



ELSEVIER

Contents lists available at SciVerse ScienceDirect

Journal of Combinatorial Theory,
Series Awww.elsevier.com/locate/jctaThe Ehrenfeucht–Silberger problem [☆]Štěpán Holub ^{a,1}, Dirk Nowotka ^b^a Department of Algebra, Charles University of Prague, Sokolovska 83, 186 75 Praha 8, Czech Republic^b Department of Computer Science, Christian-Albrechts-Universität zu Kiel, Christian-Albrechts-Platz 4, 24118 Kiel, Germany

ARTICLE INFO

Article history:

Received 29 December 2009

Available online 30 November 2011

Keywords:

Combinatorics on words

Ehrenfeucht–Silberger problem

Periodicity

Unbordered words

ABSTRACT

We consider repetitions in words and solve a longstanding open problem about the relation between the period of a word and the length of its longest unbordered factor (where factor means uninterrupted subword). A word u is called bordered if there exists a proper prefix that is also a suffix of u , otherwise it is called unbordered. In 1979 Ehrenfeucht and Silberger raised the following problem: What is the maximum length of a word w , w.r.t. the length τ of its longest unbordered factor, such that τ is shorter than the period π of w . We show that, if w is of length $\frac{7}{3}\tau$ or more, then $\tau = \pi$ which gives the optimal asymptotic bound.

© 2011 Elsevier Inc. All rights reserved.

1. Introduction

When repetitions in words are considered then two notions are central: a (the) *period*, which gives an (the least) amount by which a word has to be shifted in order to overlap with itself, and a (the shortest) *border*, which denotes a (the least) nontrivial overlap of a word with itself. Both notions, period and border, are naturally related. For every $p < |w|$ we have that p is a period of w , if, and only if, w has a border of length $|w| - p$. In particular, the period of an unbordered word is its length, and the length of the shortest border of a bordered word is not larger than its period. Moreover, a shortest border itself is always unbordered. Periodicity also restricts occurrences of long unbordered factors (uninterrupted substrings). Deeper dependencies between the period of a word and its unbordered factors have been investigated [1–6] and exploited in applications [7–9] for decades; see also references to related work below.

Let us recall the problem by Ehrenfeucht and Silberger [2]. Let w be a (finite) word of length $|w|$, let $\tau(w)$ denote the maximum length of unbordered factors of w , and let $\pi(w)$ denote the

[☆] A preliminary version of this paper appeared in the proceedings of the ICALP 2009.E-mail addresses: holub@karlin.mff.cuni.cz (Š. Holub), dn@informatik.uni-kiel.de (D. Nowotka).¹ The work on this article has been supported by the research project MSM 0021620839.

period of w . Certainly, $\tau(w) \leq \pi(w)$ since a period of w is also a period of its factors. Moreover, it is folklore that $\tau(w) = \pi(w)$ when $|w| \geq 2\pi(w)$ (it follows, for example, from the existence of Lyndon conjugates; see Chapter 5.1 in [10]). So, the relation between $\tau(w)$ and $\pi(w)$ remains interesting in cases where $|w| < 2\pi(w)$. Actually, the interesting cases are also the most common ones since a simple counting argument shows that by far most words have a period that is longer than one half of their length. This leads to a much more difficult problem, raised by Ehrenfeucht and Silberger [2] (see also Schützenberger’s comments at the end of Chapter 8 in [10]), which asks about a bound on $|w|$ depending on $\tau(w)$ – rather than on $\pi(w)$ – such that $\tau(w) = \pi(w)$ is enforced. In this paper we establish the following fact for all finite words w :

$$\text{If } |w| \geq \frac{7}{3}\tau(w) \text{ then } \tau(w) = \pi(w).$$

This multiplicative bound on the length of w is asymptotically tight; see the following example by Assous and Pouzet [11]. We do not address the additive constant in this paper (see also Conclusions).

Previous work

Ehrenfeucht and Silberger raised the problem described above in [2]. They conjectured that $|w| \geq 2\tau(w)$ implies $\tau(w) = \pi(w)$. That conjecture was falsified shortly thereafter by Assous and Pouzet [11] by the following example:

$$w = a^n ba^{n+1} ba^n ba^{n+2} ba^n ba^{n+1} ba^n \tag{1}$$

where $n \geq 0$ and $\tau(w) = 3n + 6$ (note that $ba^{n+1}ba^nba^{n+2}$ and $a^{n+2}ba^nba^{n+1}b$ are the two longest unbordered factors of w) and $\pi(w) = 4n + 7$ and $|w| = 7n + 10$, that is, $\tau(w) < \pi(w)$ and $|w| = \frac{7}{3}\tau(w) - 4 > 2\tau(w)$. Assous and Pouzet [11] in turn conjectured that $3\tau(w)$ is the bound on the length of w for establishing $\tau(w) = \pi(w)$. Duval [4] did the next step towards answering the conjecture. He established that $|w| \geq 4\tau(w) - 6$ implies $\tau(w) = \pi(w)$ and conjectured that, if w possesses an unbordered prefix of length $\tau(w)$, then $|w| \geq 2\tau(w)$ implies $\tau(w) = \pi(w)$. Despite some partial results [12–14] towards a solution, Duval’s conjecture was only solved in 2004 [15,6] with a new proof given in [5]. It turned out that the optimal bound, for Duval’s conjecture, is $2\tau(w) - 1$; note that this result lowered the bound for Ehrenfeucht and Silberger’s problem to $3\tau(w) - 2$, in accordance with the conjecture by Assous and Pouzet [11].

However, there remained a gap of $\frac{2}{3}\tau(w)$ between that bound and the largest known example which is the one given above. The bound of $\frac{7}{3}\tau(w)$ has been conjectured in [15,6]. This conjecture is proved here, and the problem by Ehrenfeucht and Silberger is finally solved.

Other related work

The result related most closely to the problem by Ehrenfeucht and Silberger is the so called critical factorization theorem (CFT).

The CFT states the following: Let $w = uv$ be a factorization of a word w into u and v . The local period of w at the point $|u|$ is the length q of the shortest square centered at $|u|$ (see p. 4 for a more formal description). It is straightforward to see that q is not larger than the period of w . The factorization uv is called critical if q equals the period of w . The CFT states that a critical factorization exists for every nonempty word w , and moreover, a critical factorization uv can always be found such that $|u|$ is shorter than the period of w . The CFT was conjectured first by Schützenberger [16], proved by Césari and Vincent [1], and brought into its current form by Duval [3]. Crochemore and Perrin [7] found a new and elegant proof of the CFT using lexicographic orders, and realized a direct application of the theorem in a new string-matching algorithm.

How does the CFT relate to the problem by Ehrenfeucht and Silberger? Observe that the shortest square x^2 centered at some point in w is always such that x is unbordered. If x results from a critical factorization and x occurs in w , then $\tau(w) = \pi(w)$. Therefore, it immediately follows from the CFT that $|w| > 2\pi(w) - 2$ implies $\tau(w) = \pi(w)$. The multiplicative constant two is optimal as shown by the words $(aba)^k abba.(aba)^k$ of length $2\pi(w) - 4$ for which $\tau(w) = \pi(w) - 1$. As already mentioned, we establish the asymptotically optimal bound on $|w|$ enforcing the equality $\tau(w) = \pi(w)$ in terms

of $\tau(w)$ instead of $\pi(w)$. This rounds off the long lasting research effort on the mutual relationship between the two basic properties of a word w , that is, $\tau(w)$ and $\pi(w)$.

2. Notation and basic facts

Let us fix a finite set A of letters, called alphabet, for the rest of this paper. Let A^* denote the monoid of all finite words over A including the *empty word* denoted by ε . Let $u, v, w \in A^*$ such that $w = uv$. Then $u^{-1}w = v$ and $wv^{-1} = u$. For all $k \geq 0$, we define $w^0 = \varepsilon$ and $w^k = ww^{k-1}$, if $k > 0$. In general, we denote variables over A by a, b, c , and d and variables over A^* are usually denoted by f, g, h, r through z , and by Greek letters, including their subscripted and primed versions. Typically, Greek variables are used to indicate a word defined as a suffix with special lexicographic properties. The letters i through q are to range over the set of nonnegative integers.

Let $w = a_1a_2 \cdots a_n$. We denote the length n of w by $|w|$, in particular $|\varepsilon| = 0$. Let $1 \leq i \leq j \leq n$. Then $u = a_i a_{i+1} \cdots a_j$ is called a *factor* of w . Let $0 \leq i \leq n$. Then $u = a_1 a_2 \cdots a_i$ is called a *prefix* of w , denoted by $u \leq_p w$, and $v = a_{i+1} a_{i+2} \cdots a_n$ is called a *suffix* of w , denoted by $v \leq_s w$. The *longest common prefix* w of two words u and v is denoted by $u \wedge_p v$ and is defined so that if $u \leq_p v$, then $w = u$, and if $v \leq_p u$, then $w = v$, and in any other case w is such that $wa \leq_p u$ and $wb \leq_p v$ for some different letters a and b . The *longest common suffix* of u and v , denoted $u \wedge_s v$, is defined similarly, as one would expect. Two words u and v , with $|u| \leq |v|$, *overlap* each other, if there is a word w , with $|v| < |w| < |uv|$, such that $u \leq_p w$ and $v \leq_s w$ or $v \leq_p w$ and $u \leq_s w$. An integer $1 \leq p \leq n$ is a *period* of w if $a_i = a_{i+p}$ for all $1 \leq i \leq n - p$. The smallest period of w is called *the period* of w , denoted by $\pi(w)$. A nonempty word u is called a *border* of a word w , if $w = uy = zu$ for some nonempty words y and z . We call w *bordered*, if it has a border, otherwise w is called *unbordered*. Let $\tau(w)$ denote the maximum length of unbordered factors of w , and $\tau_2(w)$ denote the maximum length of unbordered factors occurring at least twice in w . Let $\tau_2(w) = 0$, if no unbordered factor occurs twice in w . We have that

$$\tau(w) \leq \pi(w). \quad (2)$$

Indeed, let $u = b_1 b_2 \cdots b_{\tau(w)}$ be an unbordered factor of w . If $\tau(w) > \pi(w)$ then $b_i = b_{i+\pi(w)}$ for all $1 \leq i \leq \tau(w) - \pi(w)$ and $b_1 b_2 \cdots b_{\tau(w) - \pi(w)}$ is a border of u ; a contradiction.

Let \triangleleft be a total order on A . Then \triangleleft extends to a *lexicographic order*, also denoted by \triangleleft , on A^* with $u \triangleleft v$ if either $u \leq_p v$ or $xa \leq_p u$ and $xb \leq_p v$ and $a \triangleleft b$. Let \triangleleft^a denote an order on A where a is the maximum letter. The \triangleleft -*maximum suffix* α of a word w is defined as the suffix of w such that $v \triangleleft \alpha$ for all $v \leq_s w$.

The following remarks state some facts about maximum suffixes which are folklore. They are included in this paper to make it self-contained.

Remark 2.1. Let w be a bordered word. The shortest border u of w is unbordered, and $w = uzu$. The longest border of w has length equal to $|w| - \pi(w)$.

Indeed, if u is a border of w , then each border of u is also a border of w . Therefore u is unbordered, and it does not overlap with itself. If v is a border of w then $|w| - |v|$ is a period of w . Conversely, the prefix of w of length $|w| - \pi(w)$ is a border of w .

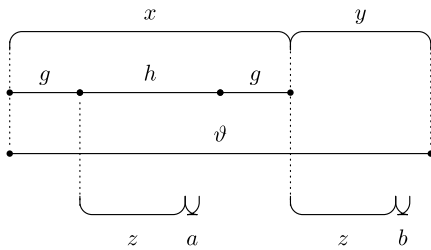
Remark 2.2. Any maximum suffix of a word w occurs only once in w and is longer than its longest border, that is, longer than $|w| - \pi(w)$.

Indeed, let α be the \triangleleft -maximum suffix of w for some order \triangleleft . Then $w = x\alpha y$ and $\alpha \triangleleft x\alpha y$ implies $y = \varepsilon$ by the maximality of α . If $w = uv\alpha$ with $|v| = \pi(w)$, then $u\alpha \leq_p w$ gives a contradiction again.

Remark 2.3. Let ϑ be its own maximum suffix w.r.t. some order \triangleleft , and let x be a prefix of ϑ of length $\pi(\vartheta)$. Then x is unbordered.

Indeed, suppose on the contrary that x is bordered, that is, $x = ghg$ for some nonempty g . Let $\vartheta = xy$. We have $gy \triangleleft \vartheta = ghgy$, by assumption, which implies $y \triangleleft hgy$. Note that gy is not a prefix of ϑ otherwise $|gh| < |x|$ is a period of ϑ contradicting the choice of x . Hence, $zb \leq_p y$ and $za \leq_p hgy$

for some different letters a and b with $b \triangleleft a$. But, $y \leq_p \vartheta$, since $|x| = \pi(w)$, implies $zb \leq_p \vartheta$ which contradicts the maximality of ϑ (since $zb \leq_p \vartheta \triangleleft za \leq_p hgy$). These arguments are illustrated by the following figure.



Let an integer q with $0 \leq q < |w|$ be called a *point* in w . A nonempty word x is called a *repetition word at point q* if $w = uv$ with $|u| = q$ and there exist words y and z such that $x \leq_s yu$ and $x \leq_p vz$. Let $\pi(w, q)$ denote the length of the shortest repetition word at point q in w . We call $\pi(w, q)$ the *local period at point q* in w . Note that the repetition word of length $\pi(w, q)$ at point q is necessarily unbordered and $\pi(w, q) \leq \pi(w)$. A factorization $w = uv$, with $u, v \neq \varepsilon$ and $|u| = q$, is called *critical*, if $\pi(w, q) = \pi(w)$, and if this holds, then q is called a *critical point*. Let \triangleleft be an order on A and \blacktriangleleft be its inverse. Then the shorter of the \triangleleft -maximum suffix and the \blacktriangleleft -maximum suffix of some word w is called a *critical suffix* of w . This terminology is justified by the following version of the so called critical factorization theorem (CFT) [7] which relates maximum suffixes and critical points.

Theorem 2.4 (CFT). *Let w be a nonempty word and γ be a critical suffix of w . Then $|w| - |\gamma|$ is a critical point.*

Remark 2.5. Let rs be an unbordered word where $|r|$ is a critical point. Then s and r do not overlap and sr is unbordered with $|s|$ as a critical point.

3. Special factorizations

Let us highlight the following definitions. They are not standard and will be central to the proof of Theorem 4.1. Let the words α and w be given. The use of Latin and Greek variables should suggest that the definitions will be typically applied in situations when w is a long word, and α is its short factor with some special lexicographic properties.

Definition 3.1. The longest prefix of α strictly shorter than α that is also a suffix of w will be called the α -*suffix* of w .

We want to note that the previous definition will be useful in situation when α is shorter than w , although it allows the other possibility too. Note also that the α -suffix is allowed to be empty.

Definition 3.2. The number $|wy^{-1}|$, where y is the α -suffix of w , is called the α -*period* of w , denoted by $\pi_\alpha(w)$.

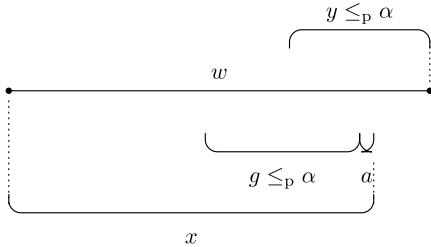
In particular, $|w| - |\alpha| < \pi_\alpha(w) \leq |w|$.

Definition 3.3. The shortest prefix x of w satisfying $\pi_\alpha(x) = \pi_\alpha(w)$ is called the α -*critical prefix* of w .

Remark 3.4. Note that the α -suffix of w can be empty, but it cannot be equal to α . For example, the abb -suffix of $aabb$ is empty. Therefore, the abb -critical prefix of $aabb$ is $aabb$ itself. In general, if α is unbordered and it is a suffix of w , then the α -suffix of w is empty.

Remark 3.5. Let x be the α -critical prefix and y the α -suffix of w . Note that $\pi_\alpha(w) \leq |x| \leq |w|$ and, in particular, $\pi_\alpha(w) = |x| = |w|$ if $y = \varepsilon$.

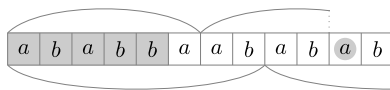
Consider the following illustration of the definitions with $ga \not\leq_p \alpha$.



Remark 3.6. Note that $za = x$, where a is a letter, is the α -critical prefix of w , if, and only if, za is the longest prefix of w satisfying $\pi_\alpha(z) < \pi_\alpha(za)$.

Example 3.7. Consider $w = ababbaababab$ of length 12 and $\alpha = ababb$. The α -suffix of w is $abab$, whence $\pi_\alpha(w) = 8$. The α -critical prefix of w is $ababbaababa$ of length 11, since

$$\pi_\alpha(ababbaababa) = 8, \quad \text{and} \quad \pi_\alpha(ababbaabab) = 6.$$



4. Solution of the Ehrenfeucht–Silberger problem

This entire section is devoted to the proof of the main result of this paper: the solution of the Ehrenfeucht–Silberger problem by Theorem 4.1.

Theorem 4.1. Let $w \in A^*$. If $|w| \geq \frac{7}{3}\tau(w)$ then $\tau(w) = \pi(w)$.

The strategy of the proof is as follows. We define a factorization of w , which allows to detect its long unbordered factors. Two main constructions leading to such factors are given in Section 4.1 (Claim 4.7) and in Section 4.2 (Claim 4.11). In Section 4.3 we show that the main assumption of the theorem, namely that the length of the constructed factors is at most $\frac{3}{7}|w|$, leads to a contradiction, unless $\tau(w) = \pi(w)$.

Note that the claim holds trivially if every letter in w occurs only once because in that case $\tau(w) = \pi(w) = |w|$. We now define the above mentioned factorization of w , which is of central importance to our approach.

Definition 4.2. Let

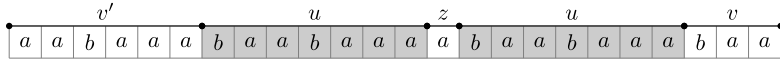
$$w = v'uzuv$$

be a factorization of w such that u is unbordered, $|u| = \tau_2(w)$ and z is of maximum length (recall that $\tau_2(w)$ denotes the maximum length of unbordered factors occurring at least twice in w). Moreover, let us fix

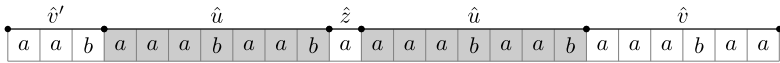
$$t = v \wedge_p zu \quad \text{and} \quad t' = v' \wedge_s uz$$

for the rest of this proof.

The example of long words where the period exceeds the length of the longest unbordered factors by Assous and Pouzet (see p. 2) turns out to highlight the most interesting cases of this proof. We therefore use its instance with $n = 2$ as a running example throughout this section. With this example, the above defined factorization is illustrated by the following figure.

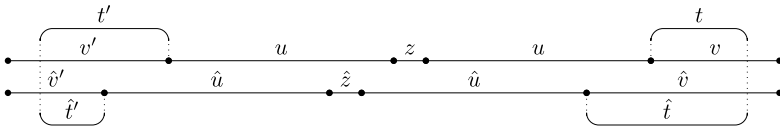


It is clear that such a factorization exists whenever a letter occurs more than once in w . However, it is not necessarily unique. For instance, the factorization on the previous picture competes with the following one.



In general, suppose that $t'u$ contains an unbordered factor \hat{u} , distinct from u but of the same length. Then we have a factorization $\hat{v}'\hat{u}\hat{z}\hat{u}\hat{v}$ of w , which also satisfies the requirements. Note, moreover, that if we define \hat{t} and \hat{t}' analogously to t and t' , then we have

$$t^{-1}v = \hat{t}^{-1}\hat{v} \quad \text{and} \quad v't'^{-1} = \hat{v}'(\hat{t}')^{-1}. \tag{3}$$



In one case (see Case 1, p. 12) it will be important to require that $t'u$ does not contain such an unbordered factor \hat{u} . That is, we shall single out the leftmost possible factorization (within bounds given by the factor $t'uzut$).

Definition 4.3. If $t'u$ does not contain an unbordered factor of length $|u|$ distinct from u , then we shall say that t' is as short as possible.

If this additional assumption is not stated explicitly, then we consider an arbitrary factorization maximizing $|u|$ and $|z|$. The assumption is helpful in view of the following claim.

Claim 4.4. Let t' be as short as possible, and let ϑ be a maximum suffix of $t'u$ w.r.t. some order \triangleleft . Then $|\vartheta| \leq |u|$.

Proof. Suppose that there is a maximum suffix ϑ of $t'u$ strictly longer than u . The prefix \hat{u} of ϑ of length $\pi(\vartheta)$ is unbordered by Remark 2.3. It is of length at least $|u|$, since otherwise u is bordered. From $|u| = \tau_2(w)$ the equality $|\hat{u}| = |u|$ follows since \hat{u} occurs at least twice in w ; a contradiction with the minimality of t' . \square

We start the proof by the following claim, which reveals a long unbordered factor in a special situation.

Claim 4.5. Let ϑ be the maximum suffix of u w.r.t. some order \triangleleft . If $v_0\vartheta$ is a prefix of $t'\vartheta$ for some nonempty word v_0 , then $uzu\vartheta^{-1}v_0\vartheta$ is unbordered.

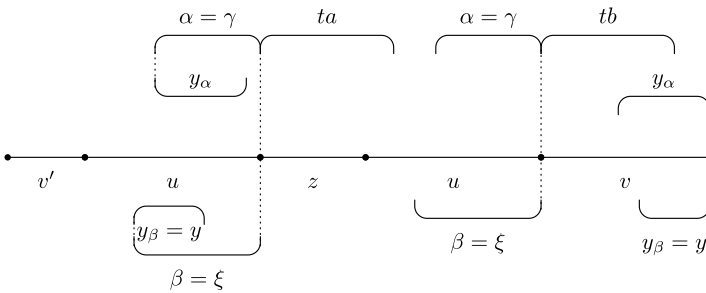
Proof. Suppose on the contrary that $uzu\vartheta^{-1}v_0\vartheta$ has a shortest border h . Note that h is, like every shortest border of a factor in w , not longer than $|u| = \tau_2(w)$. In fact $|h| < |u|$ since $|h| = |u|$ con-

tradicts the maximality of $|z|$. If $|\vartheta| < |h| < |u|$ then ϑ occurs more than once in u contradicting Remark 2.2, which states that a maximum suffix occurs only once in a word. And finally, if $|h| \leq |\vartheta|$ then u is bordered by h since then $h \leq_s \vartheta \leq_s u$; a contradiction which concludes the proof. \square

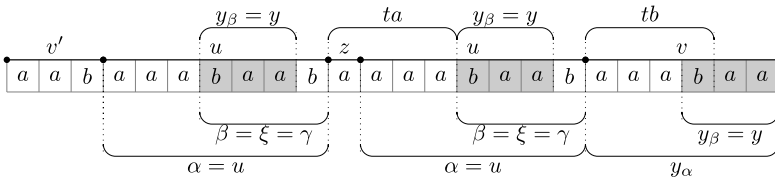
As the reader already noted, our main tool will be considering maximum suffixes w.r.t. certain lexicographic orders. Let us fix some notation.

Definition 4.6. Fix an order \triangleleft . Let α denote the \triangleleft -maximum suffix of u and β the \blacktriangleleft -maximum suffix of u , where \blacktriangleleft is the inverse order of \triangleleft . Let y_α and y_β denote the α - and β -suffix of uv . Moreover, let y be the shorter of y_α and y_β and let ξ be either α or β so that $y = y_\xi$. Let γ denote the shorter of α and β .

Note that $|y| < |\gamma|$ in any case. The following figure shall illustrate the considered setting by an example where $v \neq t$ and $|\alpha| < |\beta|$ and $|y_\alpha| > |y_\beta|$, that is, we have $y = y_\beta$ and $\xi = \beta$ and $\gamma = \alpha$.



The same notation for our running example is depicted next.



It turns out that the proof splits into two main situations according to whether or not $|v| > |ty|$. Each of the cases yields a long unbordered factor of w .

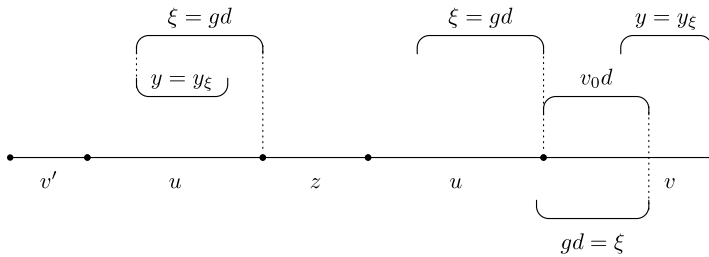
4.1. The first construction

In this subsection we shall suppose $|v| > |ty|$. We consider the ξ -critical prefix of w in order to obtain a long unbordered factor. Note that the following claim holds independently of whether or not $v \neq t$.

Claim 4.7. If $|v| > |ty|$, then $\tau(w) \geq |\gamma zuvy^{-1}|$.

Proof. Suppose $|v| > |ty|$. The inequality implies that the ξ -critical prefix of w can be written as $v'uzuv_0d$, where d is a letter and v_0 is a (possibly empty) word. Let g denote the ξ -suffix of $v'uzuv_0$.

Assume first that $gd = \xi$ as illustrated by the next figure.

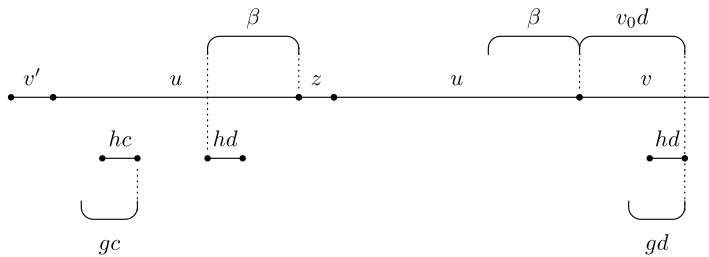


Then the word $uzuv_0d$ is unbordered, by Claim 4.5. Recall that $|\gamma| \leq |\xi| \leq |u|$ and that $|v_0d| \geq |vy^{-1}|$, since $v'uzuv_0d$ is the ξ -critical prefix of w . Therefore we have $\tau(w) \geq |uzuv_0d| \geq |\gamma zuvy^{-1}|$ as claimed.

Suppose next gc is a prefix of ξ with $c \neq d$. (Note that if $gd \neq \xi$, then $c \neq d$ is implied by the definition of the ξ -critical prefix.) We distinguish two cases on the order of c and d in \triangleleft .

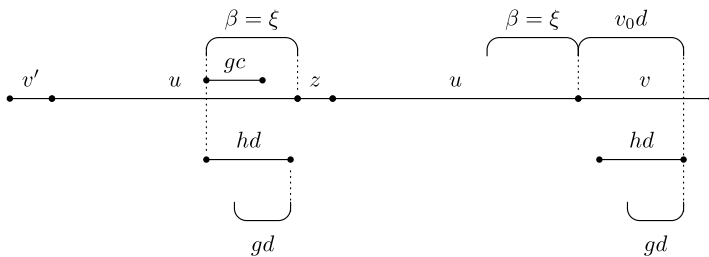
Suppose $c \triangleleft d$ and consider βzuv_0d . Recall that $|\beta| > |y|$ and $|v| \leq |v_0d| + |y|$. Hence, either βzuv_0d is unbordered and we get $\tau(w) \geq |\beta zuv_0d| \geq |\gamma zuvy^{-1}|$ and we are done, or βzuv_0d has a shortest border, say hd .

Suppose $|h| \leq |g|$ and $|h| < |\beta|$ as illustrated by the next figure.



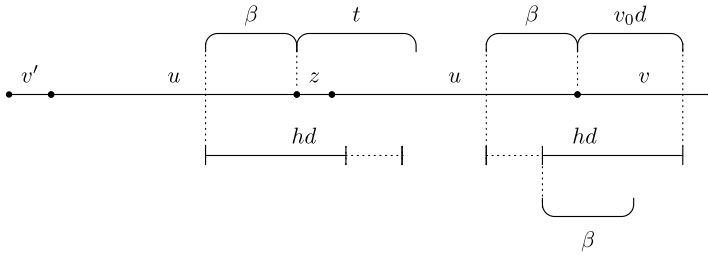
Then hd is a prefix of β and the occurrence of $hc \leq_s gc$ in ξ , and hence also in u , contradicts the maximality of β since $hd \triangleleft hc$.

Suppose $|g| < |h| < |\beta|$ as illustrated by the next figure.



Then gd occurs in u and $\xi = \beta$. Indeed, $gc \leq_p \xi$ gives a contradiction if $\xi = \alpha$ since $gc \triangleleft gd$. But now, h contradicts the assumption that g is the ξ -suffix of $v'uzuv_0$.

It remains that $|\beta| \leq |h|$ which implies $\beta \leq_p h$ as illustrated next.



The choice of u implies $|h| < |u|$. Hence, either $h = \beta v_0$ or the word $uzuv_0h^{-1}\beta$ is unbordered, by Claim 4.5. If $uzuv_0h^{-1}\beta$ is unbordered, then $|u| > |hd|$ and $|v| \leq |v_0d| + |y|$ imply $\tau(w) \geq |uzuv_0h^{-1}\beta| > |\beta zuv_0d| \geq |\gamma zuvy^{-1}|$. If $uzuv_0h^{-1}\beta$ is bordered, then $h = \beta v_0$, which implies $v_0d \leq_p t$ (recall that $t = v \wedge_p zu$), and $|v| \leq |ty|$, since $|v| \leq |v_0d| + |y|$; a contradiction. This completes the case $c \triangleleft d$.

The case $d \triangleleft c$ is similar considering αzuv_0d and the claim is thereby proved. \square

Remark 4.8. Note that we have arguments for v' mirror symmetric to those for v . That is, if we define $\alpha', \beta', \gamma', \xi'$ and γ' for v' analogously, then Claim 4.7 implies the following: If $|v'| > |t'y'|$, then $\tau(w) \geq |y'^{-1}v'uz\gamma'|$.

4.2. The second construction

In this section, we investigate the presence of long unbordered factors in w when $|v| \leq |ty|$. We shall also suppose that v is not a prefix of zu , that is, $t \neq v$.

Definition 4.9. In the rest of the paper, whenever $t \neq v$, the first letter of $t^{-1}v$ will be denoted by b and the first letter of $t^{-1}zu$ by a . In other words, the word ta is a prefix of zuv and tb a prefix of v , with $a \neq b$. Let δ denote the word such that δa is the \triangleleft^a -maximum suffix of $t'uta$ for some fixed order \triangleleft^a such that a is the maximum in A .

The word δ plays an important role in this section, similar to the role of ξ in the previous section. We first point out that every factor of $t'uv$ is strictly less than δa w.r.t. \triangleleft^a if $|v| \leq |ty|$. In particular, δa does not occur in $t'uv$ in such a case.

Claim 4.10. Let f be a factor of $t'uv$. If $|v| \leq |ty|$, then $f \triangleleft^a \delta a$ and $f \neq \delta a$.

Proof. If f occurs in $t'ut$ or y , then the claim follows from the maximality of δa .

Assuming $|v| \leq |ty|$, it remains that there is a prefix $f'b$ of f such that $f' \leq_s t'ut$. Then $f'a \leq_s t'uta$, and the maximality of δa implies $f'a \triangleleft^a \delta a$. The claim now follows from $f'b \leq_p f$ and $f'b \triangleleft^a f'a$. \square

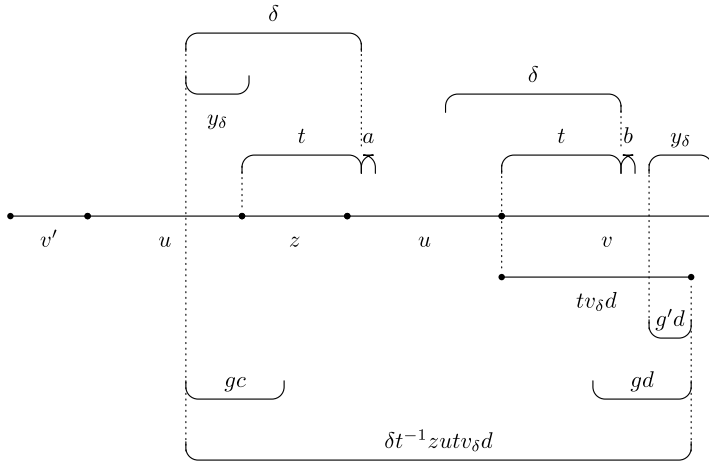
The following claim introduces a further long unbordered factor of w , namely $\delta t^{-1}zuvy_\delta^{-1}$, where y_δ is the δa -suffix of w .

Claim 4.11. The word $\delta t^{-1}zuvy_\delta^{-1}$ is unbordered, and $|y_\delta| < |v| - |t|$.

Proof. If $|y_\delta| \geq |v| - |t|$, then there is a suffix t_0 of $t'ut$ such that t_0b is a prefix of y_δ , and hence, a prefix of δ . This contradicts the maximality of δa w.r.t. \triangleleft^a since t_0a is a suffix of $t'uta$, and hence, a suffix of δa . So, we have $|y_\delta| < |v| - |t|$.

In particular, we have that the δa -critical prefix of $\delta t^{-1}zuv$ is strictly longer than $\delta t^{-1}zut$, whence it can be written as $\delta t^{-1}zutv_\delta d$, where d is a letter. The definition of the critical prefix implies that $|tv_\delta d| \geq |v| - |y_\delta|$. Let g denote the δa -suffix of $\delta t^{-1}zutv_\delta$. Since δa does not occur in $t'u v$ by Claim 4.10, we have that $gd \neq \delta a$. Therefore gc is a prefix of δa and $c \neq d$. Moreover, we deduce $d \triangleleft^a c$ from Claim 4.10.

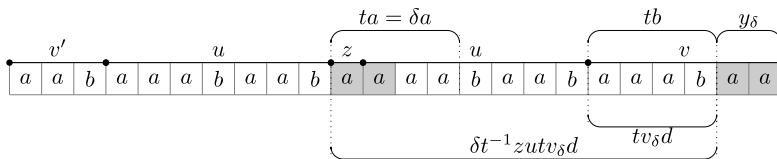
Suppose that $|tv_\delta d| > |v| - |y_\delta|$. Then there is a suffix g' of g such that $g'd$ is a prefix of y_δ , and hence, also of δ . We obtain a contradiction with the maximality of δa , since $g'c$ is a factor of δa . The situation is illustrated in the following figure.



Therefore $|tv_\delta d| = |v| - |y_\delta|$ and $\delta t^{-1}zuvy_\delta^{-1}$ is the δa -critical prefix of $\delta t^{-1}zuv$.

Suppose that $\delta t^{-1}zuvy_\delta^{-1}$ is bordered, and let h be its shortest border. The definition of the critical prefix implies that the δa -period of $\delta t^{-1}zuvy_\delta^{-1}$ is $|\delta t^{-1}zuvy_\delta^{-1}|$, whence $\delta a \leq_p h$. Since $|h| < |u|$, we have that δa occurs in uv contradicting Claim 4.10. \square

Our running example gives the following setting, with $d = b$.

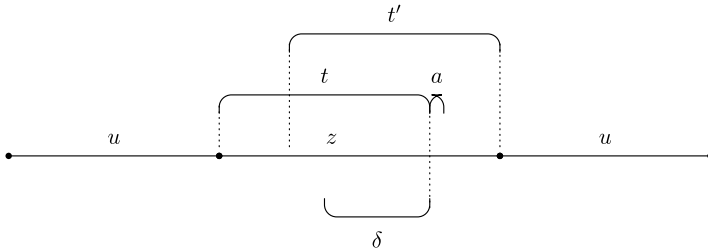


We conclude this section by some auxiliary claims.

Claim 4.12. The word δ satisfies

$$|\delta| > |t| + |t'| - |z|. \tag{4}$$

Proof. Suppose the contrary. Then δ lies within the overlap of ut and $t'u$ in uzu , as illustrated by the following figure.



This contradicts the maximality of δa since it occurs now twice in $t'uta$; see also Remark 2.2. \square

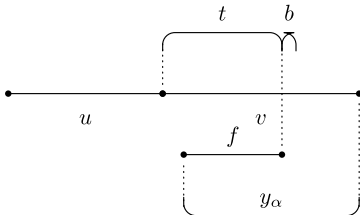
Remark 4.13. Similarly to the mirror symmetric version of Claim 4.7, see Remark 4.8, we have a mirror symmetric setting for Claim 4.11, too.

Provided that $t' \neq v'$, let $a't' \leq_s v'uz$ and $b't' \leq_s v'$ with $a' \neq b'$. Let δ' be defined analogously to δ . If $|v'| \leq |y't'|$, then Claim 4.11 translates to $y_\delta^{-1}v'uzt'^{-1}\delta'$ is unbordered and $|y'_\delta| < |v'| - |t'|$.

Claim 4.7 is formulated for an arbitrary order \triangleleft . Since we want to combine results of Section 4.1 with the present section, we shall identify \triangleleft and \triangleleft^a . In particular, we have $b \triangleleft a$, and δa is the \triangleleft -maximum suffix of $t'uta$.

The conditions $|v| \leq |ty|$ and the $t \neq v$ now imply that utb and y_α are overlapping in uv . Let f be the word such that fb is the overlap. In other words, f is a suffix of ut such that $uv = utf^{-1}y_\alpha$. Since $|y_\alpha| < |\alpha|$, we have

$$|t| > |v| - |\alpha| + |f|. \tag{5}$$



Note that fb is a prefix of y_α , and fa a suffix of uta . We have the following claim (recall Definition 4.3).

Claim 4.14. *If t' is as short as possible, then*

$$|f| > |t| + |t'| - |z|. \tag{6}$$

Proof. Similarly as above for δ , we deduce that fa cannot be a factor of the overlap of t and t' in z , otherwise α is not the \triangleleft -maximum suffix of $t'u$, a contradiction with Claim 4.4 on p. 6. \square

4.3. *Implied inequalities*

In this section, we proceed by case distinction and conclude the proof of the main Theorem 4.1. We shall suppose that $\tau(w) \leq \frac{3}{7}|w|$ and obtain in respective cases either a contradiction, or $\tau(w) = \pi(w)$. In other words, we show that the inequalities derived so far cannot hold, unless $\tau(w) > \frac{3}{7}|w|$ or $\tau(w) = \pi(w)$.

The case distinction is based on whether or not $t = v$ ($t' = v'$); in addition to the main criterion of previous sections, that is, whether or not $|v| > |ty|$ ($|v'| > |t'y'|$).

Case 1: $t \neq v$ or $t' \neq v'$, but not both.

By symmetry, we assume $t \neq v$ and $t' = v'$ in the following. We also assume that t' is as short as possible. Note that this assumption does not change the situation, that is, we still have $t \neq v$ and $t' = v'$; see (3).

Subcase 1.1: $|v| > |ty|$.

Claim 4.7 on p. 7 yields $\tau(w) \geq |\gamma zuvy^{-1}|$. If $|v'| \leq |v|$, then the inequality $|\gamma| > |y|$ implies $\tau(w) > |zuv| \geq \frac{1}{2}|w|$; a contradiction to our assumption. We therefore have $|v'| > |v|$.

Claim 4.4 implies

$$|\gamma z| > |v'|. \tag{7}$$

Indeed, if $|\gamma z| \leq |v'|$, then $\gamma z \leq_s t' = v'$, and hence, there is a maximum suffix ϑ of $t'u$ strictly longer than u contradicting Claim 4.4 (where we let ϑ be the maximum suffix of $t'u$ with respect to the same order to which γ is the maximum suffix of u).

Again, we deduce a contradiction with $\tau(w) \leq \frac{3}{7}|w|$ since $\tau(w) \geq |\gamma zuvy^{-1}| > \frac{1}{2}|w|$ by

$$\begin{aligned} 2(|\gamma| + |z| + |u| + |v| - |y|) &> |v'| + |\gamma| + |z| + 2(|u| + |v| - |y|) \quad (\text{by (7)}) \\ &> |v'| + |z| + 2(|u| + |v|) - |y| \quad (\text{by } |\gamma| > |y|) \\ &> |v'| + |z| + 2|u| + |v| \quad (\text{by } |v| > |y|) \\ &= |w|. \end{aligned}$$

Subcase 1.2: $|v| \leq |ty|$.

In this subcase, we obtain a contradiction by establishing the following set of inequalities that do not have a common solution. Inequality (4) can be transformed into

$$L_1 := |\delta| - |t| - |t'| + |z| - 1 \geq 0.$$

Claim 4.11 on p. 9 yields $|\delta zu| + 1 \leq \frac{3}{7}|w|$ which together with $|w| = |v| + |v'| + 2|u| + |z|$ gives

$$L_2 := 3|v'| + 3|v| - |u| - 4|z| - 7|\delta| - 7 \geq 0.$$

Moreover, since $|y_\delta| \leq |\delta|$, Claim 4.11 yields $|t^{-1}zuv| \leq \frac{3}{7}|w|$ and we obtain

$$L_3 := 7|t| + 3|v'| - 4|v| - 4|z| - |u| \geq 0. \tag{8}$$

The desired contradiction, that is, the fact that the above inequalities have no common solution, follows from

$$21L_1 + 4L_2 + 3L_3 = -7(|uz\delta| + 7),$$

which is obtained keeping in mind that $t' = v'$.

Case 2: $t \neq v$ and $t' \neq v'$.

By symmetry, we can suppose $|v'| \leq |v|$, which implies $\tau(w) < |\gamma zuvy^{-1}|$, see the beginning of Subcase 1.1. Claim 4.7 now yields $|v| \leq |ty|$. As above in Subcase 1.2, we obtain $L_1, L_2, L_3 \geq 0$. We need some more inequalities in this case for we assume $t' \neq v'$. Inequality (5) can be transformed into

$$L_4 := |t| + |\alpha| - |f| - |v| - 1 \geq 0,$$

and the inequality (6) into

$$L_5 := |f| - |t| - |t'| + |z| - 1 \geq 0.$$

We now exploit Remark 4.8 on p. 9. If $\tau(w) \geq |y'^{-1}v'uz\gamma'|$, then using $|y'| < |\gamma'|$ and $\tau(w) \leq \frac{3}{7}|w|$ we obtain the inequality

$$L_6 := 3|v| - 4|v'| - 4|z| - |u| - 7 \geq 0.$$

If, on the other hand, the inequality $|v'| \leq |t'y'|$ holds, then we can use Remark 4.13 and derive the mirror variant of (8), namely, the inequality

$$L_7 := 7|t'| + 3|v| - 4|v'| - 4|z| - |u| \geq 0.$$

We now get

$$\begin{aligned} & 14L_1 + 2L_2 + 2L_3 + 7L_4 + 7L_5 + 3L_7 \\ & = 14L_1 + 2L_2 + 2L_3 + 7L_4 + 7L_5 + 3L_6 + 21(|t'| + 1) = -42 - 7|zu\alpha^{-1}|. \end{aligned}$$

Once again, a sum of nonnegative values turns out to be negative; a contradiction.

Case 3: $t = v$ and $t' = v'$.

This is the only case, in which we prove $\tau(w) = \pi(w)$, instead of a contradiction.

Note that now $|uz|$ is a period of w whence $\pi(w) \leq |uz|$. We can suppose that $\pi(w) > |u|$ since otherwise $\pi(w) = \tau(w) = |u|$, and we are done. Let rs be a critical factorization of u . Then szr is unbordered of length $\pi(w)$, unless r is a prefix, and s is a suffix of z ; see Remark 2.5 on p. 4. Suppose the latter possibility. Now, either one of the words uz and zu is unbordered of length $\pi(w)$ or u is both prefix and suffix of z . We are therefore left with the case $w = v'u^i z' u^j v$, with $i, j \geq 2$, where u is not a suffix of uz' and not a prefix of $z'u$. Note that z' cannot be empty. Moreover, v' is a suffix of $uz = u^i z' u^{j-1}$ and v is a prefix of $zu = u^{i-1} z' u^j$. From the maximality of z we deduce

$$|v'| < |u| \quad \text{and} \quad |v| < |u| \tag{9}$$

which implies that v' is a suffix of u , and v is a prefix of u .

Suppose, without loss of generality, $i \leq j$. Similarly as above, we have that either $sz'u^{j-1}r$ or $z'u^i$ is unbordered, whence $|z'u^i| \leq \tau(w)$. From $|z'u^j| \leq \frac{3}{7}|w|$ we deduce

$$|v'v| \geq \left(\frac{4}{3}j - i\right)|u| + \frac{4}{3}|z'|. \tag{10}$$

If $i < j$, then we obtain from $j \geq 3$ that $|v'v| > 2|u|$; a contradiction with (9). Therefore $i = j \geq 2$.

If v' is a suffix of uz' and v a prefix of $z'u$, then we have $\pi(w) = \tau(w) = |z'u^j|$. Otherwise, we obtain from Case 1 and Case 2 an unbordered factor u_0 of $v'uz'u^jv$ longer than $\frac{3}{7}|v'uz'u^jv|$. We show that u_0 induces a long unbordered factor of w . First, suppose that u_0 is a factor of $uz'u$. Then $|u_0| \leq \pi(uz'u) = |uz'|$. Since the inequalities $|uz'| \geq |u_0| > \frac{3}{7}|v'uz'u^jv|$ imply $|u^jz'| > \frac{3}{7}|w|$, we have a contradiction with $|u^jz'| \leq \tau(w)$. Suppose now that u_0 is not a factor of $uz'u$. Without loss of generality we can suppose that $u_0 = puq$, where p is a suffix of $v'uz'$ and q is a nonempty prefix of v . Then the word $pu^j q$ is a factor of w , and, again, $|pu^j q| > \frac{3}{7}|w|$. It is obvious that $pu^j q$ is unbordered since its shortest border cannot be longer than u (by $|u| = \tau_2(w)$) and each border of length at most $|u|$ would be also a border of puq .

This concludes the proof of Theorem 4.1.

As we mentioned above, Case 3 is the only one allowing $\tau(w) = \pi(w)$. We can therefore extend Theorem 4.1 with the following claim.

Claim 4.15. *Let $|w| \geq \frac{7}{3}\tau(w)$, and let $v'uzuv$ be a factorization of w satisfying Definition 4.2. Then v' is a suffix of zu , and v is a prefix of uz .*

5. Conclusions

The relation between the period $\pi(w)$ of a word w and the length $\tau(w)$ of its longest unbordered factors has been investigated in this paper. Clearly, $\tau(w) \leq \pi(w)$. It is also not difficult to see that $\tau(w) = \pi(w)$ holds for long words, that is, for words, which are much longer than both $\tau(w)$ and $\pi(w)$. The question of interest is: When exactly is a word long enough so that $\tau(w) = \pi(w)$ is enforced? When the word length is expressed w.r.t. $\pi(w)$, it is known that

$$|w| > 2\pi(w) - 2 \text{ implies } \tau(w) = \pi(w).$$

Theorem 4.1 of the present paper makes the complementary statement

$$|w| \geq \frac{7}{3}\tau(w) \text{ implies } \tau(w) = \pi(w).$$

This solves a problem raised first by Ehrenfeucht and Silberger in 1979. Note that the result is independent of the alphabet size.

The bounds $2\tau(w)$ (see [2]) and $3\tau(w)$ (see [11]) have been previously conjectured, and several attempts in proving the latter have been made; see [4,12–14,5]. However, the bound proved above is (asymptotically) tight as demonstrated by an example in [11] with words of length $\frac{7}{3}\tau(w) - 4$ and $\tau(w) < \pi(w)$. For the sake of clarity we did not try to make the additive constant optimal in this paper. We only note that our arguments can be easily modified to obtain that already $|w| > \frac{7}{3}\tau(w) - \frac{8}{3}$ implies $\tau(w) = \pi(w)$. We do not consider this value of the additive constant to be too interesting since we conjecture that the example by Assous and Pouzet is optimal, that is

$$|w| > \frac{7}{3}\tau(w) - 3 \text{ implies } \tau(w) = \pi(w),$$

and, moreover, if $|w| = \frac{7}{3}\tau(w) - 4$ and $\tau(w) \neq \pi(w)$, then w is of the form given by (1).

Apart from the actual result, we would like to point out the proof techniques used to solve the Ehrenfeucht–Silberger problem. In particular, the notion of α -critical prefix of a word w (Definition 3.3) is used to find long unbordered factors in words with a large period, that is, words that do not have much of a global structure. We are confident that the investigation of α -critical prefixes of a word will lead to more insights in its structure, for example w.r.t. its local periods.

Acknowledgments

The authors would like to thank the anonymous referees for their valuable comments which greatly helped to improve the presentation.

References

- [1] Y. Césari, M. Vincent, Une caractérisation des mots périodiques, C. R. Acad. Sci. Paris Sér. A 286 (1978) 1175–1177.
- [2] A. Ehrenfeucht, D.M. Silberger, Periodicity and unbordered segments of words, Discrete Math. 26 (2) (1979) 101–109.
- [3] J.-P. Duval, Périodes et répétitions des mots du monoïde libre, Theoret. Comput. Sci. 9 (1) (1979) 17–26.
- [4] J.-P. Duval, Relationship between the period of a finite word and the length of its unbordered segments, Discrete Math. 40 (1) (1982) 31–44.
- [5] Š. Holub, A proof of the extended Duval's conjecture, Theoret. Comput. Sci. 339 (1) (2005) 61–67.
- [6] T. Harju, D. Nowotka, Periodicity and unbordered words: A proof of the extended Duval conjecture, J. ACM 54 (4) (2007), Article 20, 20 pp.
- [7] M. Crochemore, D. Perrin, Two-way string-matching, J. ACM 38 (3) (1991) 651–675.
- [8] D. Breslauer, Saving comparisons in the Crochemore–Perrin string-matching algorithm, Theoret. Comput. Sci. 158 (1–2) (1996) 177–192.
- [9] D. Breslauer, Z. Galil, Real-time streaming string matching, in: 22nd Annual Symposium on Combinatorial Pattern Matching, 2011, pp. 162–172.
- [10] M. Lothaire, Combinatorics on Words, Encyclopedia Math. Appl., vol. 17, Addison–Wesley, Reading, MA, 1983.
- [11] R. Assous, M. Pouzet, Une caractérisation des mots périodiques, Discrete Math. 25 (1) (1979) 1–5.
- [12] F. Mignosi, L.Q. Zamboni, A note on a conjecture of Duval and Sturmian words, Theor. Inform. Appl. 36 (1) (2002) 1–3.
- [13] J.-P. Duval, T. Harju, D. Nowotka, Unbordered factors and Lyndon words, Discrete Math. 308 (11) (2008) 2261–2264 (submitted in 2002).

- [14] T. Harju, D. Nowotka, Minimal Duval extensions, *Internat. J. Found. Comput. Sci.* 15 (2) (2004) 349–354.
- [15] T. Harju, D. Nowotka, Periodicity and unbordered words, in: *STACS 2004 (Montpellier)*, in: *Lecture Notes in Comput. Sci.*, vol. 2996, Springer-Verlag, Berlin, 2004, pp. 294–304.
- [16] M.-P. Schützenberger, A property of finitely generated submonoids of free monoids, in: *Algebraic Theory of Semigroups, Proc. Sixth Algebraic Conf., Szeged, 1976*, in: *Colloq. Math. Soc. János Bolyai*, vol. 20, North-Holland, Amsterdam, 1979, pp. 545–576.