



RESEARCH REPOSITORY

*This is the author's final version of the work, as accepted for publication following peer review but without the publisher's layout or pagination.
The definitive version is available at:*

<http://dx.doi.org/10.1016/j.dam.2016.05.003>

Bai, H., Franěk, F. and Smyth, W.F. (2016) The new periodicity lemma revisited. *Discrete Applied Mathematics*, 212 . pp. 30-36.

<http://researchrepository.murdoch.edu.au/id/eprint/31923/>

Copyright: © 2016 Elsevier B.V.
It is posted here for your personal use. No further distribution is permitted.

The New Periodicity Lemma Revisited

Haoyue Bai¹, Frantisek Franek¹, William F. Smyth^{1,2,3}

¹ *Department of Computing and Software*

McMaster University, Hamilton, Ontario, Canada

² *School of Engineering & Information Technology,*

Murdoch University, Perth, Western Australia

³ *School of Computer Science & Software Engineering*

University of Western Australia, Perth, Western Australia

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.

Keywords: string, square, canonical factorization, double square, New Periodicity Lemma

1. Introduction

In 1995 Crochemore and Rytter [3] 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.* [5] considered two squares that were prefixes of \mathbf{x} with the third square offset some distance

Email addresses: baih3@mcmaster.ca (Haoyue Bai¹), franek@mcmaster.ca (Frantisek Franek¹), smyth5@mcmaster.ca (William F. Smyth^{1,2,3})

Preprint submitted to Journal of Discrete Applied Mathematics

January 20, 2016

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

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 \mathbf{u}^2 be a prefix of \mathbf{v}^2 and consider a third square \mathbf{w}^2 ; if it contains RIS in the first \mathbf{w} , it must contain RIS in the second \mathbf{w} and vice-versa, and hence \mathbf{w} has the same size as \mathbf{v} . So the only other possibilities are that either \mathbf{w}^2 is “too small” that it does not contain RIS, or “too big” that it contains both RIS in the first \mathbf{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.

1
2
3
4
5
6
7
8
9 **2. Preliminaries**

10
11 In this section we introduce the basic notation and develop the combina-
12 torial tools that will be used to determine a canonical factorization for a dou-
13 ble square. Chief among these are the Synchronization Principle (Lemma 2)
14 and the Common Factor Lemma (Lemma 3), that lead to the Two Squares
15 Factorization Lemma (Lemma 5).
16

17 A *string* x is a finite sequence of symbols, called *letters*, drawn from a
18 (finite or infinite) set Σ , called the *alphabet*. The length of the sequence is
19 called the *length* of x , denoted $|x|$. Sometimes for convenience we represent
20 a string x of length n as an array $x[1..n]$. The string of length zero is called
21 the *empty string*, denoted ε . If a string $x = uvw$, where u, v, w are strings,
22 then u (respectively, v, w) is said to be a *prefix* (respectively, *substring*,
23 *suffix*) of x ; a *proper prefix* (respectively, *proper substring*, *proper*
24 *suffix*) if $|u| < |x|$ (respectively, $|v| < |x|$, $|w| < |x|$). An empty prefix or
25 suffix is called *trivial*. A substring is also called a *factor*. Given strings
26 u and v , $\text{lcp}(u, v)$ (respectively, $\text{lcs}(u, v)$) is the *longest common prefix*
27 (respectively, *longest common suffix*) of u and v .
28

29 If x is a concatenation of $k \geq 2$ copies of a nonempty string u , we write
30 $x = u^k$ and say that x is a *repetition*; if $k = 2$, we say that $x = u^2$ is
31 a *square*; if there exist no such integer k and no such u , we say that x is
32 *primitive*. If $x = u^k$, $k \geq 1$, and u is primitive, we call u the *primitive*
33 *root* of x . If $x = v^2$ has a proper prefix u^2 , $|u| < |v| < 2|u|$, we say that x
34 is a *double square* and write $x = \text{DS}(u, v)$. A square u^2 such that u has
35 no square prefix is said to be *regular*.
36

37 For $x = x[1..n]$, $1 \leq i < j \leq j+k \leq n$, the string $x[i+1..j+1]$ is a *right*
38 *cyclic shift* of $x[i..j]$ by 1 position if $x[i] = x[j+1]$; the string $x[i+k..j+k]$
39 is a right cyclic shift of $x[i..j]$ by k positions if $x[i+k-1..j+k-1]$ is a right
40 cyclic shift of $x[i..j]$ by $k-1$ positions and $x[i+k..j+k]$ is a right cyclic shift
41 of $x[i+k-1..j+k-1]$ of 1 position. Equivalently, we can say that $x[i..j]$ is a
42 *left cyclic shift* by k positions of $x[i+k..j+k]$. When it is clear from the
43 context, we may leave out the number of positions and just speak of a left
44 or right cyclic shift.
45

46 Strings uv and vu are *conjugates*, written $uv \sim vu$. We also say
47 that vu is the $|u|^{\text{th}}$ *rotation* of $x = uv$, written $R_{|u|}(x)$, or the $-|v|^{\text{th}}$
48 *rotation* of x , written $R_{-|v|}(x)$, while $R_0(x) = x$. As for the cyclic shift,
49 when it is clear from the context we may leave out the number of rotations
50 and just speak of a rotation. Note that whenever $x[i+k..j+k]$ is a cyclic
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

shift of $\mathbf{x}[i..j]$, the two substrings must therefore be conjugates; however, for $k > j - i + 1$, the converse does not hold (for example, in $\mathbf{x} = abacbaa$, $\mathbf{x}[5..7] = baa$ is a conjugate, but not a cyclic shift, of $\mathbf{x}[1..3] = aba$).

Lemma 1. [11, Lemma 1.4.2] *Let \mathbf{x} be a string of length n and minimum period $\pi \leq n$, and let $j \in 1..n-1$ be an integer. Then $R_j(\mathbf{x}) = \mathbf{x}$ if and only if \mathbf{x} is not primitive and π 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 \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 cyclic shift $R_j(\mathbf{x})$ 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 integers $p \geq 2$ and $q \geq 2$, $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$ 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 $\tilde{\mathbf{x}}^p$.

1
2
3
4
5
6
7
8
9 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
10 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
11 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
12 follows. \square
13
14

15 The Common Factor Lemma gives rise to the following corollary useful
16 for showing the uniqueness of the canonical factorization of a double square
17 presented in Section 3:
18

19 **Lemma 4** ([4]). *Suppose that \mathbf{x} and \mathbf{y} are primitive strings, and that p and*
20 *q are positive integers.*
21

- 22
23 (a) *If $\mathbf{x}^p = \mathbf{y}^q$, then $\mathbf{x} = \mathbf{y}$ and $p = q$.*
24
25 (b) *If \mathbf{x}_1 (respectively, \mathbf{y}_1) is a proper prefix of \mathbf{x} (respectively, \mathbf{y}) and*
26 *$\mathbf{x}^p\mathbf{x}_1 = \mathbf{y}^q\mathbf{y}_1$ for $p \geq 2$, $q \geq 2$, then $\mathbf{x} = \mathbf{y}$, $\mathbf{x}_1 = \mathbf{y}_1$ and $p = q$.*
27
28

29 *Proof.* For (a), first consider $p = 1$, thus $\mathbf{x} = \mathbf{y}^q$. Since \mathbf{x} is primitive,
30 therefore $q = 1$ and $\mathbf{x} = \mathbf{y}$, as required. Similarly for $q = 1$. Suppose then
31 that $p, q \geq 2$. This means that \mathbf{x}^p and $\mathbf{y}^q = \mathbf{x}^p$ have a common factor of
32 length $p|\mathbf{x}| = q|\mathbf{y}| \geq |\mathbf{x}| + |\mathbf{y}|$, so that by Lemma 3 $\mathbf{x} \sim \mathbf{y}$. Hence $|\mathbf{x}| = |\mathbf{y}|$
33 and so $\mathbf{x} = \mathbf{y}$.
34

35 For (b), since again $p \geq 2$, $q \geq 2$, it follows as in (a) that $\mathbf{x}^p\mathbf{x}_1 = \mathbf{y}^q\mathbf{y}_1$
36 has a common factor of length at least $|\mathbf{x}| + |\mathbf{y}|$, hence the result. \square
37

38 Note that in Lemma 4(b) the requirement $p \geq 2$, $q \geq 2$ is essential. For
39 instance, $\mathbf{x} = aabb$, $\mathbf{x}_1 = aa$ and $p = 2$ yields $\mathbf{x}^p\mathbf{x}_1 = aabbaabbaa$, identical
40 to $\mathbf{y}^q\mathbf{y}_1$ produced by $\mathbf{y} = aabbaabba$, $\mathbf{y}_1 = a$ and $q = 1$ — but of course
41 $\mathbf{x} \neq \mathbf{y}$.
42
43
44

45 3. Canonical Factorization of Double Squares

46
47 The most general unique factorization of any double square $\text{DS}(\mathbf{u}, \mathbf{v})$ dis-
48 cussed in this section was presented in [1]. For the sake of completeness, we
49 not only quote the results from [1], but include the proofs as well.
50

51 The uniqueness of the factorization allows us to speak of the *canonical*
52 *factorization of* $\text{DS}(\mathbf{u}, \mathbf{v})$. This structure has been described before [4, 5, 8,
53 7, 9], but not as precisely and with more assumptions required or in a weaker
54 form; above all, Lemma 5 establishes the uniqueness of the breakdown with
55 no additional assumptions.
56
57
58

Lemma 5 ([1], Two Squares Factorization Lemma). *If $\mathbf{x} = \text{DS}(\mathbf{u}, \mathbf{v})$, 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 $e_1 > e_2$ (thus $e_1 \geq 2$);*

(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. In this proof, for a tandem repeat $\mathbf{w}\mathbf{w}$, we use $\mathbf{w}_{[1]}$ to refer to the first \mathbf{w} , while $\mathbf{w}_{[2]}$ to refer to the second \mathbf{w} .

Let \mathbf{z} be the nonempty proper prefix of $\mathbf{u}_{[2]}$ that is in addition a suffix 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 prefix \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 $e_1 \geq e_2k$ 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:

$|\mathbf{u}_2| = 0$: 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 4(a).

$|\mathbf{u}_2| > 0$: 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 4(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

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

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 Lemma 3 with variables

$$(\mathbf{x}, \mathbf{y}, p, q) \equiv (\mathbf{w}, \mathbf{u}_1, 2k, e_1 + e_2)$$

to conclude that $\mathbf{w} \sim \mathbf{u}_1$, thus by (1) that $\mathbf{w} = \mathbf{u}_1$. Let $\bar{\mathbf{u}}_2$ be a suffix of \mathbf{u}_1 so that $\mathbf{u}_1 = \mathbf{u}_2 \bar{\mathbf{u}}_2$. By the Synchronization Principle (Lemma 2), (1) implies that $\bar{\mathbf{u}}_2$ 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 5, \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 Lemma 3 with variables $(\mathbf{x}, \mathbf{y}, p, q) \equiv (\mathbf{w}, \mathbf{u}_1, 2k, e_1 + e_2)$, 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 5(b) requires that $e_1 \geq 2$. To see that this condition is not necessary, consider $\mathbf{x} = \mathbf{abaababab}$, where $\mathbf{u} = (\mathbf{ab})\mathbf{a}$, $\mathbf{v} = (\mathbf{ab})\mathbf{a}(\mathbf{ab})$, so that $\mathbf{u}_1 = \mathbf{ab}$, $\mathbf{u}_2 = \mathbf{a}$, $e_1 = e_2 = 1$, but \mathbf{u} is primitive.
- (b) On the other hand, consider $\mathbf{x} = \mathbf{abaabaabaabaabaabaab}$, where $\mathbf{u} = (\mathbf{aba})^2 = (\mathbf{abaab})\mathbf{a}$, $\mathbf{v} = (\mathbf{abaab})\mathbf{a}(\mathbf{abaab})$, so that $\mathbf{u}_1 = \mathbf{abaab}$, $\mathbf{u}_2 = \mathbf{a}$, $e_1 = e_2 = 1$, where now \mathbf{u} is *not* primitive.

Lemma 5 gives credence to the following definition of terminology and notation:

Definition 6. 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 5 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 5 gives rise to a number of important observations:

Observation 7. In Lemma 5, $|\mathbf{u}_2| > 0$ 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 ;
- (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 $e_1 = e_2 = 1$.

Proof. (a) $|\mathbf{u}_2| = 0$ implies that $\mathbf{v} = \mathbf{u}_1^{e_1+e_2}$ and since $e_1+e_2 \geq 2$, it follows that \mathbf{v} is not primitive.

(b) We show that $|\mathbf{u}_2| = 0$ implies \mathbf{u} not primitive. First note that $u = u_1^{e_1}$. If $e_1 \geq 2$, then it follows directly that u is not primitive. If $e_1 = 1$, then $e_2 = 1$ and so $v = u_1^2$ and $u = u_1$ which is a contradiction as $|v| < 2|u|$.

(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 5, $|\mathbf{u}_2| > 0$ 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 $e_1 = e_2 = 1$. □

In the context of Observation 7(f), consider the double square

$$\text{DS}(\mathbf{u}, \mathbf{v}) = aabaaaabaabaaaab,$$

with $\mathbf{u} = aabaa$, $\mathbf{v} = aabaaaab$. 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.

4. Rare Factors in Double Squares

In this section we consider a double square $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with nonempty \mathbf{u}_2 . We present three kinds of factors of \mathbf{v}^2 that have highly restricted number of occurrences in \mathbf{v}^2 . The first of these, the so-called *inversion factor* (IF for short) was introduced in [4]. The additional two factors introduced here, *right inversion subfactor* (RIS for short) and *left inversion subfactor* (LIS for short), are in fact subfactors of IF. Note that A. Thierry [12] examines rare factors in configurations $\mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_2}$.

Using the canonical factorization of $\text{DS}(\mathbf{u}, \mathbf{v})$, we have

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

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

Lemma 8. *Consider a double square $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with nonempty \mathbf{u}_2 . Then the inversion factor IF has exactly two occurrences in \mathbf{v}^2 that are distance $|\mathbf{v}|$ apart, as shown in (2).*

Proof. If IF occurs elsewhere in \mathbf{v}^2 , then by the Synchronization Principle (Lemma 2) its primitive subfactor $\mathbf{u}_2\bar{\mathbf{u}}_2$ can only align with another occurrence of $\mathbf{u}_2\bar{\mathbf{u}}_2$. Therefore its other subfactor $\bar{\mathbf{u}}_2\mathbf{u}_2$ must align with $\mathbf{u}_2\bar{\mathbf{u}}_2$, thus by Lemma 1 contradicting the primitiveness of $\mathbf{u}_2\bar{\mathbf{u}}_2$. □

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

The quantity $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ gives the maximum 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 maximum number of positions $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1}$ and $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}$ can be cyclically shifted to the right. In [4], 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 proved:

Lemma 9 ([4]). *Consider $\mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2\mathbf{u}_1^{e_2}$, where \mathbf{u}_1 is primitive and \mathbf{u}_2 is a nonempty 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 [4] 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 $|\bar{\mathbf{w}}| = |\bar{\mathbf{u}}_2|$; then a stronger result is given (re-phrased in the terminology of this paper):

Lemma 10 ([4]). *Consider a double square $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with nonempty \mathbf{u}_2 , and let $p = \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$, $s = \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$. Then any inversion factor in \mathbf{v}^2 is either $R_i(\text{IF})$ or $R_{-j}(\text{IF})$ for some $i \in 0..p$ or some $j \in 0..s$. Moreover, for every $i \in 0..p$ (respectively, $j \in 0..s$), every $R_i(\text{IF})$ (respectively, $R_{-j}(\text{IF})$) occurs exactly twice in \mathbf{v}^2 with occurrences distance exactly $|\mathbf{v}|$ apart.*

The following simple lemma will be used for determining a different type of rare factor in a double square. It says that if a substring \mathbf{u} of a string \mathbf{x} and its rotation \mathbf{u}' completely overlap except for one symbol, then \mathbf{u} can be cyclically shifted one position to the right, or, equivalently, \mathbf{u}' can be cyclically shifted one position to the left.

Lemma 11. *If the substrings $\mathbf{x}[1..n]$ and $\mathbf{x}[2..n+1]$ of $\mathbf{x} = \mathbf{x}[1..n+1]$ are conjugates, then $\mathbf{x}[1] = \mathbf{x}[n+1]$.*

Proof. (Due to A. Thierry) Since $\mathbf{x}[1..n]$ and $\mathbf{x}[2..n+1]$ are conjugates, the frequency of the alphabet symbols in both must be the same. Let $\mathbf{x}[1] = a$. Then $\mathbf{x}[1..n]$ must have the same number of a 's as $\mathbf{x}[2..n+1]$, and so $\mathbf{x}[n+1] = a$. \square

Note that Lemma 11 does not hold if \mathbf{u} and \mathbf{u}' overlap less. For

$$\underbrace{abbbb}_{\mathbf{u}}ba\dots \text{ and } ab\underbrace{bbba}_{\mathbf{u}'}\dots,$$

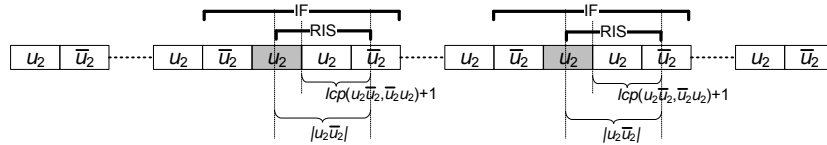
\mathbf{u} cannot be cyclically shifted to the right nor \mathbf{u}' to the left, yet \mathbf{u}' is a rotation of \mathbf{u} .

Definition 12. Consider a double square $\mathbf{x} = \text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with nonempty \mathbf{u}_2 . The **right inversion subfactor** (or RIS) is defined to be a factor $\mathbf{x}[i..j]$ of length $|\mathbf{u}_1|$ where $i = (e_1 - 1)|\mathbf{u}_1| + |\mathbf{u}_2| + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ and $j = e_1|\mathbf{u}_1| + |\mathbf{u}_2| + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + 1$. Similarly, the **left inversion subfactor** (or LIS) is a factor $\mathbf{x}[i..j]$ of length $|\mathbf{u}_1|$ where $i = e_1|\mathbf{u}_1| + |\mathbf{u}_2| - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ and $j = (e_1 + 1)|\mathbf{u}_1| + |\mathbf{u}_2| - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) - 1$.

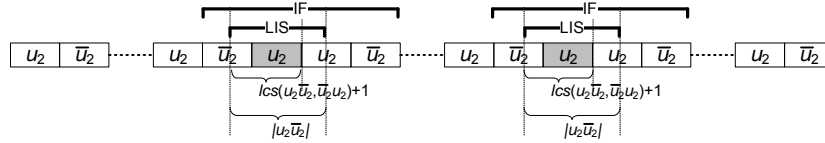
Note that another RIS also naturally occurs at position $(e_1 - 1)|\mathbf{u}_1| + |\mathbf{u}_2| + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + |\mathbf{v}|$, and another LIS at position $e_1|\mathbf{u}_1| + |\mathbf{u}_2| - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + |\mathbf{v}|$.

It is possible to view RIS in this way: let $\mathbf{x}[i..j]$ be the maximum right cyclic shift of the rightmost $\mathbf{u}_2\bar{\mathbf{u}}_2$ of $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1}$, respectively of $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}$. Then RIS is $\mathbf{x}[i+1, j+1]$. Similarly, let $\mathbf{x}[i, j]$ be the maximum left cyclic shift of the first $\mathbf{u}_2\bar{\mathbf{u}}_2$ of $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}$, respectively of $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2}$. Then LIS is $\mathbf{x}[i-1, j-1]$.

For a better understanding, we illustrate the two natural occurrences of RIS in the following diagram:



and the two natural occurrences of LIS as follows:



Lemma 13. Consider a double square $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ with nonempty \mathbf{u}_2 . Let \mathbf{p} be the longest common prefix and \mathbf{s} the longest common suffix of $\mathbf{u}_2\bar{\mathbf{u}}_2$ and $\bar{\mathbf{u}}_2\mathbf{u}_2$. Then there is no occurrence of RIS, respectively LIS, that is a factor of $\mathbf{s}(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}\mathbf{u}_2\mathbf{p}$.

Proof. Let us assume that there is an occurrence of RIS that is a factor of $\mathbf{s}(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}\mathbf{u}_2\mathbf{p}$. Therefore RIS must be a rotation of $\mathbf{u}_2\bar{\mathbf{u}}_2$ since its length is $|\mathbf{u}_1|$. Consider \mathbf{w} , the maximum right cyclic shift of the rightmost $\mathbf{u}_2\bar{\mathbf{u}}_2$ of $(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1}$: \mathbf{w} overlaps with RIS except for one symbol and RIS is its rotation,

so we can apply Lemma 11 to conclude that \mathbf{w} can be right cyclically shifted one more position, a contradiction.

The proof for LIS follows the same line of reasoning. \square

5. New Periodicity Lemma Revisited

Let us first state the original New Periodicity Lemma (NPL) in terms used in this paper:

Lemma 14 ([5], 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 7(d)-(e), the requirement that \mathbf{v} be primitive is redundant: the fact that \mathbf{u}^2 is regular forces $|\mathbf{u}_2| > 0$ as well as the primitiveness of \mathbf{v} . Also by Observation 7(f), the regularity of \mathbf{u}^2 implies that in the canonical factorization $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$, $e_1 = e_2 = 1$; in other words, NPL applies only to a small subset of possible double squares. Therefore $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2$, $\mathbf{v} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$, and $|\mathbf{v}| - |\mathbf{u}| = |\mathbf{u}_1|$. Thus in loose terms, NPL forbids a square \mathbf{w}^2 starting in \mathbf{u}_1 with size $|\mathbf{u}_1| < |\mathbf{w}| < |\mathbf{v}|$, with the possible exception of size $|\mathbf{u}|$.

Here we present a theorem that extends the result to all possible double squares; the meaning of the suffix u' of v will be illuminated in the proof of the theorem.

Theorem 15. *Consider a double square $\text{DS}(\mathbf{u}, \mathbf{v})$ and let \mathbf{u}' be a suffix of \mathbf{v} such that $\mathbf{v} = \mathbf{u}\mathbf{u}'$. Let \mathbf{w}^2 be any square that is a factor of \mathbf{v}^2 . Then exactly one of the following mutually exclusive cases holds:*

- (a) $\mathbf{w} = \mathbf{v}$, or
- (b) $|\mathbf{w}| < |\mathbf{u}|$, or
- (c) $|\mathbf{u}| \leq |\mathbf{w}| < |\mathbf{v}|$ and the primitive root of \mathbf{w} is a conjugate of the primitive root of \mathbf{u}' .

Before we embark on the proof of the theorem, let us discuss how it relates to the original NPL. As mentioned above, for a very specific double square, NPL forbids squares starting in the first \mathbf{u}_1 of lengths bigger than $|\mathbf{u}_1|$ but smaller than $|\mathbf{v}|$ with a possible exception of length $|\mathbf{u}|$. Theorem 15 for such a double square forbids squares starting anywhere if their length is bigger

than $|\mathbf{u}|$ and smaller than $|\mathbf{v}|$. So the “forbidding power” of the theorem is slightly less than that of NPL with respect to the sizes of \mathbf{w}^2 ; however it covers a larger range of possible starts for “forbidden” squares (anywhere instead of in the first \mathbf{u}_1), and above all, it applies to all double squares without any additional conditions or constraints — more than a fair trade-off in our opinion.

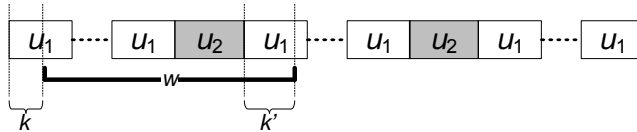
Proof. If $|w| \geq |v|$, then $w = v$ since w^2 is a factor of v^2 , and thus case (a) holds. Hence, for the remainder of the proof we can assume that $|w| < |v|$. Let $\text{DS}(\mathbf{u}, \mathbf{v}) = (\mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$. Then $\mathbf{u}' = \mathbf{u}_1^{e_2}$ and the primitive root of \mathbf{u}' is \mathbf{u}_1 .

We first deal with the case $|\mathbf{u}_2| = 0$. By Lemma 5, \mathbf{u}_1 is the primitive root of $\mathbf{u} = \mathbf{u}_1^{e_1}$ and of $\mathbf{v} = \mathbf{u}_1^{e_1+e_2}$ with $e_1 > e_2 \geq 1$. If case (b) does not hold, we must have $|w| \geq |\mathbf{u}| = |\mathbf{u}_1^{e_1}| > |\mathbf{u}_1|$. Thus \mathbf{w}^2 and $\mathbf{u}_1^{2e_1+2e_2}$ have a common factor of size $|\mathbf{u}_1| + |w|$, so that by Lemma 3, the primitive root of \mathbf{w} is a conjugate of \mathbf{u}_1 , i.e. case (c) holds.

Now let us deal with case $|\mathbf{u}_2| > 0$. Suppose that (b) does not hold and so $|\mathbf{u}| \leq |w| < |v|$. We employ the same notation as in the proof of Lemma 5: thus, for example, $\mathbf{w}_{[1]}$ refers to the first occurrence of \mathbf{w} in \mathbf{w}^k , $\mathbf{w}_{[2]}$ to the second, etc.

Let us assume that there is a square \mathbf{w}^2 starting in $\mathbf{u}_1^{e_1}$ such that $|w| > |\mathbf{u}|$. Since for $|w| = |\mathbf{u}|$, \mathbf{w} can only be a conjugate of \mathbf{u} , and hence the primitive root of \mathbf{w} must be a conjugate of the primitive root of \mathbf{u} , i.e. \mathbf{u}_1 , so that (c) holds, we may suppose $|w| > |\mathbf{u}|$. First note that due to the virtual left-right symmetry of the canonical factorization $\mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2\mathbf{u}_1^{e_2}$ (only the exponents e_1 and e_2 may differ), and to the fact that the arguments presented below can be applied either from the left or from the right, we therefore need only prove the assertion for \mathbf{w}^2 starting in $\mathbf{v}_{[1]}$. Several cases need to be discussed:

- (1) \mathbf{w}^2 starts in the first \mathbf{u}_1 of $\mathbf{u}_1^{e_1}$ and ends in the first \mathbf{u}_1 of $\mathbf{u}_1^{e_1+e_2}$



Since $|w| > |\mathbf{u}| = |\mathbf{u}_1^{e_1}| + |\mathbf{u}_2|$, $k < k'$.

- (i) $k' \leq \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$

Then $\mathbf{w}_{[1]}$ has as a prefix a k -th rotation of \mathbf{u}_1 and $\mathbf{w}_{[2]}$ has as a

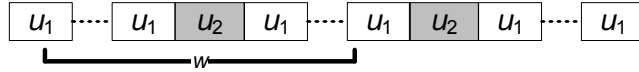
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

prefix a k' -th rotation of \mathbf{u}_1 . By the Synchronization Principle, $k = k'$, a contradiction. This case is not possible.

(ii) $k' > \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$

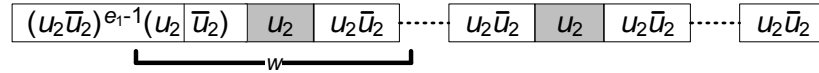
Here $\mathbf{w}_{[1]}$ contains the first RIS, and so $\mathbf{w}_{[2]}$ must contain an occurrence of RIS. Since $\mathbf{w}_{[2]}$ is a factor in $\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2$, therefore by Lemma 13 $\mathbf{w}_{[2]}$ must contain the second RIS and so $|\mathbf{w}| \geq |\mathbf{v}|$, a contradiction. This case is not possible.

(2) \mathbf{w}^2 starts in the first \mathbf{u}_1 of $\mathbf{u}_1^{e_1}$ and ends past the first \mathbf{u}_1 of $\mathbf{u}_1^{e_1+e_2}$.



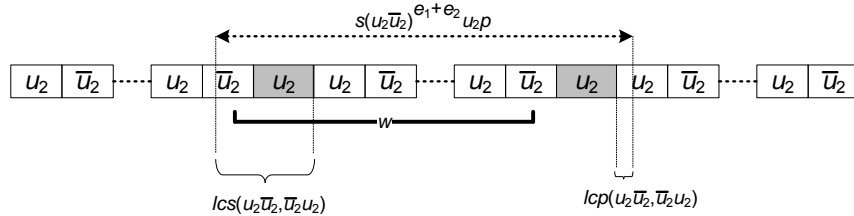
The same argument as in (1)(ii) gives $|\mathbf{w}| \geq |\mathbf{v}|$, a contradiction. This case is not possible.

(3) \mathbf{w}^2 starts in $\mathbf{u}_1^{e_1-1}\mathbf{u}_2$ but not in the first \mathbf{u}_1 .



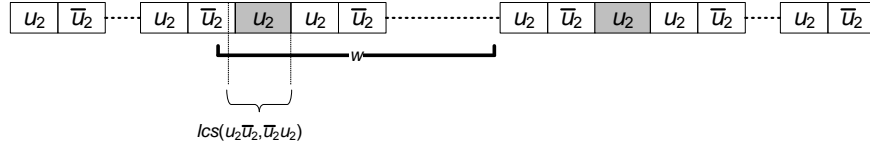
Then $e_1 > 1$ and since $|\mathbf{w}| > |\mathbf{u}| = |\mathbf{u}_1^{e_1}\mathbf{u}_2|$, $\mathbf{w}_{[1]}$ ends past the first \mathbf{u}_1 of $\mathbf{u}_1^{e_1+e_2}$. Therefore, $\mathbf{w}_{[1]}$ contains RIS and so $|\mathbf{w}| \geq |\mathbf{v}|$, i.e. this case is not possible.

(4) \mathbf{w}^2 starts in the suffix $\bar{\mathbf{u}}_2\mathbf{u}_2$ of \mathbf{u} whose length $\leq \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$.



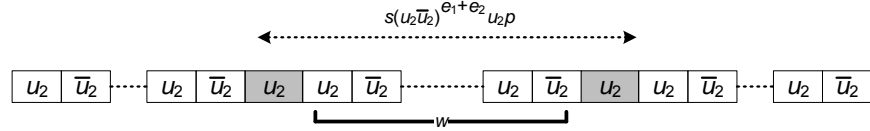
Here $\mathbf{w}_{[1]}$ is a factor in $s\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2p$, where s is the maximal common suffix and p the maximal common prefix of $\mathbf{u}_2\bar{\mathbf{u}}_2$ and $\bar{\mathbf{u}}_2\mathbf{u}_2$. Thus \mathbf{w}^2 and $s\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2p$ have a common factor of size $|\mathbf{u}_1 + \mathbf{w}|$ and by the Common Factor Lemma (Lemma 3), the primitive root of \mathbf{w} is a conjugate of \mathbf{u}_1 , i.e. case (c) holds.

(5) w^2 starts in the suffix $\bar{u}_2 u_2$ of u whose length $> \text{lcs}(u_2 \bar{u}_2, \bar{u}_2 u_2)$.



Then $w_{[1]}$ contains LIS and thus $w_{[2]}$ must contain an occurrence of LIS and so $|w| \geq |v|$, and so this case is not possible.

(6) w^2 starts past the first u .



The same argument as in (4) shows that the primitive root of w is a conjugate of u_1 and so case (c) holds.

□

6. Conclusion

We presented a unique factorization, referred to as the *canonical factorization*, of a double square; that is, a configuration of two squares starting at the same position of “comparable” lengths. Utilizing the canonical factorization we discussed so-called rare factors; that is, factors occurring in a few well-defined positions in a double square. The existence of the rare factors RIS and LIS is then used to establish existential limits for a third square (Theorem 15), greatly generalizing the New Periodicity Lemma. The theorem in comparison to NPL vastly extends the range of applicability from a very specific type of double square $DS(\mathbf{u}, \mathbf{v})$ (\mathbf{u}^2 must be regular) to any type of double square.

References

- [1] H. Bai, F. Franek, and W.F. Smyth. Two squares canonical factorization. In *Proceedings of the Prague Stringology Conference 2014*, Czech Technical University in Prague, Czech Republic, 2014.

- 1
2
3
4
5
6
7
8
9 [2] W. Bland and W.F. Smyth. Three overlapping squares: the general case
10 characterized & applications. *Theoretical Computer Science*, 596-6:23–
11 40, 2015.
12
13 [3] M. Crochemore and W. Rytter. Squares, cubes, and time-space efficient
14 string searching. *Algorithmica*, 13(5):405–425, 1995.
15
16 [4] A. Deza, F. Franek, and A. Thierry. How many double squares can a
17 string contain? *Discrete Applied Mathematics*, 180:52–69, 2015.
18
19 [5] K. Fan, S.J. Puglisi, W.F. Smyth, and A. Turpin. A new periodicity
20 lemma. *SIDMA*, 20-3:656–668, 2006.
21
22 [6] N. J. Fine and H. S. Wilf. Uniqueness theorems for periodic functions.
23 *Proceedings of the American Mathematical Society*, 16(1):pp. 109–114,
24 1965.
25
26 [7] F. Franek, R.C.G. Fuller, J. Simpson, and W.F. Smyth. More results
27 on overlapping squares. *Journal of Discrete Algorithms*, 17:2–8, 2012.
28
29 [8] E. Kopylova and W.F. Smyth. The three squares lemma revisited. *Jour-
30 nal of Discrete Algorithms*, 11:3–14, 2012.
31
32 [9] N.H. Lam. *On the number of squares in a string*. AdvOL-Report,
33 2013/2, McMaster University, 2013.
34
35 [10] J. Simpson. Intersecting periodic words. *Theoretical Computer Science*,
36 374:58–65, 2007.
37
38 [11] B. Smyth. *Computing Patterns in Strings*. ACM Press Bks.
39 Pearson/Addison-Wesley, 2003.
40
41 [12] A. Thierry. Combinatorics of the interrupted period. *Proceedings of the
42 Prague Stringology Conference*, 2015.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58