



# The New Periodicity Lemma revisited

Haoyue Bai<sup>a</sup>, Frantisek Franek<sup>a,\*</sup>, William F. Smyth<sup>a,b,c</sup>

<sup>a</sup> Department of Computing and Software, McMaster University, Hamilton, Ontario, Canada

<sup>b</sup> School of Engineering & Information Technology, Murdoch University, Perth, Western Australia, Australia

<sup>c</sup> School of Computer Science & Software Engineering, University of Western Australia, Perth, Western Australia, Australia

## ARTICLE INFO

### Article history:

Received 30 September 2014

Received in revised form 20 January 2016

Accepted 2 May 2016

Available online 1 June 2016

### Keywords:

String

Square

Canonical factorization

Double square

New Periodicity Lemma

## ABSTRACT

In 2006, the *New Periodicity Lemma* (NPL) was published, showing that the occurrence of two squares starting at a position  $i$  in a string necessarily precludes the occurrence of other squares of specified period in a specified neighbourhood of  $i$ . The proof of this lemma was complex, breaking down into 14 subcases, and requiring that the shorter of the two squares be *regular*. In this paper we significantly relax the conditions required by the NPL and removing the need for regularity altogether, and we establish a more precise result using a simpler proof based on lemmas that expose new combinatorial structures in a string, in particular a *canonical factorization* for any two squares that start at the same position.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

In 1995 Crochemore and Rytter [3] considered three distinct squares, all prefixes of a given string  $x$ , 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. [5] considered two squares that were prefixes of  $x$  with the third square offset some distance to the right; they proved a *New Periodicity Lemma* (NPL), describing conditions under which the third square could not exist. Since that time there has been considerable work done [2,7,8,10] in an effort to specify more precisely the combinatorial structure of the string in the neighbourhood of such squares.

In this paper we first discuss a *canonical factorization*, a unique breakdown into primitive strings of what we call a *double square*, i.e. a pair of two squares starting at the same position and being of “comparable” lengths. A weaker form of the factorization was instrumental in the improved upper bound for the number of distinct squares [4]; weaker in the sense that it only applied to FS-double squares, i.e. two squares that start at the same position and both are rightmost occurrences. Note that our notion of double squares is weaker and hence every FS-double square is a double square, but not the other way around. The canonical factorization indicates that double squares indeed have an intricate highly periodic intrinsic structure. This structure has two factors that are unique in their occurrences within the double square. They were introduced in [4] and referred to as *inversion factors* due to their structure. In [12], Thierry discusses the *core of the period interrupt*, a very similar concept to the one we introduce here as RIS (Right Inversion Subfactor) and LIS (Left Inversion Subfactor). RIS has only two occurrences in the double square, and so does LIS. The usage of RIS, respective LIS, is straightforward as it significantly limits the size and the placement of a possible third square: let  $u^2$  be a prefix of  $v^2$  and consider a third square  $w^2$ ; if it contains RIS in the first  $w$ , it must contain RIS in the second  $w$  and vice-versa, and hence  $w$  has the same size as  $v$ . So the only other

\* Corresponding author.

E-mail addresses: [baih3@mcmaster.ca](mailto:baih3@mcmaster.ca) (H. Bai), [franek@mcmaster.ca](mailto:franek@mcmaster.ca) (F. Franek), [bill@arg.cas.mcmaster.ca](mailto:bill@arg.cas.mcmaster.ca) (W.F. Smyth).

possibilities are that either  $w^2$  is “too small” that it does not contain RIS, or “too big” that it contains both RIS in the first  $w$ . The restrictions imposed by RIS or LIS allow us to prove a new version of the NPL that is much more general in its application while at the same time being more precise in its result.

The paper is structured as follows: in Section 2 we discuss the basic facts and notations. In Section 3 we present and prove the *Two Squares Factorization Lemma* giving what we refer to as the canonical factorization of a double square. In Section 4 we discuss the inversion factors and their refinements RIS and LIS. The new formulation of the NPL is then presented and proved in Section 5. Finally, Section 6 presents a brief conclusion of the research described.

## 2. Preliminaries

In this section we introduce the basic notation and develop the combinatorial tools that will be used to determine a canonical factorization for a double square. Chief among these are the Synchronization Principle (Lemma 2) and the Common Factor Lemma (Lemma 3), that lead to the Two Squares Factorization Lemma (Lemma 5).

A **string**  $x$  is a finite sequence of symbols, called **letters**, drawn from a (finite or infinite) set  $\Sigma$ , called the **alphabet**. The length of the sequence is called the **length** of  $x$ , denoted  $|x|$ . Sometimes for convenience we represent a string  $x$  of length  $n$  as an array  $x[1..n]$ . The string of length zero is called the **empty string**, denoted  $\varepsilon$ . If a string  $x = uvw$ , where  $u, v, w$  are strings, then  $u$  (respectively,  $v, w$ ) is said to be a **prefix** (respectively, **substring**, **suffix**) of  $x$ ; a **proper prefix** (respectively, **proper substring**, **proper suffix**) if  $|u| < |x|$  (respectively,  $|v| < |x|$ ,  $|w| < |x|$ ). An empty prefix or suffix is called **trivial**. A substring is also called a **factor**. Given strings  $u$  and  $v$ ,  $\text{lcp}(u, v)$  (respectively,  $\text{lcs}(u, v)$ ) is the **longest common prefix** (respectively, **longest common suffix**) of  $u$  and  $v$ .

If  $x$  is a concatenation of  $k \geq 2$  copies of a nonempty string  $u$ , we write  $x = u^k$  and say that  $x$  is a **repetition**; if  $k = 2$ , we say that  $x = u^2$  is a **square**; if there exist no such integer  $k$  and no such  $u$ , we say that  $x$  is **primitive**. If  $x = u^k$ ,  $k \geq 1$ , and  $u$  is primitive, we call  $u$  the **primitive root** of  $x$ . If  $x = v^2$  has a proper prefix  $u^2$ ,  $|u| < |v| < 2|u|$ , we say that  $x$  is a **double square** and write  $x = \text{DS}(u, v)$ . A square  $u^2$  such that  $u$  has no square prefix is said to be **regular**.

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

Strings  $uv$  and  $vu$  are **conjugates**, written  $uv \sim vu$ . We also say that  $vu$  is the  $|u|$ th **rotation** of  $x = uv$ , written  $R_{|u|}(x)$ , or the  $-|v|$ th **rotation** of  $x$ , written  $R_{-|v|}(x)$ , while  $R_0(x) = x$ . 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 whenever  $x[i+k..j+k]$  is a cyclic shift of  $x[i..j]$ , the two substrings must therefore be conjugates; however, for  $k > j-i+1$ , the converse does not hold (for example, in  $x = abacba$ ,  $x[5..7] = baa$  is a conjugate, but not a cyclic shift, of  $x[1..3] = aba$ ).

**Lemma 1** ([11, Lemma 1.4.2]). Let  $x$  be a string of length  $n$  and minimum period  $\pi \leq n$ , and let  $j \in 1..n-1$  be an integer. Then  $R_j(x) = x$  if and only if  $x$  is not primitive and  $\pi$  is divisible by  $j$ .

The following results (Lemmas 2–3) were first stated in [4] without proof as they all follow from the Periodicity lemma of Fine and Wilf, [6]. Although proofs were later provided in [1], we repeat them here for completeness.

**Lemma 2** (Synchronization Principle). The primitive string  $x$  occurs exactly  $p$  times in  $x_2x^p x_1$ , where  $p$  is a nonnegative integer and  $x_1$  (respectively,  $x_2$ ) is a proper prefix (respectively, proper suffix) of  $x$ .

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

**Lemma 3** (Common Factor Lemma). Suppose that  $x$  and  $y$  are primitive strings, where  $x_1$  (respectively,  $y_1$ ) is a proper prefix and  $x_2$  (respectively,  $y_2$ ) a proper suffix of  $x$  (respectively,  $y$ ). If for integers  $p \geq 2$  and  $q \geq 2$ ,  $x_2x^p x_1$  and  $y_2y^q y_1$  have a common factor of length  $|x| + |y|$ , then  $x \sim y$ .

**Proof.** First consider the special case  $x_1 = x_2 = y_1 = y_2 = \varepsilon$ , where  $x^p, y^q$  have a common prefix  $f$  of length  $|x| + |y|$ . We show that in this case  $x = y$ .

Observe that  $f$  has prefixes  $x$  and  $y$ , so that if  $|x| = |y|$ , then  $x = y$ , as required. Therefore suppose WLOG that  $|x| < |y|$ . Note that  $y \neq x^k$  for any integer  $k \geq 2$ , since otherwise  $y$  would not be primitive, contradicting the hypothesis of the lemma. Hence there exists  $k \geq 1$  such that  $k|x| < |y|$  and  $(k+1)|x| > |y|$ . But since  $f = yx$ , it follows that

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

again by Lemma 1 contrary to the assumption that  $x$  is primitive. We conclude that  $|x| \neq |y|$ , hence that  $|x| = |y|$  and  $x = y$ , as required.

Now consider the general case, where  $f$  of length  $|x| + |y|$  is a common factor of  $x_2x^p x_1$  and  $y_2y^q y_1$ . Then  $x_2x^p x_1 = ufu'$  for some  $u$  and  $u'$ . If  $|u| \geq |x|$ , then  $f$  is a factor of  $x_1x^{p-1}x_2$ , and so we can assume WLOG that  $|u| < |x|$ . Setting  $\tilde{x} = R_{|u|}(x)$ , we see that  $f$  is a prefix of  $\tilde{x}^p$ .

Download English Version:

<https://daneshyari.com/en/article/4950000>

Download Persian Version:

<https://daneshyari.com/article/4950000>

[Daneshyari.com](https://daneshyari.com)