ | = | babd | 4 |
| $g(\mathbf{cb})$ | = | eabdfbabd | 9 |
| $g(\mathbf{bd})$ | = | babdfbace | 9 |
| $g(\mathbf{ce})$ | = | eabdfbace | 9 |
| $g(\mathbf{bdfb})$ | = | babdfbaceabdfbabd | 17 |

To prove that the infinite word **g** is square-free we first show that it contains no square with less than four occurrences of the word $g(\mathbf{a}) = \mathbf{abac}$. Then, we show it contains no square with at least four occurrences of it. The word **abac** is chosen because its occurrences in **g** correspond to $g(\mathbf{a})$ only, so they are used to synchronise the parsing of the word according to the codewords $g(\mathbf{a})$'s.

Lemma 3 *No square in **g** can contain less than four occurrences of **abac**.*

Proof. Assume by contradiction that a square ww in **g** contains less than four occurrences of **abac**. Let x be the shortest word whose image by g contains ww .

Then x is a factor of **g** that belongs to the set $\mathbf{a}((A_6 \setminus \{\mathbf{a}\})^* \mathbf{a})^5$. Since two consecutive occurrences of **a** in **g** are separated by a string of length at most 4 (the largest such string is indeed **bdfb** as a consequence of Lemma 2), the set is finite.

The square-freeness of all these factors has been checked via an elementary implementation of the test, which proves the result. ■

Proposition 1 *No square in **g** can contain at least four occurrences of **abac**.*

Proof. The proof is by contradiction: let k be the maximal integer for which $g^k(\mathbf{a})$ is square-free and let ww be a square occurring in $g^{k+1}(\mathbf{a})$ and containing at least 4 occurrences of **abac**. Distinguishing several cases according to the words between consecutive occurrences of **abac** (see Table 1), we deduce that $g^k(\mathbf{a})$ is not square-free, the contradiction.

The square ww can be written

$$\underbrace{v_0(\mathbf{abac} \cdots \mathbf{abac})u_1}_{v_0(\mathbf{abac} \cdots \mathbf{abac})u_1} \underbrace{v_1(\mathbf{abac} \cdots \mathbf{abac})u_2}_{v_1(\mathbf{abac} \cdots \mathbf{abac})u_2}$$

where v_0, u_1, v_1 , and u_2 contain no occurrence of **abac**. It occurs in the image of a factor of **g**. The central part of w starting and ending with **abac** is the image of a unique word U factor of $g^k(\mathbf{a})$ due to the code property:

$$g(U) = v_0^{-1} w u_1^{-1} = v_1^{-1} w u_2^{-1}.$$

We split the proof in two parts according to whether **abac** occurs in $u_1 v_1$ or not.

No abac in u_1v_1 . We consider five cases according to the value of u_1v_1 , the gap of **abac** (see Table 1).

1. $u_1v_1 = \mathbf{babd}$ corresponds to $g(\mathbf{b})$ only. If either u_1 or v_1 is empty, then v_0 or u_2 is $g(\mathbf{b})$, in either case we get \mathbf{bUbU} or \mathbf{UbUb} that are squares. Else v_0 has a suffix \mathbf{d} so it belongs to $g(\mathbf{b})$, and again \mathbf{bUbU} is a square in \mathbf{g} .
2. $u_1v_1 = \mathbf{eabdfbabd}$ corresponds to $g(\mathbf{cb})$ only. An occurrence of \mathbf{cb} always belongs to $g(\mathbf{ab})$ therefore U has a prefix \mathbf{abd} and a suffix \mathbf{aba} , and the letter after \mathbf{aba} is \mathbf{c} . If v_1 is empty, u_2 has a prefix $\mathbf{eabdfbabd}$ so it is $g(\mathbf{cb})$ and again \mathbf{UcbUcb} is a square. If v_1 is not empty then v_0 has a suffix \mathbf{d} , suffix of $g(\mathbf{b})$, therefore \mathbf{bUcbUc} is a square.
3. $u_1v_1 = \mathbf{babdfbace}$ corresponds to $g(\mathbf{bd})$. The word \mathbf{abda} is a factor of $g(\mathbf{ba})$ only so U has a prefix \mathbf{aba} and a suffix \mathbf{ba} . If $|u_1| = 0$, $v_0 = \mathbf{babdfbace}$ can only be $g(\mathbf{bd})$ so \mathbf{bdUbdU} is a square. Otherwise u_2 must have a prefix \mathbf{b} ; since U has a suffix \mathbf{ba} the next letter after it is either \mathbf{b} or \mathbf{c} ; as only $g(\mathbf{b})$ is prefixed by \mathbf{b} the letter is \mathbf{b} so u_2 has a prefix or is a prefix of $g(\mathbf{b})$, and we know that \mathbf{bab} is always followed by \mathbf{d} thus \mathbf{UbdUbd} is a square.
4. $u_1v_1 = \mathbf{eabdfbace}$ corresponds to $g(\mathbf{ce})$ only. If u_1 is empty, v_0 is $g(\mathbf{ce})$ so \mathbf{ceUceU} is a square. Otherwise, u_2 has a prefix or is a prefix of $g(\mathbf{c})$; the next letter after $g(\mathbf{c})$ is either \mathbf{b} or \mathbf{e} ; (see Lemma 1); if it is \mathbf{b} the right-most U has a suffix \mathbf{aba} but the left-most U has a suffix \mathbf{fba} , which cannot be. Therefore the letter after \mathbf{c} is \mathbf{e} and \mathbf{UceUce} is a square.
5. $u_1v_1 = \mathbf{babdfbaceabdfbabd}$. If $|v_1| > 12$, v_0 has a suffix $g(\mathbf{dfb})$ and the letter before it is \mathbf{b} , so $\mathbf{bdfbUbdfbU}$ is a square. If $0 < |v_1| \leq 12$, then $|u_1| \geq 5$, so u_2 has a prefix or is a prefix of $g(\mathbf{bd})$ so the next letter is either \mathbf{a} or \mathbf{f} . If it is \mathbf{a} the right-most U has a suffix \mathbf{ba} but v_0 is a suffix of or has a suffix $g(\mathbf{b})$; the letter before it is either $g(\mathbf{c})$ or $g(\mathbf{f})$; if it is \mathbf{c} then U has a prefix \mathbf{abd} and $\mathbf{bdfbabd}$ is from the concatenation of $g(\mathbf{c})$ and $g(\mathbf{b})$ or $g(\mathbf{dfb})$; in either case the left occurrence of U will have \mathbf{ea} as a suffix, a contradiction since $\mathbf{fbUbdfbUbd}$ and $\mathbf{UbdfbUbdfb}$ are both squares.

An occurrence of abac in u_1v_1 . Then the suffix of u_1 is either \mathbf{aba} , \mathbf{ab} or \mathbf{a} while the respective prefix of v_1 is \mathbf{c} , \mathbf{ac} or \mathbf{bac} .

Note that \mathbf{c} is followed either by \mathbf{b} or \mathbf{e} (Lemma 1) and that \mathbf{cb} occurs only in the image of \mathbf{ab} . Then if the occurrence of **abac** is followed by \mathbf{b} , the occurrence of \mathbf{cb} in v_0 is preceded by \mathbf{aba} , and then there is a square starting 1, 2 or 3 positions before the occurrence of ww , which brings us back to the first case. Therefore, **abac** is followed by \mathbf{e} .

The occurrence of **abace** comes from $g(\mathbf{ac})$, and by Lemma 2 u_1v_1 contains an occurrence of $g(\mathbf{bac})$. So, the occurrence of **abace** is preceded by \mathbf{d} , and since \mathbf{da} occurs only in the image of \mathbf{ba} , the occurrence of \mathbf{da} in u_2 is followed by \mathbf{bac} , which yields a square starting 1, 2 or 3 positions after the occurrence of ww . Again this takes us back to the first case.

In all cases we deduce the existence of a square in $g^k(\mathbf{a})$, which is a contradiction with the definition of k . Therefore there is no square in \mathbf{g} containing at least four occurrences of \mathbf{abac} . ■

The next corollary is a direct consequence of Lemma 3 and Proposition 1.

Corollary 1 *The infinite word \mathbf{g} is square-free, or equivalently, the morphism g is weakly square-free.*

4 Binary translation

The second part of the proof of Theorem 2 consists in showing that the special infinite square-free word on 6 letters introduced in the previous section can be transformed into the desired binary word. This is done with a second morphism h from A_6^* to B^* defined by

$$\begin{cases} h(\mathbf{a}) = 10011, \\ h(\mathbf{b}) = 01100, \\ h(\mathbf{c}) = 01001, \\ h(\mathbf{d}) = 10110, \\ h(\mathbf{e}) = 0110, \\ h(\mathbf{f}) = 1001. \end{cases}$$

Note that the codewords of h do not form a prefix code, nor a suffix code, nor even a uniquely-decipherable code! We have for example $g(\mathbf{ae}) = 10011 \cdot 0110 = 1001 \cdot 10110 = g(\mathbf{fd})$. However, parsing the word $h(y)$ when y is a factor of \mathbf{g} is unique due to the absence of some doublets and triplets in it (see Lemmas 1 and 2). For example \mathbf{fd} does not occur, which induces the unique parsing of 100110110 as $10011 \cdot 0110$.

Proposition 2 *The infinite word $\mathbf{h} = h(g^\infty(\mathbf{a}))$ contains no factor of exponent larger than $7/3$. It contains the 12 squares 0^2 , 1^2 , $(01)^2$, $(10)^2$, $(001)^2$, $(010)^2$, $(011)^2$, $(100)^2$, $(101)^2$, $(110)^2$, $(01101001)^2$, $(10010110)^2$ only. Words 0110110 and 1001001 are the only factors with an exponent larger than 2.*

The proof is based on the fact that occurrences of 10011 in \mathbf{h} identify occurrences of \mathbf{a} in \mathbf{g} and on the unique parsing mentioned above. It proceeds by considering several cases according to the gaps between consecutive occurrences of 10011 (see Table 2), associated with gaps between consecutive occurrences of \mathbf{a} in \mathbf{g} , which leads to analyse paths in the graph of Figure 1.

Proof. We show that if \mathbf{h} would contain a square not in the list it would come from a square in \mathbf{g} , which cannot be since \mathbf{g} is square-free (Corollary 1).

Suppose \mathbf{h} contains the square w^2 . It is a factor of $h(g^k(\mathbf{a}))$, for some integer k and can be written $\underbrace{v_0(h(\mathbf{a}) \cdots h(\mathbf{a}))u_1}_{\text{}} \underbrace{v_1(h(\mathbf{a}) \cdots h(\mathbf{a}))u_2}_{\text{}}$. The central part of

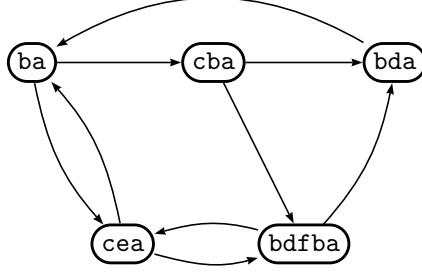


Figure 1: Graph showing immediate successors of gaps in the word \mathbf{g} : a suffix of it following an occurrence of \mathbf{a} is the label of an infinite path.

Table 2: Gaps between consecutive occurrences of 10011 in \mathbf{h} .

| | | | |
|--------------------|---|---------------------|----|
| $h(\mathbf{b})$ | = | 01100 | 5 |
| $h(\mathbf{cb})$ | = | 0100101100 | 10 |
| $h(\mathbf{bd})$ | = | 0110010110 | 10 |
| $h(\mathbf{ce})$ | = | 010010110 | 9 |
| $h(\mathbf{bdfb})$ | = | 0110010110100101100 | 19 |

w is the image of a unique square-free factor U of $g^k(\mathbf{a})$ due to the unique parsing mentioned above:

$$h(U) = (h(\mathbf{a}) \cdots h(\mathbf{a})) = v_0^{-1} w u_1^{-1} = v_1^{-1} w u_2^{-1}.$$

We proceed through different cases as in the proof of Proposition 1.

No $h(\mathbf{a})$ in $u_1 v_1$.

1. $u_1 v_1 = 01100$ corresponds to $h(\mathbf{b})$ only.

If $|v_1| > 1$, then v_0 belongs to $h(\mathbf{b})$, $\mathbf{b}U\mathbf{b}U$ is a square. Else $|u_1| \geq 4$ so u_2 belongs to $h(\mathbf{b})$, it cannot belong to $h(\mathbf{e})$ since \mathbf{ae} is not a factor of \mathbf{g} , therefore $U\mathbf{b}U\mathbf{b}$ is a square of \mathbf{g} .

2. $u_1 v_1 = 0110010110$ corresponds to $h(\mathbf{bd})$.

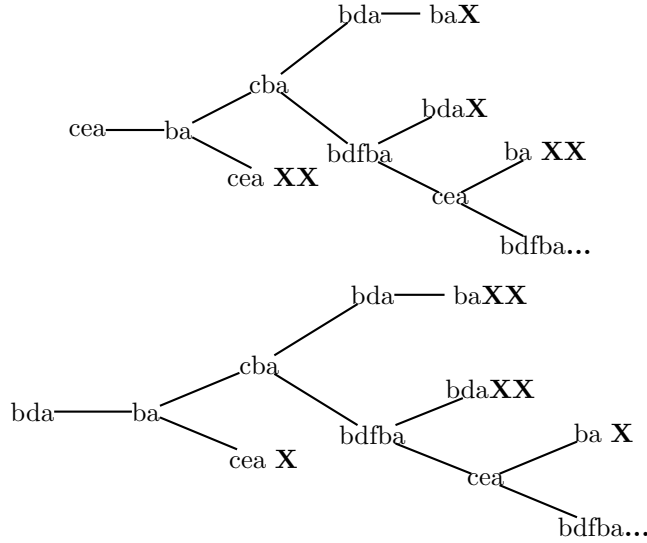
$$v_0 \underbrace{(h(\mathbf{a}) \cdots h(\mathbf{a}))}_{h(\mathbf{bd})} \underbrace{(h(\mathbf{a}) \cdots h(\mathbf{a}))}_{u_2}$$

the word \mathbf{abda} is a factor of $g(\mathbf{ba})$ only, so U has a prefix \mathbf{abac} and a suffix \mathbf{ba} (Note that U cannot be \mathbf{aba} since $\mathbf{ababdaba}$ is not a factor of \mathbf{g}).

$$v_0 \underbrace{(h(\mathbf{abac}) \cdots h(\mathbf{ba}))}_{h(\mathbf{bd})} \underbrace{(h(\mathbf{abac}) \cdots h(\mathbf{ba}))}_{u_2}$$

If u_2 comes from or has a prefix $h(\mathbf{b})$ then the letter after \mathbf{bab} is always \mathbf{d} so we have the square $U\mathbf{bd}U\mathbf{bd}$. Then u_2 is a prefix of or has a prefix $h(\mathbf{c})$, the longest common prefix (LCP) of $h(\mathbf{c})$ and $h(\mathbf{b})$ is 01 , so v_0 has a suffix 10010110 , which is a suffix of $h(\mathbf{bd})$ or $h(\mathbf{ce})$. If v_0 comes from $h(\mathbf{bd})$ then we have the square $\mathbf{bd}U\mathbf{bd}U$. So v_0 is a suffix of $h(\mathbf{ce})$

$$h(\mathbf{ce}) \underbrace{(h(\mathbf{abac}) \cdots h(\mathbf{ba}))}_{h(\mathbf{bd})} \underbrace{(h(\mathbf{abac}) \cdots h(\mathbf{ba}))}_{h(\mathbf{c})}.$$



The sign XX shows that the particular branch of the trie terminates because either a square occurs or the sequence is not a factor of \mathbf{g} . The sign X on the other hand represents the termination of a particular branch as a consequence of the discontinuation of the corresponding branch in the other trie. If we continue these tries we will have:

$$\begin{array}{l} \mathbf{ce} \underbrace{\mathbf{abac\ babd\ fbace\ abdf\ babd\ abac\ eabdf\ bace\ abac\ babd\ abac\ eabdf \dots ba}} \\ \mathbf{bd} \underbrace{\mathbf{abac\ babd\ fbace\ abdf\ babd\ abac\ eabdf\ bace\ abac\ babd\ abac\ eabdf \dots ba} \mathbf{c} \end{array}$$

which is the image of

$$\mathbf{eabdf\ bace} \underbrace{\mathbf{abac \dots abac}}_{h(\mathbf{bd})} \mathbf{babd\ fbace} \underbrace{\mathbf{abac \dots abac}}_{h(\mathbf{c})} \mathbf{e}$$

itself image of

$$\mathbf{ce} \underbrace{\mathbf{a \dots a}}_{h(\mathbf{bd})} \mathbf{bd} \underbrace{\mathbf{a \dots a}}_{h(\mathbf{c})} \mathbf{c}$$

so we have the same situation as at the starting point; but U is shorter in this case, therefore if we continue this process we should have

$$\mathbf{ce\ abac\ babd\ fbace\ abdf\ babd\ abac\ babd\ fbace\ abdf\ bace\ a}$$

but $\mathbf{abdf\ bace}$ is the image of \mathbf{fe} that is not in D (Lemma 1).

3. $u_1v_1 = 0100101100$ corresponds to $h(\mathbf{cb})$.

The word \mathbf{acba} is a factor of $g(\mathbf{ab})$ only, so U has a prefix \mathbf{abd} and a suffix \mathbf{aba} :

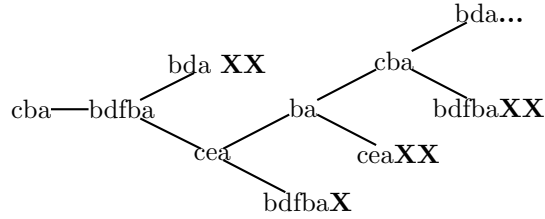
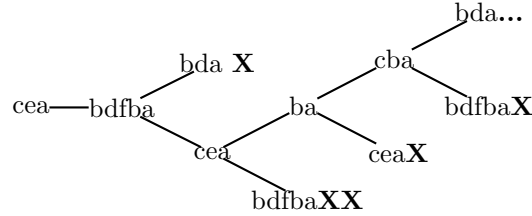
$$v_0 \underbrace{(h(\mathbf{abd}) \dots h(\mathbf{aba}))} h(\mathbf{cb}) \underbrace{(h(\mathbf{abd}) \dots h(\mathbf{aba}))} u_2$$

The word u_2 comes from or has a prefix $h(\mathbf{c})$. If the letter after it is \mathbf{b} , we have the square $U\mathbf{cb}U\mathbf{cb}$.

Otherwise u_2 comes from or has a prefix $h(\mathbf{ce})$. If v_0 comes from or has a suffix $h(\mathbf{b})$ then we have the square $\mathbf{b}U\mathbf{cb}U\mathbf{c}$.

Therefore the letter before U is \mathbf{e} preceded by \mathbf{c} , i.e. the string before the left U is \mathbf{ce} :

$$h(\mathbf{ce}) \underbrace{(h(\mathbf{abd}) \dots h(\mathbf{aba}))} h(\mathbf{cb}) \underbrace{(h(\mathbf{abd}) \dots h(\mathbf{aba}))} h(\mathbf{ce}).$$



Now we have the same situation as in the previous case

$$h(g(\mathbf{ce})) \underbrace{(h(g(\mathbf{abac})) \dots h(g(\mathbf{ba})))} h(g(\mathbf{bd})) \underbrace{(h(g(\mathbf{abac})) \dots h(g(\mathbf{ba})))} h(g(\mathbf{c})).$$

4. $u_1v_1 = 010010110$ corresponds to $h(\mathbf{ce})$ only.

Before \mathbf{c} is always \mathbf{ba} (Lemma 2) and after \mathbf{e} is \mathbf{ab} (Lemma 2), so \mathbf{ab} is a prefix of U and \mathbf{ba} is a suffix of U :

$$v_0 \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{ba}))} h(\mathbf{ce}) \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{ba}))} u_2.$$

(i): u_2 belongs to $h(\mathbf{cb})$ since we cannot have $U\mathbf{ce}U\mathbf{ce}$ and the letter after \mathbf{c} is \mathbf{b} or \mathbf{e} (Lemma 1):

$$v_0 \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{ba}))} h(\mathbf{ce}) \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{ba}))} h(\mathbf{cb})$$

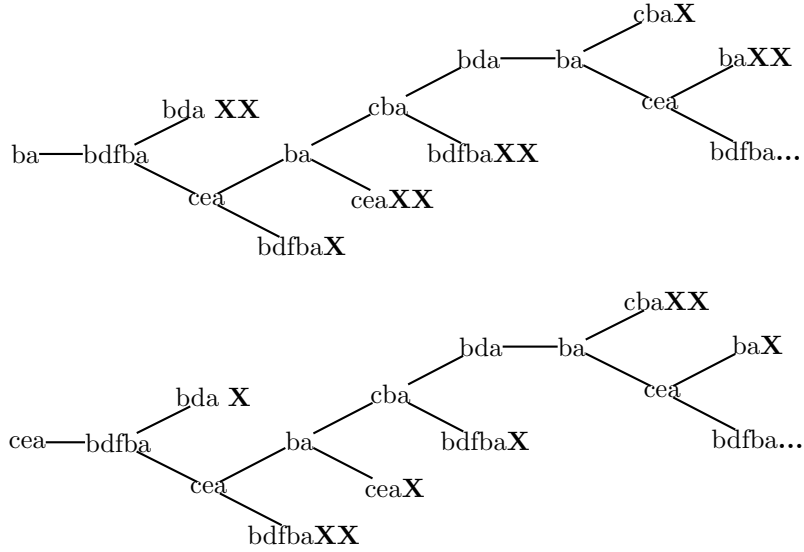
The letter before **bacb** is **a** so:

$$v_0 \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{aba}))}_{\mathbf{a}} h(\mathbf{ce}) \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{aba}))}_{\mathbf{a}} h(\mathbf{cb}).$$

NOTE: U is not **aba** since **abaceabacb** is not a factor of $g^k(\mathbf{a})$.

Now **abace** is a prefix of the image of **ac** so U has a prefix **abdf** and the word before it is either **ce** or **b**; the first choice gives the square **ceUceU** and the second choice:

$$h(\mathbf{b}) \underbrace{(h(\mathbf{abdf}) \dots h(\mathbf{aba}))}_{\mathbf{b}} h(\mathbf{ce}) \underbrace{(h(\mathbf{abdf}) \dots h(\mathbf{aba}))}_{\mathbf{b}} h(\mathbf{cb}).$$



Now if we continue the above tries we get:

$$\mathbf{b} \underbrace{\mathbf{abdfbaceabacbabdabaceabdfbabdabacbabdfbaceabdfba \dots ba}}_{\mathbf{b}} \\ \mathbf{ce} \underbrace{\mathbf{abdfbaceabacbabdabaceabdfbabdabacbabdfbaceabdfba \dots ba}}_{\mathbf{ce}} \mathbf{cb}$$

which is the image of

$$\mathbf{bd} \underbrace{\mathbf{abacbabdfbaceabdf \dots ba}}_{\mathbf{bd}} \mathbf{ce} \underbrace{\mathbf{abacbabdfbaceabdf \dots ba}}_{\mathbf{ce}} \mathbf{b}.$$

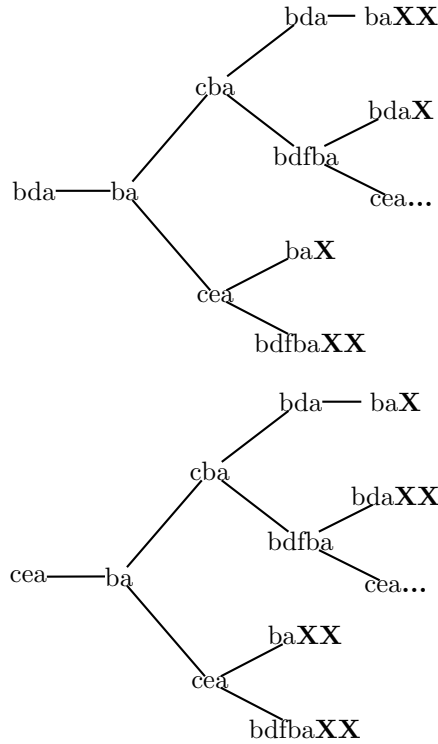
This is the same situation as the next case and we will see that after going one step back it brings us back to this case again. Now we are exactly in the same situation as at the beginning except that the length of the word $X = \mathbf{abdf \dots a}$ is shorter than U . Repeating this process enough times we should see that the word

$$\mathbf{babdfbaceabacbabdabaceabdfbaceabacbabdaba}$$

which is the image of $\mathbf{bdabaceaba}$, is not a factor of $g^k(\mathbf{a})$.

(ii): u_2 belongs to $h(\mathbf{b})$ (the LCP of $h(\mathbf{c})$ and $h(\mathbf{b})$ is 01) so v_0 must have a suffix 0010110 , which belongs to $h(\mathbf{bd})$ because if it belongs to $h(\mathbf{ce})$ then \mathbf{ceUceU} is a square.

$$h(\mathbf{bd}) \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{ba}))}_{\text{...}} h(\mathbf{ce}) \underbrace{(h(\mathbf{ab}) \dots h(\mathbf{ba}))}_{\text{...}} h(\mathbf{b}).$$



Continuing this trie we have

$$\mathbf{bd} \underbrace{\mathbf{abac babd fbace a \dots ba ce}}_{\text{...}} \underbrace{\mathbf{abac babd fbace a \dots ba bd}}_{\text{...}}.$$

This is factor of $g(\mathbf{b} \underbrace{\mathbf{abdf} \dots \mathbf{a ce}}_{\text{...}} \underbrace{\mathbf{abdf} \dots \mathbf{a cb}}_{\text{...}})$ which is the previous case.

5. $u_1v_1 = 0110010110100101100$ corresponds to $h(\mathbf{bdfb})$ only. This case is dealt with the same method.

$$u_0 \underbrace{(h(\mathbf{a}) \dots h(\mathbf{a}))}_{\text{...}} h(\mathbf{bdfb}) \underbrace{(h(\mathbf{a}) \dots h(\mathbf{a}))}_{\text{...}} u_2.$$

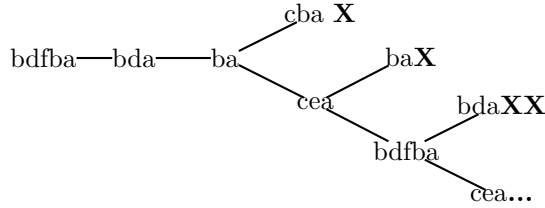
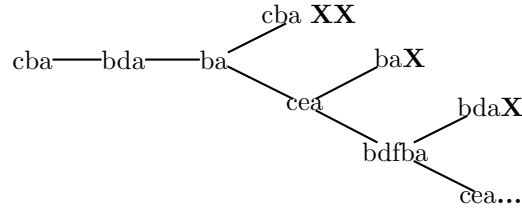
If u_2 belongs to $h(\mathbf{c})$, the LCP of $h(\mathbf{c})$ and $h(\mathbf{b})$ is 01 so u_0 must have a suffix 10010110100101100 , therefore u_0 belongs to $h(\mathbf{bdfb})$. But $\mathbf{bdfbU bdfbU}$

is a square and a factor of $g^k(a)$; a contradiction, so u_2 belongs to or has a prefix $h(b)$. We have two choices here.

(i): the next word after the right occurrence of U is ba . The LCP of $h(bd)$ and $h(ba)$ is 10 , u_0 has suffix of 110100101100 , so it either belongs to $h(dfb)$ or $h(acb)$. The first case gives that $dbfUdbfU$ is a square and a factor of $g^k(a)$, a contradiction. So u_0 belongs to $h(acb)$:

$$h(acb) \underbrace{(h(abda) \dots h(a))}_{h(bdfb)} \underbrace{(h(abda) \dots h(a))}_{h(ba)}.$$

Prefixes and suffixes of U are determined only by looking at D and T .



We have:

$$\begin{aligned} & abac\ babd\ \underbrace{abac\ eabdf\ bace \dots abac\ babd\ fbace}_{abdf\ babd\ abac\ eabdf\ bace \dots abac\ babd\ fbace\ abac} \\ & abdf\ babd\ \underbrace{abac\ eabdf\ bace \dots abac\ babd\ fbace}_{abac\ eabdf\ bace \dots abac\ babd\ fbace\ abac} \end{aligned}$$

which is the image of

$$\underbrace{ab\ ace \dots a\ bdfb\ ace \dots a\ bda}.$$

Now this is the next case so if we go back enough steps we should see that the length of U decreases and at the end we get

$$ac\ babd\ abac\ eabdf\ babd\ abac\ eaba$$

but this is not a factor of $g^k(a)$, a contradiction.

(ii): the word after U is bd . Now here the only possible letter after abd is a since if it is f it is a prefix of fb so we have $UdbfUdbf$, a contradiction. As the LCP of $h(bdfb)$ and $h(bda)$ is 01100101101001 u_0 must have a suffix 01100 so it can belong to $h(ab)$ or $h(acb)$.

(I):

$$h(\mathbf{ab}) \underbrace{(h(\mathbf{a}) \dots h(\mathbf{a}))}_{h(\mathbf{bdfb})} \underbrace{(h(\mathbf{a}) \dots h(\mathbf{a}))}_{h(\mathbf{bda})}.$$

Only using D , T and the Figure 1 we can continue building U ,

$$h(\mathbf{ab}) \underbrace{(h(\mathbf{ace}) \dots h(\mathbf{ba}))}_{h(\mathbf{bdfb})} \underbrace{(h(\mathbf{acea}) \dots h(\mathbf{ba}))}_{h(\mathbf{bda})}.$$

Continuing further we get:

$$h(\mathbf{abac eabdf babd} \underbrace{\mathbf{abac} \dots \mathbf{abac}}_{h(\mathbf{bdfb})} \mathbf{babd fbase abdf babd} \underbrace{\mathbf{abac} \dots \mathbf{abac}}_{h(\mathbf{bda})} \mathbf{babda}).$$

This is the image of

$$h(g(\mathbf{acb} \underbrace{\mathbf{a} \dots \mathbf{a}}_{h(\mathbf{bdfb})} \mathbf{a} \dots \mathbf{a} \mathbf{ba}))$$

and we are back to the case above.

(II):

$$h(\mathbf{acb}) \underbrace{(h(\mathbf{a}) \dots h(\mathbf{a}))}_{h(\mathbf{bdfb})} \underbrace{(h(\mathbf{a}) \dots h(\mathbf{a}))}_{h(\mathbf{bda})}.$$

Using the same method we build the word U :

$$\mathbf{ac b} \underbrace{\mathbf{abd} \dots \mathbf{ba}}_{h(\mathbf{bdfb})} \underbrace{\mathbf{ace} \dots \mathbf{ba}}_{h(\mathbf{bda})} \mathbf{bd a}.$$

Here we cannot go further as U cannot have \mathbf{abd} nor \mathbf{ace} as prefixes at the same time.

An occurrence of $h(\mathbf{a})$ in u_1v_1 . Looking at Figure 1, the image of the concatenation of two connected nodes (distance 1 arrow) are the possibilities for $u_1v_1h(\mathbf{a})$, but note that the second period of the square must start within $h(\mathbf{a})$, starting point of the arrow, otherwise it is one of the cases above. If the lengths of both nodes are larger than 2 then by unique parsing we are bound to have a square in $g^k(a)$ and get a contradiction. So we have to consider only the four cases where one of the nodes is \mathbf{ba} :

1. $u_1v_1 = h(\mathbf{bacb}) = 01100\mathbf{10011}0100101100$, so u_2 must have a prefix $h(\mathbf{b})$ and u_0 a suffix of $h(\mathbf{cb})$, before \mathbf{cb} is always \mathbf{a} , so $acbUbacbUb$ is a square in $g^k(a)$.
2. $u_1v_1 = h(\mathbf{bace}) = 01100\mathbf{10011}010010110$, so u_2 must have a prefix $h(\mathbf{b})$ and u_0 a suffix $h(\mathbf{ce})$, before \mathbf{ce} is always \mathbf{a} , so $aceUbaceUb$ is a square in $g^k(a)$.
3. $u_1v_1 = h(\mathbf{ceab}) = 010010110\mathbf{10011}01100$, so u_2 must have a prefix of $h(\mathbf{ce})$ and u_0 a suffix of $h(\mathbf{b})$, after \mathbf{ce} is always \mathbf{a} , so $bUceabUcea$ is a square in $g^k(a)$.

4. $u_1v_1 = h(\text{bdab}) = 01100101101\mathbf{001}101100$, so using tries as before shows that after enough backward iteration we should have

fbace abdf babd abac babd abac eabdf babd abac babd

which contains a square.

In all cases the conclusion is that we get a square in $g^k(\mathbf{a})$, a contradiction with the definition of k . This completes the proof of Proposition 2. ■

Theorem 2 follows immediately from Proposition 2.

5 Conclusion

The constraint on the number of squares imposed on binary words slightly differs from the constraint considered by Shallit [10]. The squares occurring in his word have period smaller than 7. Our word contains less squares but their maximal period is 8.

Looking at repetitions in words on larger alphabets, the subject introduces a new type of threshold, that we call the *finite-repetitions threshold* (FRt). For the alphabet of a letters, $\text{FRt}(a)$ is defined as the smallest rational number for which there exists an infinite word avoiding $\text{FRt}(a)^+$ -powers and containing a finite number of r -powers, where r is Dejean's repetitive threshold. Karhumäki and Shallit results as well as ours show that $\text{FRt}(2) = 7/3$. Our result additionally proves that the associated minimal number of squares is 12.

Computation shows that the maximal length of $(7/4)^+$ -free ternary word with only one $7/4$ -repetition is 102. This leads us state the following conjecture, which has been tested up to length 20000.

Conjecture 1 *The finite-repetitions threshold of 3-letter alphabet is $\frac{7}{4}$ and the associated number of $\frac{7}{4}$ -powers is 2.*

Values for larger alphabets remain to be explored.

References

- [1] G. Badkobeh and M. Crochemore. An infinite binary word containing only three distinct squares. 2010. Submitted.
- [2] J. D. Currie and N. Rampersad. A proof of Dejean's conjecture. *Math. Comput.*, 80(274):1063–1070, 2011.
- [3] F. Dejean. Sur un théorème de Thue. *J. Comb. Theory, Ser. A*, 13(1):90–99, 1972.
- [4] A. S. Fraenkel and J. Simpson. How many squares must a binary sequence contain? *Electr. J. Comb.*, 2, 1995.

- [5] T. Harju and D. Nowotka. Binary words with few squares. *Bulletin of the EATCS*, 89:164–166, 2006.
- [6] J. Karhumäki and J. Shallit. Polynomial versus exponential growth in repetition-free binary words. *J. Comb. Theory, Ser. A*, 105(2):335–347, 2004.
- [7] M. Lothaire, editor. *Combinatorics on Words*. Cambridge University Press, second edition, 1997.
- [8] N. Rampersad, J. Shallit, and M. Wei Wang. Avoiding large squares in infinite binary words. *Theor. Comput. Sci.*, 339(1):19–34, 2005.
- [9] M. Rao. Last cases of Dejean’s conjecture. *Theor. Comput. Sci.*, 412(27):3010–3018, 2011.
- [10] J. Shallit. Simultaneous avoidance of large squares and fractional powers in infinite binary words. *Intl. J. Found. Comput. Sci.*, 15(2):317–327, 2004.
- [11] A. Thue. Über unendliche Zeichenreihen. *Norske vid. Selsk. Skr. I. Mat. Nat. Kl. Christiania*, 7:1–22, 1906.