

Proceedings of the Prague Stringology Conference 2014

Edited by Jan Holub and Jan Žďárek



September 2014



Prague Stringology Club
<http://www.stringology.org/>

Two Squares Canonical Factorization*

Haoyue Bai¹, Frantisek Franek¹, and William F. Smyth^{1,2}

¹ Department of Computing and Software
McMaster University, Hamilton, Ontario, Canada
{baih3,franek,smyth}@mcmaster.ca

² School of Computer Science & Software Engineering
University of Western Australia

Abstract. We present a new combinatorial structure in a string: a canonical factorization for any two squares that occur at the same position and satisfy some size restrictions. We believe that this canonical factorization will have application to related problems such as the New Periodicity Lemma, Crochemore-Rytter Three Squares Lemma, and ultimately the maximum-number-of-runs conjecture.

Keywords: string, primitive string, square, double square, factorization

1 Introduction

In 1995 Crochemore and Rytter [2] considered three distinct squares, all prefixes of a given string \mathbf{x} , and proved the *Three Squares Lemma* stating that, subject to certain restrictions, the largest of the three was at least the length of the sum of the other two. In 2006 Fan *et al.* [4] considered a special case of such two squares prefixes of \mathbf{x} with a third square possibly offset some distance to the right; they proved a *New Periodicity Lemma* describing conditions under which the third square could not exist. Since that time there has been considerable work done [1,5,6,8] in an effort to specify more precisely the combinatorial structure of the string in the neighbourhood of such two squares.

In this paper we present a unique *canonical factorization* into primitive strings of what we call *double squares* – i.e. two squares starting at the same position and satisfying some size restrictions. The notion of double squares and their unique factorization can be traced to Lam [7]. A version of the factorization for more specific double squares was presented in [3]. Here we present it in full generality. In conclusion we indicate how this result can be applied to the proof of New Periodicity Lemma.

2 Preliminaries

In this section we develop the basic combinatorial tools that will be used to determine a canonical factorization for a double square. Chief among these are the Synchronization Principle (see Lemma 2), and the Common Factor Lemma (see Lemma 3), that lead to the main result, the Two Squares Factorization Lemma (see Lemma 6).

A *string* \mathbf{x} is a finite sequence of symbols, called *letters*, drawn from a (finite or infinite) set Σ , called the *alphabet*. The length of the sequence is called the *length* of \mathbf{x} , denoted $|\mathbf{x}|$. Sometimes for convenience we represent a string \mathbf{x} of length n as an array $\mathbf{x}[1..n]$. The string of length zero is called the *empty string*, denoted ε . If a string $\mathbf{x} = \mathbf{u}\mathbf{v}\mathbf{w}$, where \mathbf{u} , \mathbf{v} , \mathbf{w} are strings, then \mathbf{u} (respectively, \mathbf{v} , \mathbf{w}) is said to

* This work was supported by the *Natural Sciences and Engineering Research Council of Canada*

be a **prefix** (respectively, **substring**, **suffix**) of \mathbf{x} ; a **proper prefix** (respectively, **proper substring**, **proper suffix**) if $|\mathbf{u}| < |\mathbf{x}|$ (respectively, $|\mathbf{v}| < |\mathbf{x}|$, $|\mathbf{w}| < |\mathbf{x}|$). A substring is also called a **factor**. Given strings \mathbf{u} and \mathbf{v} , $\text{lcp}(\mathbf{u}, \mathbf{v})$ (respectively, $\text{lcs}(\mathbf{u}, \mathbf{v})$) is the **longest common prefix** (respectively, **longest common suffix**) of \mathbf{u} and \mathbf{v} .

If \mathbf{x} is a concatenation of $k \geq 2$ copies of a nonempty string \mathbf{u} , we write $\mathbf{x} = \mathbf{u}^k$ and say that \mathbf{x} is a **repetition**; if $k = 2$, we say that $\mathbf{x} = \mathbf{u}^2$ is a **square**; if there exist no such integer k and no such \mathbf{u} , we say that \mathbf{x} is **primitive**. If $\mathbf{x} = \mathbf{v}^2$ has a proper prefix \mathbf{u}^2 , $|\mathbf{su}| < |\mathbf{v}| < 2|\mathbf{u}|$, we say that \mathbf{x} is a **double square** and write $\mathbf{x} = \text{DS}(\mathbf{u}, \mathbf{v})$. A square \mathbf{u}^2 such that \mathbf{u} has no square prefix is said to be **regular**.

For $\mathbf{x} = \mathbf{x}[1..n]$, $1 \leq i < j \leq j+k \leq n$, the string $\mathbf{x}[i+k..j+k]$ is a **right cyclic shift** by k positions of $\mathbf{x}[i..j]$ if $\mathbf{x}[i] = \mathbf{x}[j+1]$, \dots , $\mathbf{x}[i+k-1] = \mathbf{x}[j+k]$. Equivalently, we can say that $\mathbf{x}[i..j]$ is a **left cyclic shift** by k positions of $\mathbf{x}[i+k..j+k]$. When it is clear from the context, we may leave out the number of positions and just speak of a cyclic shift.

Strings \mathbf{uv} and \mathbf{vu} are **conjugates**, written $\mathbf{uv} \sim \mathbf{vu}$. We also say that \mathbf{vu} is the $|\mathbf{u}|^{\text{th}}$ **rotation** of \mathbf{x} , written $R_{|\mathbf{u}|}(\mathbf{x})$, or the $-|\mathbf{v}|^{\text{th}}$ **rotation** of \mathbf{x} , written $R_{-|\mathbf{v}|}(\mathbf{x})$, while $R_0(\mathbf{x}) = R_{-|\mathbf{x}|}(\mathbf{x}) = \mathbf{x}$ is a **primitive rotation**. Similarly as for the cyclic shift, when it is clear from the context, we may leave out the number of rotations and just speak of a rotation. Note that all cyclic shifts are conjugates, but not the other way around.

In the following lemma, the symbol $|$ denotes *divisibility*, i.e. $a | b$ means that a is divisible by b .

Lemma 1 [9, Lemma 1.4.2] *Let \mathbf{x} be a string of length n and minimum period $\pi \leq n$, and let $j = 1, \dots, n-1$ be an integer. Then $R_j(\mathbf{x}) = \mathbf{x}$ if and only if \mathbf{x} is not primitive ($\pi < n$, $\pi | n$) and $j | \pi$.*

The following results (Lemmas 2–6) are based on the development given in [3]. Though Lemmas 2 and 3 are folklore, we include their proofs.

Lemma 2 (Synchronization Principle) *The primitive string \mathbf{x} occurs exactly p times in $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$, where p is a nonnegative integer and \mathbf{x}_1 (respectively, \mathbf{x}_2) is a proper prefix (respectively, proper suffix) of \mathbf{x} .*

Proof. From Lemma 1 a rotation $R_j(\mathbf{x})$ of \mathbf{x} can equal \mathbf{x} only if \mathbf{x} is not primitive. Since here \mathbf{x} is primitive, the only occurrences of \mathbf{x} are exactly those determined by \mathbf{x}^p . \square

Lemma 3 (Common Factor Lemma) *Suppose that \mathbf{x} and \mathbf{y} are primitive strings, where \mathbf{x}_1 (respectively, \mathbf{y}_1) is a proper prefix and \mathbf{x}_2 (respectively, \mathbf{y}_2) a proper suffix of \mathbf{x} (respectively, \mathbf{y}). If for nonnegative integers p and q , $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$ and $\mathbf{y}_2\mathbf{y}^q\mathbf{y}_1$ have a common factor of length $|\mathbf{x}|+|\mathbf{y}|$, then $\mathbf{x} \sim \mathbf{y}$.*

Proof. First consider the special case $\mathbf{x}_1 = \mathbf{x}_2 = \mathbf{y}_1 = \mathbf{y}_2 = \varepsilon$, where \mathbf{x}^p , \mathbf{y}^q have a common prefix \mathbf{f} of length $|\mathbf{x}|+|\mathbf{y}|$. We show that in this case $\mathbf{x} = \mathbf{y}$.

Observe that \mathbf{f} has prefixes \mathbf{x} and \mathbf{y} , so that if $|\mathbf{x}| = |\mathbf{y}|$, then $\mathbf{x} = \mathbf{y}$, as required. Therefore suppose WLOG that $|\mathbf{x}| < |\mathbf{y}|$. Note that $\mathbf{y} \neq \mathbf{x}^k$ for any integer $k \geq 2$, since otherwise \mathbf{y} would not be primitive, contradicting the hypothesis of the lemma.

Hence there exists $k \geq 1$ such that $k|\mathbf{x}| < |\mathbf{y}|$ and $(k+1)|\mathbf{x}| > |\mathbf{y}|$. But since $\mathbf{f} = \mathbf{y}\mathbf{x}$, it follows that

$$R_{|\mathbf{y}| - k|\mathbf{x}|}(\mathbf{x}) = \mathbf{x},$$

again by Lemma 1 contrary to the assumption that \mathbf{x} is primitive. We conclude that $|\mathbf{x}| \not< |\mathbf{y}|$, hence that $|\mathbf{x}| = |\mathbf{y}|$ and $\mathbf{x} = \mathbf{y}$, as required.

Now consider the general case, where \mathbf{f} of length $|\mathbf{x}| + |\mathbf{y}|$ is a common factor of $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$ and $\mathbf{y}_2\mathbf{y}^q\mathbf{y}_1$. Then $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1 = \mathbf{u}\mathbf{f}\mathbf{u}'$ for some \mathbf{u} and \mathbf{u}' . If $|\mathbf{u}| \geq |\mathbf{x}|$, then \mathbf{f} is a factor of $\mathbf{x}_1\mathbf{x}^{p-1}\mathbf{x}_2$, and so we can assume WLOG that $|\mathbf{u}| < |\mathbf{x}|$. Setting $\tilde{\mathbf{x}} = R_{|\mathbf{u}|}(\mathbf{x})$, we see that \mathbf{f} is a prefix of $\tilde{\mathbf{x}}^p$.

Similarly, by setting $\mathbf{y}_2\mathbf{y}^q\mathbf{y}_1 = \mathbf{v}\mathbf{f}\mathbf{v}'$, we can assume that $|\mathbf{v}| < |\mathbf{y}|$, hence that \mathbf{f} is also a prefix of $\tilde{\mathbf{y}}^q$ for $\tilde{\mathbf{y}} = R_{|\mathbf{v}|}(\mathbf{y})$. But this is just the special case considered above, for which $\tilde{\mathbf{x}} = \tilde{\mathbf{y}}$. Since $\mathbf{x} \sim \tilde{\mathbf{x}}$ and $\mathbf{y} \sim \tilde{\mathbf{y}}$, the result follows. \square

Note that Lemma 3 could be equivalently stated in a more general form:

Lemma 4 *Suppose that \mathbf{x} and \mathbf{y} are strings where \mathbf{x}_1 (respectively, \mathbf{y}_1) is a proper prefix and \mathbf{x}_2 (respectively, \mathbf{y}_2) a proper suffix of \mathbf{x} (respectively, \mathbf{y}). If for nonnegative integers p and q , $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$ and $\mathbf{y}_2\mathbf{y}^q\mathbf{y}_1$ have a common factor of length $|\mathbf{x}| + |\mathbf{y}|$, then the primitive root $\bar{\mathbf{x}}$ of \mathbf{x} and the primitive root $\bar{\mathbf{y}}$ of \mathbf{y} are conjugates.*

The Common Factor Lemma gives rise to the following useful corollary:

Lemma 5 *Suppose that \mathbf{x} and \mathbf{y} are primitive strings, and that p and q are positive integers.*

- (a) *If $\mathbf{x}^p = \mathbf{y}^q$, then $\mathbf{x} = \mathbf{y}$ and $p = q$.*
- (b) *If \mathbf{x}_1 (respectively, \mathbf{y}_1) is a proper prefix of \mathbf{x} (respectively, \mathbf{y}) and $\mathbf{x}^p\mathbf{x}_1 = \mathbf{y}^q\mathbf{y}_1$ for $p \geq 2$, $q \geq 2$, then $\mathbf{x} = \mathbf{y}$, $\mathbf{x}_1 = \mathbf{y}_1$ and $p = q$.*

Proof. For (a), first consider $p = 1$, thus $\mathbf{x} = \mathbf{y}^q$. Since \mathbf{x} is primitive, therefore $q = 1$ and $\mathbf{x} = \mathbf{y}$, as required. Similarly for $q = 1$. Suppose then that $p, q \geq 2$. This means that \mathbf{x}^p and $\mathbf{y}^q = \mathbf{x}^p$ have a common factor of length $p|\mathbf{x}| = q|\mathbf{y}| \geq |\mathbf{x}| + |\mathbf{y}|$, so that by Lemma 3 $\mathbf{x} \sim \mathbf{y}$. Hence $|\mathbf{x}| = |\mathbf{y}|$ and so $\mathbf{x} = \mathbf{y}$.

For (b), since again $p \geq 2$, $q \geq 2$, it follows as in (a) that $\mathbf{x}^p\mathbf{x}_1 = \mathbf{y}^q\mathbf{y}_1$ has a common factor of length at least $|\mathbf{x}| + |\mathbf{y}|$, hence the result. \square

Note that in Lemma 5(b) the requirement $p \geq 2$, $q \geq 2$ is essential. For instance, $\mathbf{x} = aabb$, $\mathbf{x}_1 = aa$ and $p = 2$ yields $\mathbf{x}^p\mathbf{x}_1 = aabbaabbaa$, identical to $\mathbf{y}^q\mathbf{y}_1$ produced by $\mathbf{y} = aabbaabba$, $\mathbf{y}_1 = a$ and $q = 1$ — but of course $\mathbf{x} \neq \mathbf{y}$.

3 Main Result – Two Squares Factorization Lemma

The next lemma specifies the structure imposed by the occurrence of two squares at the same position in a string. This structure has been described before, see [3,4,5,6,7], but not as precisely and with more assumptions required; above all, Lemma 6 establishes the uniqueness of the breakdown.

Lemma 6 (Two Squares Factorization Lemma) *For a double square $DS(\mathbf{u}, \mathbf{v})$, there exists a unique primitive string \mathbf{u}_1 such that $\mathbf{u} = \mathbf{u}_1^{e_1}\mathbf{u}_2$ and $\mathbf{v} = \mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_2}$, where \mathbf{u}_2 is a possibly empty proper prefix of \mathbf{u}_1 and e_1, e_2 are integers such that $e_1 \geq e_2 \geq 1$. Moreover,*

- (a) if $|\mathbf{u}_2| = 0$, then $e_1 > e_2 \geq 1$;
- (b) if $|\mathbf{u}_2| > 0$, then \mathbf{v} is primitive, and if in addition $e_1 \geq 2$, then \mathbf{u} also is primitive.

In both cases, the factorization is unique.

Proof. If we have \mathbf{u}^k , $k \geq 2$, we refer to the first copy of \mathbf{u} as $\mathbf{u}_{[1]}$, to the second copy of \mathbf{u} as $\mathbf{u}_{[2]}$ etc.

Let \mathbf{z} be the nonempty proper prefix of $\mathbf{u}_{[2]}$ that is in addition a suffix \mathbf{z} of $\mathbf{v}_{[1]}$. But then \mathbf{z} is also a prefix of $\mathbf{v}_{[1]}$, hence of $\mathbf{v}_{[2]}$; thus if $|\mathbf{u}| \geq 2|\mathbf{z}|$, it follows that \mathbf{z}^2 is a prefix of \mathbf{u} . In general, there exists an integer $k = \lfloor |\mathbf{u}|/|\mathbf{z}| \rfloor \geq 1$ such that $\mathbf{u} = \mathbf{z}^k \mathbf{z}'$ for some proper suffix \mathbf{z}' of \mathbf{z} . Let \mathbf{u}_1 be the primitive root of \mathbf{z} , so that $\mathbf{z} = \mathbf{u}_1^{e_2}$ for some integer $e_2 \geq 1$. Therefore, for some $e_1 \geq e_2 k$ and some prefix \mathbf{u}_2 of \mathbf{u}_1 , $\mathbf{u} = \mathbf{u}_1^{e_1} \mathbf{u}_2$ and $\mathbf{v} = \mathbf{u} \mathbf{z} = \mathbf{u}_1^{e_1} \mathbf{u}_2 \mathbf{u}_1^{e_2}$, as required. To prove uniqueness we consider two cases:

- (i) $|\mathbf{u}_2| = 0$
Here $\mathbf{u} = \mathbf{u}_1^{e_1}$ and $\mathbf{v} = \mathbf{u}_1^{e_1+e_2}$, so that $\mathbf{x} = \mathbf{u}_1^{2(e_1+e_2)}$. Since $|\mathbf{v}| < 2|\mathbf{u}|$ and $e_1 \geq e_2$, it follows that $e_1 > e_2$. The uniqueness of \mathbf{u}_1 is a consequence of Lemma 5(a).
- (ii) $|\mathbf{u}_2| > 0$
Suppose the choice of \mathbf{u}_1 is not unique. Then there exists some primitive string \mathbf{w}_1 with proper prefix \mathbf{w}_2 , together with integers $f_1 \geq f_2 \geq 1$, such that $\mathbf{u} = \mathbf{w}_1^{f_1} \mathbf{w}_2$ and $\mathbf{v} = \mathbf{w}_1^{f_1} \mathbf{w}_2 \mathbf{w}_1^{f_2}$. If both $e_1 \geq 2$ and $f_1 \geq 2$, it follows from Lemma 5(b) that $\mathbf{u}_1 = \mathbf{w}_1$ and $e_1 = f_1$. If $e_1 = f_1 = 1$, we observe that $\mathbf{v} = \mathbf{u} \mathbf{u}_1 = \mathbf{u} \mathbf{w}_1$, so that again $\mathbf{u}_1 = \mathbf{w}_1$. In the only remaining case, exactly one of e_1, f_1 equals 1: therefore suppose WLOG that $f_1 > e_1 = 1$. Then $\mathbf{u} = \mathbf{u}_1 \mathbf{u}_2 = \mathbf{w}_1^{f_1} \mathbf{w}_2$ and $\mathbf{v} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 = \mathbf{w}_1^{f_1} \mathbf{w}_2 \mathbf{w}_1^{f_2}$, so that $\mathbf{u}_1 = \mathbf{w}_1^{f_2}$. But since \mathbf{u}_1 is primitive, this forces $f_2 = 1$ and $\mathbf{u}_1 = \mathbf{w}_1$, which, since $\mathbf{u}_1 \mathbf{u}_2 = \mathbf{w}_1^{f_1} \mathbf{w}_2 = \mathbf{u}_1^{f_1} \mathbf{w}_2$, implies that $f_1 = 1$, a contradiction. Thus all cases have been considered, and \mathbf{u}_1 is unique.

We now show that \mathbf{v} is primitive. Suppose the contrary, so there exists some primitive \mathbf{w} and an integer $k \geq 2$ such that $\mathbf{v} = \mathbf{w}^k$. It follows that $|\mathbf{w}| \leq |\mathbf{v}|/2 \leq |\mathbf{u}_1^{e_1}| + |\mathbf{u}_2|$. Note that

$$\mathbf{w}^{2k} = \mathbf{v}^2 = \mathbf{u}_1^{e_1} \mathbf{u}_2 \mathbf{u}_1^{e_1+e_2} \mathbf{u}_2 \mathbf{u}_1^{e_2}, \tag{1}$$

so that \mathbf{w}^{2k} and $\mathbf{u}_1^{e_1+e_2} \mathbf{u}_2$ have a common factor $\mathbf{u}_1^{e_1+e_2} \mathbf{u}_2$ of length

$$(|\mathbf{u}_1^{e_1}| + |\mathbf{u}_2|) + |\mathbf{u}_1^{e_2}| \geq |\mathbf{w}| + |\mathbf{u}_1|.$$

Thus we can apply Common Factor Lemma 3 to conclude that $\mathbf{w} \sim \mathbf{u}_1$, thus by (1) that $\mathbf{w} = \mathbf{u}_1$. But (1) then requires that the primitive string $\mathbf{u}_1 = \mathbf{u}_2 \bar{\mathbf{u}}_2$ aligns with $\mathbf{u}_2 \mathbf{u}_1$, and so $\bar{\mathbf{u}}_2$ is a prefix of \mathbf{u}_1 , in contradiction to Lemma 1. We conclude that \mathbf{v} is primitive.

Now suppose in addition that $e_2 \geq 2$, but that \mathbf{u} is not primitive. Then there exists some primitive \mathbf{w} and some integer $k \geq 2$ such that $\mathbf{u} = \mathbf{w}^k$. Hence $|\mathbf{w}| \leq |\mathbf{u}|/2 = (|\mathbf{u}_1^{e_1}| + |\mathbf{u}_2|)/2 < |\mathbf{u}_1^{e_1-1}| + |\mathbf{u}_2|$, since $e_1 \geq 2$ and $|\mathbf{u}_2| > 0$. Therefore, since $\mathbf{u}_1^{e_1} \mathbf{u}_2$ is a prefix of $\mathbf{u}^2 = \mathbf{w}^{2k}$, and since $e_2 \geq 1$ by Lemma 6, \mathbf{w}^{2k} and $\mathbf{u}_1^{e_1+e_2}$ have a common prefix $\mathbf{u}_1^{e_1} \mathbf{u}_2$. Note that $|\mathbf{u}_1^{e_1} \mathbf{u}_2| \geq |\mathbf{v}| + |\mathbf{u}_1|$, so that again applying Common Factor Lemma 3, we conclude that $\mathbf{u}_1 = \mathbf{w}$. This in turn implies $\mathbf{u} = \mathbf{u}_1^{e_1} \mathbf{u}_2 = \mathbf{u}_1^k$, impossible since $0 < |\mathbf{u}_2| < |\mathbf{u}_1|$. Therefore \mathbf{u} is primitive, as required.

Finally we remark that since \mathbf{u}_1 is a uniquely determined primitive string, therefore \mathbf{u}_2 , e_1 and e_2 are also uniquely determined. \square

The following examples show that the statement of the lemma is sharp:

- (a) The second part of Lemma 6(b) requires that $e_1 \geq 2$. To see that this condition is not necessary, consider $\mathbf{v}^2 = abaababaab$, where $\mathbf{u} = (ab)a$, $\mathbf{v} = (ab)a(ab)$, so that $\mathbf{u}_1 = ab$, $\mathbf{u}_2 = a$, $e_1 = e_2 = 1$, but \mathbf{u} is primitive.
- (b) On the other hand, consider $\mathbf{v}^2 = abaabaabaabaabaabaab$, where $\mathbf{u} = (aba)^2 = (abaab)a$, $\mathbf{v} = (abaab)a(abaab)$, so that $\mathbf{u}_1 = abaab$, $\mathbf{u}_2 = a$, $e_1 = e_2 = 1$, where now \mathbf{u}_1 is *not* primitive.

Lemma 6 gives credence to the following definition of terminology and notation:

Definition 7 For a double square $DS(\mathbf{u}, \mathbf{v})$ we call the unique factorization $\mathbf{v}^2 = \mathbf{u}_1^{e_1} \mathbf{u}_2 \mathbf{u}_1^{e_1+e_2} \mathbf{u}_2 \mathbf{u}_1^{e_2}$ guaranteed by Lemma 6, the **canonical factorization** of $DS(\mathbf{u}, \mathbf{v})$ and denote it by $DS(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$. The symbol $\bar{\mathbf{u}}_2$ denotes the suffix of \mathbf{u}_1 such that $\mathbf{u}_1 = \mathbf{u}_2 \bar{\mathbf{u}}_2$.

Lemma 6 also gives rise to a number of important observations:

Observation 8 In Lemma 6, $|\mathbf{u}_2| > 0$ if any one of the following conditions holds:

- (a) \mathbf{v} is primitive;
- (b) \mathbf{u} is primitive;
- (c) there is no other occurrence of \mathbf{u}^2 farther to the right in \mathbf{v}^2 (\mathbf{u}^2 is rightmost);
- (d) \mathbf{u}^2 is regular.

Moreover:

- (e) $|\mathbf{u}_2| > 0$ if and only if \mathbf{v} is primitive;
- (f) If \mathbf{u}^2 is regular, then $e_1 = e_2 = 1$ and \mathbf{u}_1 is regular.

Proof.

- (a) $|\mathbf{u}_2| = 0$ implies \mathbf{v} not primitive.
- (b) $|\mathbf{u}_2| = 0$ implies \mathbf{u} not primitive.
- (c) $|\mathbf{u}_2| = 0$ implies $\mathbf{u}^2 = \mathbf{u}_1^{2e_1}$, which occurs twice in $\mathbf{v}^2 = \mathbf{u}_1^{2(e_1+e_2)}$, in particular as a suffix.
- (d) Since \mathbf{u}^2 is regular, therefore \mathbf{u} is primitive, so that by (b), $|\mathbf{u}_2| > 0$.
- (e) By (a), primitive \mathbf{v} implies $|\mathbf{u}_2| > 0$; by Lemma 6, $|\mathbf{u}_2| > 0$ implies that \mathbf{v} is primitive.
- (f) By (d), regular \mathbf{u}^2 implies $|\mathbf{u}_2| > 0$, so that $\mathbf{u} = \mathbf{u}_1^{e_1} \mathbf{u}_2$, which is regular only if $e_1 = e_2 = 1$ and \mathbf{u}_1 is regular. \square

In the context of Observation 8(f), consider the double square $DS(\mathbf{u}, \mathbf{v})$ where $\mathbf{u} = aabaa$, $\mathbf{v} = aabaaaab$. In this case, we find $\mathbf{u}_1 = aab$, $\mathbf{u}_2 = aa$, $e_1 = e_2 = 1$, but observe that \mathbf{u} has prefix a^2 , so \mathbf{u}^2 is not regular. Thus the condition $e_1 = 1$ is more general than the requirement that \mathbf{u}^2 be regular.

Now, following [3], consider the case $|\mathbf{u}_2| > 0$ of Lemma 6 and set $\mathbf{u}_1 = \mathbf{u}_2 \bar{\mathbf{u}}_2$. Thus \mathbf{v}^2 becomes

$$\begin{aligned} \mathbf{v}^2 &= (\mathbf{u}_2 \bar{\mathbf{u}}_2)^{e_1} \mathbf{u}_2 (\mathbf{u}_2 \bar{\mathbf{u}}_2)^{e_1+e_2} \mathbf{u}_2 (\mathbf{u}_2 \bar{\mathbf{u}}_2)^{e_2} \\ &= (\mathbf{u}_2 \bar{\mathbf{u}}_2)^{e_1-1} \mathbf{u}_2 (\text{IF}) (\mathbf{u}_2 \bar{\mathbf{u}}_2)^{e_1+e_2-2} \mathbf{u}_2 (\text{IF}) (\mathbf{u}_2 \bar{\mathbf{u}}_2)^{e_2-1} \end{aligned} \quad (2)$$

where $\text{IF} = \bar{\mathbf{u}}_2 \mathbf{u}_2 \mathbf{u}_2 \bar{\mathbf{u}}_2 = R_{|\mathbf{u}_2|}(\mathbf{u}_1) \mathbf{u}_1$ is called the **inversion factor**.

Lemma 9 Consider a double square $DS(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with a non-empty \mathbf{u}_2 . Then the inversion factor IF have exactly two occurrences in \mathbf{v}^2 exactly a distance of $|\mathbf{v}|$ apart as shown in (2).

Proof. If IF occurs elsewhere in \mathbf{v}^2 , by the Synchronization principle its subfactor $\mathbf{u}_2\bar{\mathbf{u}}_2$ must align with an occurrence of $\mathbf{u}_2\bar{\mathbf{u}}_2$ as it is primitive. Thus, its subfactor $\bar{\mathbf{u}}_2\mathbf{u}_2$ must align with $\mathbf{u}_2\bar{\mathbf{u}}_2$, contradicting the primitiveness of $\mathbf{u}_2\bar{\mathbf{u}}_2$, see Lemma 1. \square

The quantity $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ gives the maximal number of positions the structures $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}$ and $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2}$ can be cyclically shifted to the left in \mathbf{v}^2 , while $\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ gives the maximal number of positions the structures $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1}$ and $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}$ can be cyclically shifted to the right. In [3], the following lemma limiting the size of $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ was given.

Lemma 10 ([3]) *Considering $\mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2\mathbf{u}_1^{e_2}$, where \mathbf{u}_1 is primitive and \mathbf{u}_2 is a non-empty proper prefix of \mathbf{u}_1 , $e_1 \geq e_2 \geq 1$, and $\bar{\mathbf{u}}_2$ a suffix of \mathbf{u}_1 so that $\mathbf{u}_1 = \mathbf{u}_2\bar{\mathbf{u}}_2$, then $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) \leq |\mathbf{u}_1| - 2$.*

In fact, in [3] the inversion factor is defined more generally as any factor $\bar{\mathbf{w}}\mathbf{w}\mathbf{w}\bar{\mathbf{w}}$ of \mathbf{v}^2 such that $|\mathbf{w}| = |\mathbf{u}_2|$ and $|\bar{\mathbf{w}}| = |\bar{\mathbf{u}}_2|$ and a stronger result is given (re-phrased in the terminology of this paper):

Lemma 11 ([3]) *Consider a double square $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with a non-empty \mathbf{u}_2 and let $p = \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ and $s = \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$. Then any inversion factor in \mathbf{v}^2 is either $R_i(\text{IF})$ or $R_{-j}(\text{IF})$ for some $i \in 0, \dots, p$ or some $j \in 0, \dots, s$. Moreover, every $R_i(\text{IF})$ or $R_{-j}(\text{IF})$ appear exactly twice in \mathbf{v}^2 exactly a distance $|\mathbf{v}|$ apart for every $i \in 0, \dots, p$ and every $j \in 0, \dots, s$.*

4 Possible application to New Periodicity Lemma

Some years ago a New Periodicity Lemma was published [4], showing that the occurrence of two special squares at a position i in a string, necessarily precludes the occurrence of other squares of specific period in a specific neighbourhood of i . The proof of this lemma was complex, breaking down into 14 subcases, and required a very strong condition that the shorter of the two squares be regular.

Lemma 12 ([4], New Periodicity Lemma) *Let $\mathbf{x} = \text{DS}(\mathbf{u}, \mathbf{v})$, where we require that \mathbf{u}^2 be regular and that \mathbf{v} be primitive. Then for all integers k and w such that $0 \leq k < |\mathbf{v}| - |\mathbf{u}|$ and $|\mathbf{v}| - |\mathbf{u}| < w < |\mathbf{v}|$, $w \neq |\mathbf{u}|$, $\mathbf{x}[k+1..k+2w]$ is not a square.*

First note that by Observation 8, the requirement that v be primitive is redundant; the fact that u^2 is regular necessarily forces the primitiveness of v . Also note that the regularity of u^2 necessarily implies that in the canonical factorization of $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$, $e_1 = e_2 = 1$.

Consider $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, 1, 1)$. Let $\bar{\mathbf{u}}_2$ be a suffix of \mathbf{u}_1 such that $\mathbf{u}_1 = \mathbf{u}_2\bar{\mathbf{u}}_2$. The canonical factorization thus has the form

$$(\mathbf{u}_2\bar{\mathbf{u}}_2)\mathbf{u}_2(\mathbf{u}_2\bar{\mathbf{u}}_2)(\mathbf{u}_2\bar{\mathbf{u}}_2)\mathbf{u}_2(\mathbf{u}_2\bar{\mathbf{u}}_2).$$

Let us consider a square \mathbf{w}^2 such that $|\mathbf{u}_1| < |\mathbf{w}| < |\mathbf{v}|$ and $|\mathbf{w}| \neq |\mathbf{u}|$. We want to show that this is not possible.

If for instance \mathbf{w} starts in the first \mathbf{u}_2 and ends in the fourth \mathbf{u}_2 , then \mathbf{w} contains fully the IF, so the second \mathbf{w} has to as well, and so $|\mathbf{w}| \geq |\mathbf{v}|$, a contradiction.

If \mathbf{w} ends in the second $\bar{\mathbf{u}}_2$ we cannot argue using IF, but still knowing that $\mathbf{u}_2\bar{\mathbf{u}}_2$ is

primitive and also all its rotations are primitive, using the Synchronization principle can be applied to obtain a contradiction.

Almost all possible cases for w^2 except two can be easily shown impossible using only the properties of the canonical factorization. Thus, we believe, and it is our immediate goal for future research, that the canonical factorization will not only provide us with a significantly simplified proof of New Periodicity Lemma, but will also allow us to significantly reduce the conditions on u^2 from u being regular to just being primitive. We also believe that the canonical factorization in the same way will not only provide a simpler proof of Crochemore-Rytter Three Squares Lemma, but will extend the applicability of the lemma to three squares when any of the squares is primitive (the original lemma requires that the smallest square be primitive).

5 Conclusion and future work

We presented a unique factorization of a double square, i.e. a configuration of two squares u^2 and v^2 starting at the same position and satisfying $|u| < |v| < 2|u|$. We call this factorization the *canonical factorization*. It has very strong combinatorial properties as it is an almost periodic repetition of a primitive string. We indicated that we would like to use this new insight into the structure of double squares in improving the New Periodicity Lemma [4] and Crochemore-Rytter's Three Squares Lemma [2] and simplifying their proofs. As of preparing this final version of the Prague Stringology Conference 2014 proceedings, we are happy to report that the canonical factorization presented here indeed greatly simplified and generalized both. The follow-up work will focus on presenting of these results in a near future.

References

1. W. BLAND AND W. F. SMYTH: *Overlapping squares: the general case characterized & applications*. submitted for publication, 2014.
2. M. CROCHEMORE AND W. RYTTER: *Squares, cubes, and time-space efficient string searching*. *Algorithmica*, 13 1995, pp. 405–425.
3. A. DEZA, F. FRANEK, AND A. THIERRY: *How many double squares can a string contain?* submitted for publication, 2013.
4. K. FAN, S. PUGLISI, W. F. SMYTH, AND A. TURPIN: *A new periodicity lemma*. *SIAM Journal on Discrete Mathematics*, 20 2006, pp. 656–668.
5. F. FRANEK, R. C. G. FULLER, J. SIMPSON, AND W. F. SMYTH: *More results on overlapping squares*. *Journal of Discrete Algorithms*, 17 2012, pp. 2–8.
6. E. KOPYLOVA AND W. F. SMYTH: *The three squares lemma revisited*. *Journal of Discrete Algorithms*, 11 2012, pp. 3–14.
7. N. H. LAM: *On the number of squares in a string*. *AdvOL-Report 2013/2*, McMaster University, 2013.
8. J. SIMPSON: *Intersecting periodic words*. *Theoretical Computer Science*, 374 2007, pp. 58–65.
9. B. SMYTH: *Computing Patterns in Strings*, Pearson Addison-Wesley, 2003.

Proceedings of the Prague Stringology Conference 2014

Edited by Jan Holub and Jan Žďárek

Published by: Prague Stringology Club

Department of Theoretical Computer Science
Faculty of Information Technology
Czech Technical University in Prague
Thákurova 9, Praha 6, 160 00, Czech Republic.

ISBN 978-80-01-05547-2

URL: <http://www.stringology.org/>

E-mail: psc@stringology.org Phone: +420-2-2435-9811

Printed by Česká technika – Nakladatelství ČVUT
Thákurova 550/1, Praha 6, 160 41, Czech Republic

© Czech Technical University in Prague, Czech Republic, 2014