

# STRUCTURAL FACTORIZATION OF SQUARES IN STRINGS

Haoyue Bai



# STRUCTURAL FACTORIZATION OF SQUARES IN STRINGS

By Haoyue BAI

A Thesis Proposal Submitted to the Department of Computing and  
Software of McMaster University in Partial Fulfilment of the  
Requierelement for the Degree of Doctor of Philosophy



McMaster University © Copyright by Haoyue Bai, May 2017

Ph.D. Thesis  
Department Computing and Software

McMaster University  
Hamilton, Ontario, Canada

TITLE: Structural Factorization of Squares in Strings  
AUTHOR: Haoye Bai  
M.Eng. McMaster University  
Bachelor. Beijing University of Post and Telecommunications  
SUPERVISOR: Dr. Antoine Deza & Dr. Frantisek Franek  
NUMBER OF PAGES: V,73,III

# Abstract

A balanced double square in a string  $\mathbf{x}$  consists of two squares starting in the same position and of comparable lengths. We present a unique factorization of the longer square into primitive components referred to as the canonical factorization and analyze its properties. In particular, we examine the inversion factors and the right and left inversion subfactors. All three substrings are collectively referred to as rare factors as they occur only twice in a significant portion of the larger square. The inversion factors were essential for determining the classification of mutual configurations of double squares and thus providing the best-to-date upper bound of  $11n/6$  for the number of distinct squares in a string of length  $n$  in [8] by Deza, Franek, and Thierry. The right and left inversion subfactors have the advantage of being half the length of the inversion factors, thus providing a stronger discrimination property for a possible third square. This part of the thesis was published in [2] by Bai, Franek, and Smyth.

The canonical factorization and the right and left inversion subfactors are used to formulate and prove a significantly stronger version of the New Peri-

odicity Lemma [9] by Fan, Puglisi, Smyth, and Turpin, 2006, that basically restricts what kind of a third square can exist in a balanced double square. This part of the thesis was published in [3] by Bai, Franek, and Smyth.

The canonical factorization and the inversion factors are applied to formulate and prove a stronger version of the Three Squares Lemma [7] by Crochemore and Rytter. This part of the thesis was published in [1] by Bai, Deza, and Franek.

# Acknowledgments

I would like to express my gratitude to my supervisors, Dr. Antoine Deza and Dr. Frantisek Franek, for their invaluable guidance, generous support and continuous encouragement to my research studies and my life.

My special thanks go to the members of the supervisory and defence committees: Dr. Antoine Deza, Dr. Frantisek Franek, Dr. Kai Huang, Dr. Fred Hoppe, Dr. Dalibor Froncek.

I appreciated the help and moral support from all my colleagues including the members of the Advanced Optimization Laboratory and the members of the school of Computational Science and Engineering.

Furthermore, I am grateful for the financial aid provided by McMaster University.

*To my family and friends.*

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Background of the research . . . . .	5
1.2	Preliminaries . . . . .	8
<b>2</b>	<b>Canonical factorization of double squares</b>	<b>26</b>
2.1	Two Squares Factorization . . . . .	28
2.2	Inversion Factors in a Balanced Double Square . . . . .	37
2.3	Rare Factors in Balanced Double Squares . . . . .	40
<b>3</b>	<b>Application of canonical factorization to New Periodicity</b>	
	Lemma	50
<b>4</b>	<b>Application of canonical factorization to Three Squares</b>	
	Lemma	62
<b>5</b>	<b>Conclusion and Future work</b>	<b>72</b>
	<b>Bibliography</b>	<b>I</b>

## List of Definition

Symbol	Meaning	Page
$\mathbf{x}$	a <i>string</i> $\mathbf{x}$ over an alphabet $\Sigma$	8
$\varepsilon$	empty string	8
$\mathbf{x}[1..n]$	a finite string $\mathbf{x}$ with $n$ symbols	8
$ \mathbf{x} $	the length of string $\mathbf{x}$	8
$\mathbf{xy}$	concatenation of two strings $\mathbf{x}$ and $\mathbf{y}$	9
$\prec$	an <i>order</i> of an alphabet $\Sigma$	10
$\mathbf{u}^k$	concatenation of $k$ copies of $\mathbf{u}$	11
$\mathbf{u}_{[i]}$	denotes the $i$ -th occurrence of $\mathbf{u}$ in $\mathbf{u}^k$ , where $1 \leq i \leq k$	11
$R_k(\mathbf{x})$	is a <i>right <math>k</math>-rotation</i> of a string $\mathbf{x}$	16
$\mathbf{u} \sim \mathbf{v}$	strings $\mathbf{u}$ and $\mathbf{v}$ are <i>conjugates</i>	17
$\text{lcp}(\mathbf{u}, \mathbf{v})$	the length of the <i>longest common prefix</i>	18
$\text{lcs}(\mathbf{u}, \mathbf{v})$	the length of the <i>longest common suffix</i>	18
$DS(\mathbf{u}, \mathbf{v})$	a double square that consist of two squares $\mathbf{u}^2$ and $\mathbf{v}^2$ start at the same position and that $ \mathbf{u}  <  \mathbf{v} $ .	27
$DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$	the <i>canonical factorization</i> of the double square $DS(\mathbf{u}, \mathbf{v})$	28
$\mathbb{IF}$	inversion factor	37
$\mathbb{IF}_1$	the first occurrence of $\mathbb{IF}$ in a double square	37
$\mathbb{IF}_2$	the second occurrence of $\mathbb{IF}$ in a double square	37
$\mathbb{RIS}$	the right inversion subfactor	41
$\mathbb{LIS}$	the left inversion subfactor	41
$\mathbf{R}_1$	the maximal right cyclic shift of the left half of $\mathbb{IF}_1$	45
$\mathbf{R}_2$	the maximal right cyclic shift of the left half of $\mathbb{IF}_2$	45
$\mathbf{L}_1$	the maximal left cyclic shift of the right half of $\mathbb{IF}_1$	45
$\mathbf{L}_2$	the maximal left cyclic shift of the right half of $\mathbb{IF}_2$	45



# Chapter 1

## Introduction

### 1.1 Background of the research

The structural properties of several squares starting in the same proximity have been studied by researchers interested in two major problems concerning periodicities in strings: the *maximum number of runs* and the *maximum number of distinct squares* in a string. In 1995, Crochemore and Rytter [7] described the relationship among the lengths of three distinct squares starting at the same position in a string in the *Three Squares Lemma*:

*Let  $\mathbf{u}^2 \neq \mathbf{v}^2$  be proper prefixes of  $\mathbf{w}^2$  and let  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  be primitive, then  $|\mathbf{u}| + |\mathbf{v}| < |\mathbf{w}|$ .*

In [10], Fraenkel and Simpson observed that in fact a slightly different version was proven as only the primitivness of  $\mathbf{u}$  was actually used in the proof, they also gave a counterexample for the sharp inequality:  $\mathbf{u} = aba$ ,  $\mathbf{v} = abaab$

and  $\mathbf{w} = \text{abaababa}$ . Thus, they used the lemma in this form:

*Let  $\mathbf{u}^2$  and  $\mathbf{v}^2$  be prefixes of  $\mathbf{w}^2$  so that  $\mathbf{u}$  is primitive and  $|\mathbf{u}| < |\mathbf{v}| < |\mathbf{w}|$ .*

*Then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$ .*

and used it for obtaining a bound of  $2|\mathbf{x}|$  for the number of distinct squares in a string  $\mathbf{x}$ . We will discuss the Three Squares Lemma in more details in Chapter 4.

In 2006, Fan *et al.* [9] considered a case of two squares starting at the same position with a third square possibly offset some distance to the right; they presented the *New Periodicity Lemma* describing conditions under which the third square could not exist. Since that time there has been considerable work done [4, 11, 14, 16] in an effort to specify more precisely the combinatorial structure of a string in the neighbourhood of such two squares. We will discuss the New Periodicity Lemma in more details in Chapter 3.

In [8], the canonical factorization of two rightmost squares starting at the same position, which is referred there as an FS-double square, proved essential for improving the Frankel-Simpson upper bound for the maximum number of distinct squares in a string  $\mathbf{x}$  to  $\frac{11}{6}|\mathbf{x}|$ .

This was the motivation of the research described here and it lead to three publications that anchor this thesis: *Two squares canonical factorization* [2], *The New Periodicity Lemma Revisited* [3], and *On a lemma of Crochemore and Rytter* [1].

In the first of the publications, *Two squares canonical factorization* [2], the canonical factorization for all balanced double squares is analyzed. In

a sense, it was a surprising result as no strong assumptions such as primitiveness of the squares or being the rightmost occurrences, are needed. This result is described in Chapter 2.

The realization that the inversion factors of the canonical factorization are such prominent structures as to completely remove the need for using the Three Squares Lemma for improving the upper bound for the maximum number of distinct squares in [8], lead to a renewed interest in the structure of three squares starting at the same proximity.

The New Periodicity Lemma [9] has a complicated proof as the combinatorial structure of three squares two of which start at the same position is quite forbidding. However, since when the two squares starting at the same position form a balanced double square can be dealt with via the inversion factors imposed by the canonical factorization, we embarked on the research. It required a refinement of the notion of the inversion factor, but lead to a new theorem with significantly reduced assumptions about the squares: originally, the smallest square was required to be regular – a very stringent requirement, and the bigger square was required to be primitive. This research and results lead to the publication *The New Periodicity Lemma Revisited* [3] and it is described in Chapter 3.

The canonical factorization also lead to a more straightforward and streamlined proof of the Three Squares Lemma with weaker assumptions as it deals simply with the cases when two smallest or the two biggest squares form a balanced double square. This research and results lead to the publication

On a lemma of Crochemore and Rytter [1] and are described in Chapter 4.

The thesis is concluded with Chapter 5 containing a conclusion and ideas for future research.

## 1.2 Preliminaries

**Definition 1.1.** An alphabet  $\Sigma$  is a non-empty set of elements referred to as symbols or letters. Frequently, it is required to be finite, but it can also be infinite. A string  $\mathbf{x}$  over an alphabet  $\Sigma$  is a contiguous sequence of symbols drawn from  $\Sigma$ . A string without any symbol is called the empty string and denoted by  $\epsilon$ . A string is also often called a word.

**Definition 1.2.** The alphabet of a string is the set of all of the symbols that occur in that string.

string	Example	Non-example
$aabba$	$\{a, b\}$	$\{a, c\}$
12445	$\{1, 2, 4, 5\}$	$\{3\}$
$abcca$	$\{a, b, c\}$	$\{a, b\}$
$ab^+$	$\{a, b\}$	$\{a\}$
123323	$\{1, 2, 3\}$	$\{1, 2, 3, 4, \dots\}$

Table 1.1: Examples and Non-examples of alphabet of a string

**Definition 1.3.** The length of a finite string  $\mathbf{x}$  is the number of symbols it is composed of, denoted as  $|\mathbf{x}|$ . The empty string has length zero. A finite string  $\mathbf{x}$  with  $n$  symbols can be represented as an array  $\mathbf{x}[1..n]$  or as a sequence  $\mathbf{x}[1]\mathbf{x}[2]\dots\mathbf{x}[n]$ .

one-sided infinite string	bi-infinite string
has a first element, no final element	neither a first nor a final element
$aab\dots$	$\dots aab\dots$
$\dots bbbb$	$\dots bbbb\dots$

Table 1.2: Examples of infinite strings

**Definition 1.4.** Concatenation is a basic operation of strings that joins two strings  $\mathbf{x}$  and  $\mathbf{y}$  together into one, which is denoted in the order that the two strings are concatenated as  $\mathbf{xy}$ . For instance,  $\mathbf{xy}$  represents a string that is the concatenation of string  $\mathbf{x}$  followed by string  $\mathbf{y}$ ; that is,  $\mathbf{xy} = \mathbf{x}[1][2]\dots\mathbf{x}[i]\mathbf{y}[1]\mathbf{y}[2]\dots\mathbf{y}[j]$ , where  $\mathbf{x} = \mathbf{x}[1\dots i]$  and  $\mathbf{y} = \mathbf{y}[1\dots j]$  and  $i, j \geq 1$ .

**Definition 1.5.** If a string  $\mathbf{x} = \mathbf{uvw}$ , where  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  are strings, then  $\mathbf{u}$  (respectively,  $\mathbf{v}$ ,  $\mathbf{w}$ ) is said to be a prefix (respectively, substring, suffix) of  $\mathbf{x}$  if  $0 \leq |\mathbf{u}| \leq |\mathbf{x}|$  (respectively,  $0 \leq |\mathbf{v}| \leq |\mathbf{x}|$ ,  $0 \leq |\mathbf{w}| \leq |\mathbf{x}|$ ); a proper prefix (respectively, proper substring, proper suffix) if  $0 < |\mathbf{u}| < |\mathbf{x}|$  (respectively,  $0 < |\mathbf{v}| < |\mathbf{x}|$ ,  $0 < |\mathbf{w}| < |\mathbf{x}|$ ). A substring is also called a factor or subword.

string	Example	Non-example
$aabba$	$a, aa, aab, aabb, ab, abb, abba, b, bb, bba, ba$	$\epsilon$
$12445$	$1, 12, 124, 1244, 2, 24, 244, 2445, 4, 44, 445, 45, 5$	$12445$
$abcca$	$a, ab, abc, abcc, b, bc, bcc, bcca, c, cc, cca, ca$	$ccc$
$abbb$	$a, ab, abb, b, bb, bbb$	$aba$
$323$	$3, 32, 2, 23$	$3234$

Table 1.3: Examples and Non-examples of proper substrings of a string

**Definition 1.6.** The symbol  $\Sigma^*$  denotes the set of all finite strings over the alphabet  $\Sigma$ , including the empty string  $\epsilon$ .

**Definition 1.7.** An order (or ordering)  $\prec$  of an alphabet  $\Sigma$  is a transitive, reflexive, and antisymmetric relation  $\prec$  over the alphabet that is total, i.e. that for every pair of distinct symbols  $x, y \in \Sigma$ , either  $x \prec y$  or  $y \prec x$ . Thus, the alphabet can be diagrammed as a line of symbols.

**Definition 1.8.** The lexicographic ordering  $\prec$  is an ordering of  $\Sigma^*$  induced by the ordering  $\prec$  of the alphabet  $\Sigma$ : For every pair of strings  $x = x_0x_1\dots x_{k-1}$  and  $y = y_0y_1\dots y_{l-1}$  in  $\Sigma^*$ ,  $x \prec y$  if there is an index  $i \geq 0$  such that  $i < k$  and  $i < l$  and  $x_i \prec y_i$  and  $x_j = y_j$  for all  $0 \leq j < i$ ; or if  $l > k$  and for all  $0 \leq j < k$  we have  $x_j = y_j$ .

order	lexicographic order
$a \prec b \prec c$	$ab \prec abb \prec abbc \prec abc \prec acbc \prec cbc$
$a \prec c \prec b$	$acbc \prec ab \prec abc \prec abb \prec abbc \prec cbc$
$b \prec a \prec c$	$ab \prec abb \prec abbc \prec abc \prec acbc \prec cbc$
$b \prec c \prec a$	$cbc \prec ab \prec abb \prec abbc \prec abc \prec acbc$
$c \prec a \prec b$	$cbc \prec acbc \prec ab \prec abc \prec abb \prec abbc$
$c \prec b \prec a$	$cbc \prec acbc \prec ab \prec abc \prec abb \prec abbc$

Table 1.4: lexicographic order of strings  $ab$ ,  $abb$ ,  $abc$ ,  $cbc$ ,  $abbc$ ,  $acbc$ , over alphabet  $\{a, b, c\}$

**Observation 1.9.** For an alphabet of size  $n$ , there are  $C_n^1 \cdot C_{n-1}^1 \cdot C_{n-2}^1 \cdot C_{n-3}^1 \dots \cdot 1 = n \cdot (n-1) \cdot (n-2) \cdot (n-3) \cdot \dots \cdot 3 \cdot 2 \cdot 1 = n!$  different orderings. I.e. there are as many different orderings of  $\Sigma^*$  as there are permutations of  $|\Sigma|$ .

**Definition 1.10.** A string  $\mathbf{x}$  is primitive if there are no string  $\mathbf{u}$  and an integer  $k \geq 2$  so that  $\mathbf{x} = \mathbf{u}^k = \underbrace{\mathbf{u}\mathbf{u}\dots\mathbf{u}}_{k \text{ times}}$ . For every string  $\mathbf{x}$  there is a unique smallest string  $\mathbf{u}$  so that  $\mathbf{x} = \mathbf{u}^k$  for some integer  $k \geq 1$ ;  $\mathbf{u}$  is necessarily primitive and is called the primitive root of  $\mathbf{x}$ .

**Definition 1.11.** For a concatenation  $\mathbf{u}^k$ ,  $k \geq 2$ ,  $\mathbf{u}_{[1]}$  denotes the first occurrence of  $\mathbf{u}$  in  $\mathbf{u}^k$ ,  $\mathbf{u}_{[2]}$  denotes the second occurrence of  $\mathbf{u}$  in  $\mathbf{u}^k$ , ...,  $\mathbf{u}_{[k]}$  denotes the  $k$ -th occurrence of  $\mathbf{u}$  in  $\mathbf{u}^k$ .

Example	non-Example
$ab$	$abab$
$a$	$aaaa$
$aab$	$aabaab$
$abac$	$acac$
$aabca$	$abaabaaba$

Table 1.5: Examples and Non-examples of primitive strings

**Definition 1.12.** A repeat is a collection of identical substrings of  $\mathbf{x}[1..n]$  described by  $\mathbf{u} = \mathbf{x}[i_1..i_1 + p - 1] = \mathbf{x}[i_2..i_2 + p - 1] = \dots = \mathbf{x}[i_q..i_q + p - 1]$ , where  $1 \leq i_1 < i_2 < \dots < i_q \leq n$ ,  $q \geq 2$ , and  $p \geq 1$ .  $\mathbf{u}$  is called the generator (or the root) of the repeat.

**Definition 1.13.** A repetition is a tandem repeat; that is, for every two consecutive identical substrings in the repeat, the gap between their starting positions is equal to the size of the generator. A repetition with generator  $\mathbf{u}$  repeating  $q$  times ( $q$  an integer) can be presented as  $\mathbf{u}^q$ , where  $\mathbf{u}$  is a non-empty string, and  $q \geq 2$ . The integer value  $|\mathbf{u}|$  is called the period and  $q$  the exponent or power of the repetition.

Example	Non-example
$abab$	$ab$
$aaaa$	$a$
$abaaba$	$abaabaa$
$aabaab$	$baabaab$
$abcabcabc$	$cabcabcabc$

Table 1.6: Examples and Non-examples of repetitions

**Definition 1.14.** A square is a repetition with power of 2. A cube is a repetition with power 3.

Example	Non-example
$abab$	$ab$
$aaaa$	$a$
$abaaba$	$aabaaba$
$aa$	$aaa$
$aabaaaba$	$aabab$

Table 1.7: Examples and Non-examples of squares

Example	Non-example
$ababab$	$abab$
$aaaaaa$	$aa$
$bbb$	$abaaba$
$abcabcabc$	$abc$
$aabaabaab$	$aabab$

Table 1.8: Examples and Non-examples of cubes

**Definition 1.15.** If  $\mathbf{u}$  is primitive, then the square  $\mathbf{u}^2$  is called primitively rooted.

The square  $\mathbf{u}^2$  is regular if no prefix of  $\mathbf{u}$  is a square.

Example	Non-example
$(ab)ab$	$(abab)abab$
$(a)a$	$(aa)aa$
$(abc)abc$	$(abcabc)abcabc$
$(aba)aba$	$(aaa)aaa$
$(aab)aab$	$(bbbb)bbbb$

Table 1.9: Examples and Non-examples of primitively rooted squares – ‘()’ shows the root of the square

**Definition 1.16.** *If a repetition  $\mathbf{u}^q = \mathbf{x}[s..s + qp - 1]$  where  $p = |\mathbf{u}|$ , can be extended by another copy of  $\mathbf{u}$  to the left of  $\mathbf{x}$ , that is,  $\mathbf{x}[s - p..s - 1] = \mathbf{u}$ , then we say the repetition can be extended to the left. A repetition is called a left-maximal repetition if it cannot be extended to the left. Similarly, if we cannot extend by another copy of  $\mathbf{u}$  to the right, then the repetition is a right-maximal repetition. A maximal repetition refers to a repetition that can be extended neither to the left nor to the right.*

string	Example
$\mathbf{x} = cccababccc$	$ccc(abab)ccc$
$\mathbf{x} = ababaaaababa$	$abab(aaaa)baba$
$\mathbf{x} = bbabaababb$	$bb(abaaba)bb$
$\mathbf{x} = abaaabaabbab$	$aba(aabaab)bab$
$\mathbf{x} = abbabcabccba$	$abb(abcabc)cba$

Table 1.10: Examples of maximal repetitions in a string  $\mathbf{x}$  – ‘()’ shows the maximal repetition

Non-example
left extendable: $cccb(abab)cccc$
right extendable: $abab(aaaa)abab$
two-side extendable : $aaa(abaaba)aaa$
right extendable: $aba(aabaab)aba$
two-side extendable : $cbc(abcabc)aba$

Table 1.11: Non-examples of maximal repetitions – ‘()’ shows the repetition

**Definition 1.17.** A run in a string  $\mathbf{x} = \mathbf{x}[1..n]$  starting at a position  $s$ , ending at a position  $e$  with a period  $p$  is a substring  $\mathbf{x}[s..e]$  so that

1.  $\mathbf{x}[s..e] = \mathbf{x}[s..s+p-1]^r \mathbf{t}$  where the integer  $r \geq 2$  and  $\mathbf{t}$  called tail is a proper prefix of the root  $\mathbf{x}[s..s+p-1]$ ;
2. the root  $\mathbf{x}[s..s+p-1]$  is primitive;
3. the root  $\mathbf{x}[s..s+p-1]$  cannot be shifted left, i.e. either  $s = 1$  or  $\mathbf{x}[s-1] \neq \mathbf{x}[s+p-1]$ ;
4.  $\mathbf{x}[e-p+1..e]$  cannot be shifted right, i.e. either  $e = n$  or  $\mathbf{x}[e+1] \neq \mathbf{x}[e-p+1]$ .

A run can thus be encoded by a triple  $r = (s, e, p)$  where  $s$ ,  $e$ , and  $p$  specify the starting position, the ending position, and the period of the run, respectively. Note that the run without its tail, e.i.  $\mathbf{x}[s..s+p-1]^r$ , is a maximal repetition, but not every maximal repetition is a run.

string	Example	Non-example
<i>ccababccc</i>	<i>cc(abab)ccc</i>	<i>ccabab(cc)c</i>
<i>baaaab</i>	<i>b(aaaa)b</i>	<i>ba(aa)ab</i>
<i>aaaabaabaaaa</i>	<i>aa(abaabaa)aa</i>	<i>(aaa)abaabaaaa</i>
<i>abaabbabbab</i>	<i>ab(aa)bbabbab</i>	<i>abaa(bbabba)b</i>
<i>acaacaacabcabcac</i>	<i>acaaca(cabcabca)c</i>	<i>(acaacaac)abcabcac</i>

Table 1.12: Examples and Non-examples of runs in a string  $x$ , '()' shows the run and non-run

**Definition 1.18.** *The maximum-number-of-runs problem counts the maximum number of runs in a string, i.e. the occurrences of the runs are counted rather than their types.*

**Definition 1.19.** *The maximum-number-of-distinct-squares problem counts the maximum number of distinct squares in a string, i.e. the types of the squares in the string are counted rather than their occurrences.*

string	square-type	run-occurrence
<i>ababcabab</i>	$(ab)^2 \times 1$	$(ab)^2 \times 2$
<i>abababab</i>	$(ab)^2 \times 1, (ba)^2 \times 1$	$(ab)^4 \times 1$
<i>aaaaaa</i>	$a^2 \times 1$	$a^6 \times 1$

Table 1.13: Examples of maximum-number-of-distinct-squares *vs* maximum-number-of-runs when we only count distinct primitively rooted squares

string	square-type	run-occurrence
<i>aaaabaacaabaaaa</i>	$a^2 \times 1, (aa)^2 \times 1$	$a^4 \times 2, a^2 \times 2$
<i>aaaaaa</i>	$a^2 \times 1, (aa)^2 \times 1, (aaa)^2 \times 1$	$a^6 \times 1$
<i>aaaa</i>	$a^2 \times 1, (aa)^2 \times 1$	$a^4 \times 1$

Table 1.14: Examples of maximum-number-of-distinct-squares *vs* maximum-number-of-runs when we count distinct squares no matter whether their roots are primitive or not

**Definition 1.20.** *The factorization or decomposition of a string is a collection of disjoint substrings such that their concatenation gives the whole string.*

**Definition 1.21.** [Lempel-Ziv factorization]

Let  $\mathbf{x} = \mathbf{x}[1..n]$  be a string of length  $n$  over an alphabet  $\Sigma$ . The Lempel-Ziv factorization of  $\mathbf{x}$  is a factorization  $\mathbf{x} = \mathbf{w}_1\mathbf{w}_2\dots\mathbf{w}_m$  such that each  $\mathbf{w}_k$ ,  $1 \leq k \leq m$ , is either:

- a letter  $c \in \Sigma$  that does not occur in  $\mathbf{w}_1\mathbf{w}_2\dots\mathbf{w}_{k-1}$ , or
- the longest substring of  $\mathbf{x}$  that occurs at least twice in  $\mathbf{w}_1\mathbf{w}_2\dots\mathbf{w}_k$ .

string	L-Z factorization	Non-example
aaaa	$a, aaa$	$aa, aa$
abaaaba	$a, b, a, aa, ba$	$a, b, aaa, ba$
aabaab	$a, a, b, aab$	$a, a, ba, a, b$
aaabcac	$a, aa, b, c, a, c$	$a, a, a, b, c, a, c$
abababa	$a, b, ababa$	$a, b, ab, ab, a$

Table 1.15: Examples and Non-examples of Lempel-Ziv factorization

**Definition 1.22.** *The notion of right rotation is defined by induction:*

- A string  $\mathbf{x}$  is a right 1-rotation of a string  $\mathbf{y}$  iff  $\mathbf{y} = \mathbf{y}[1..n]$  and  $\mathbf{x} = \mathbf{y}[2..n]\mathbf{y}[1]$ ,  $\mathbf{x}$  is also denoted as  $R_1(\mathbf{y})$ .
- A string  $\mathbf{x}$  is a right  $k$ -rotation,  $k \geq 2$ , of a string  $\mathbf{y}$  iff  $\mathbf{x}$  is a right 1-rotation of  $\mathbf{z}$  which is a right  $(k-1)$ -rotation of  $\mathbf{y}$ ,  $\mathbf{x}$  is also denoted as  $R_k(\mathbf{y})$ .

*Similarly for left rotation.*

*A string  $\mathbf{x}$  is referred to as a trivial rotation of itself. A string  $\mathbf{x}$  is a rotation of a string  $\mathbf{y}$  if it is a left or right  $k$ -rotation of  $\mathbf{y}$  for some  $k \geq 1$  or if  $\mathbf{x} = \mathbf{y}$ . If  $\mathbf{x}$  is a rotation of  $\mathbf{y}$ , we may also say that  $\mathbf{x}$  and  $\mathbf{y}$  are conjugates or that  $\mathbf{x}$  is a conjugate of  $\mathbf{y}$ .*

**Definition 1.23.** *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.*

Note a run can be described as a maximal right cyclic shift of the generator which itself cannot be left cyclic shifted; the length of the right cyclic shift must be at least twice the length of the generator. Similarly, a ran can be described as a maximal right shift of the leading square of the run which itself cannot be left cyclic shifted; the right cyclic shift can be of any length, including 0. The leading square of a run is its leftmost square.

**Definition 1.24.** *Strings  $\mathbf{uv}$  and  $\mathbf{vu}$  are conjugates, written  $\mathbf{uv} \sim \mathbf{vu}$ . We also say that  $\mathbf{vu}$  is the  $|\mathbf{u}|^{th}$  rotation of  $\mathbf{x}$ , written  $R_{|\mathbf{u}|}(\mathbf{x})$ , or the  $-|\mathbf{v}|^{th}$  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.*

**Example**

In a string  $\mathbf{x} = abbaababbaabb$ :

$bbaa$  is the cyclic shift of  $abba$ , and  $bba|a \sim a|bba$ ,

while  $abb|a \sim a|abb$ , but  $aabb$  is not a cyclic shift of  $abba$

**Definition 1.25.** Given strings  $\mathbf{u}$  and  $\mathbf{v}$ ,  $lcp(\mathbf{u}, \mathbf{v})$  (respectively,  $lcs(\mathbf{u}, \mathbf{v})$ ) denotes the length of the longest common prefix (respectively, longest common suffix) of  $\mathbf{u}$  and  $\mathbf{v}$ .

**Definition 1.26.** A string  $\mathbf{x}$  over an ordered alphabet  $(\mathcal{A}, \prec)$  is Lyndon w.r.t. (with respect to)  $\prec$  iff  $\mathbf{x} \prec \mathbf{y}$  for any non-trivial rotation  $\mathbf{y}$  of  $\mathbf{x}$ .

A string of length one is thus Lyndon, and it is referred to as trivial Lyndon word.

**Observation 1.27.** Let  $\mathbf{x} = \mathbf{x}[1..n]$ ,  $n > 1$  be a string over an ordered alphabet  $(\mathcal{A}, \prec)$ . Then the following are equivalent:

1.  $\mathbf{x}$  is a non-trivial Lyndon word w.r.t.  $\prec$ .
2.  $\mathbf{x}[1..j-1] \prec \mathbf{x}[j..n]$  for any  $1 < j \leq n$ .
3.  $\mathbf{x}[1..n] \prec \mathbf{x}[j..n]$  for any  $1 < j \leq n$ .
4. there exists  $1 < j \leq n$  so that both  $\mathbf{x}[1..j-1]$  and  $\mathbf{x}[j..n]$  are Lyndon.

Item 4 is the basis of the definition of the standard factorization of a Lyndon word  $\mathbf{x}$ : it is a pair of Lyndon words  $\mathbf{u}, \mathbf{v}$  so that  $\mathbf{x} = \mathbf{uv}$  where  $\mathbf{v}$  is maximal such (or, equivalently  $\mathbf{u}$  is minimal such).

**Definition 1.28.** A Lyndon substring  $\mathbf{x}[i..j]$  of a string  $\mathbf{x} = \mathbf{x}[1..n]$  is maximal Lyndon substring of  $\mathbf{x}$  if either  $j = n$  or  $\mathbf{x}[i..k]$  is not Lyndon for any  $j < k \leq n$ .

A Lyndon substring  $\mathbf{x}[i..j]$  of  $\mathbf{x}$  is a non-extensible Lyndon substring  $\mathbf{x}$  if  $\mathbf{x}[i..j]$  is a maximal Lyndon substring of  $\mathbf{x}\mathbf{y}$  for any  $\mathbf{y}$ .

string	Example of Lyndon factorization	Non-example
<i>abbcc</i>	<i>a, bbcc</i>	<i>ab, acc</i>
<i>aab</i>	<i>a, ab</i>	<i>aa, b</i>
<i>aaabab</i>	<i>a, aabab</i>	<i>aa, abab</i>
<i>abacc</i>	<i>ab, acc</i>	<i>aba, cc</i>
<i>aababb</i>	<i>a, ababb</i>	<i>aab, abb</i>

Table 1.16: Examples and Non-examples of Standard factorization of a Lyndon word

The following theorem is due to the Chen, Fox, and Lyndon, [5], though it is not stated there in this form:

**Theorem 1.29.** [Lyndon factorization]

For any string  $\mathbf{x}$  over an ordered alphabet  $(\mathcal{A}, <)$ , there is a unique factorization of  $\mathbf{x}$  into  $m$  factors,  $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2\dots\mathbf{x}_m$ , such that each factor  $\mathbf{x}_i$ ,  $1 \leq i \leq m$ , is a maximal Lyndon word, and  $\mathbf{x}_1 \succeq \mathbf{x}_2 \succeq \dots \succeq \mathbf{x}_m$ .

string	Example of Lyndon factorization	Non-example
<i>aaaa</i>	<i>a, a, a, a</i>	<i>aa, aa</i>
<i>aba</i>	<i>ab, a</i>	<i>a, ba</i>
<i>aabaab</i>	<i>aab, aab</i>	<i>a, abaab</i>
<i>abaaba</i>	<i>ab, aab, a</i>	<i>aba, aba</i>
<i>abccba</i>	<i>abccb, a</i>	<i>abc, cba</i>

Table 1.17: Examples and Non-examples of Lyndon factorization

The Synchronization Principle and Common Factor Lemma are important tools for investigation of canonical factorizations, so we state them here. Lemma 1.30 will be used in the proof of the Synchronization Principle. The symbol  $|$  represents *divides*, i.e.  $a|b$  means that  $a$  divides  $b$ .

**Lemma 1.30** ([17], Lemma 1.4.2). *Let  $\mathbf{x}$  be a string of length  $n$  and minimum period  $\pi \leq n$ , let  $j$  be an integer and  $1 \leq j < n$ . Then  $R_j(\mathbf{x}) = \mathbf{x}$  if and only if  $\mathbf{x}$  is not primitive ( $\pi < n$ ,  $\pi | n$ ) and  $\pi | j$ .*

string	$n$	$\pi$	$j$	$R_j(\mathbf{x})$
<i>bababa</i>	6	2	4	<i>bab<u>aba</u></i>
<i>aaa</i>	3	1	3	<i><u>aaa</u></i>

Table 1.18: Example of strings with  $R_j(\mathbf{x}) = \mathbf{x}$ , where underlined letters are those that were rotated

string	$n$	$\pi$	$j$	$R_j(\mathbf{x})$	comment
<i>ababa</i>	5	5	1	<i>bab<u>aa</u></i>	$\mathbf{x}$ is primitive, $\pi \nmid j$ , $R_j(\mathbf{x}) \neq \mathbf{x}$
<i>ababab</i>	6	2	3	<i>bab<u>aba</u></i>	$\mathbf{x}$ is not primitive but $\pi \nmid j$ , $R_j(\mathbf{x}) \neq \mathbf{x}$

Table 1.19: Examples of strings with  $R_j(\mathbf{x}) \neq \mathbf{x}$ , where underlined letters are those that were rotated. Moreover, when  $\mathbf{x}$  is primitive,  $n = \pi$  and  $j < n$ ,  $\pi$  can not divide  $j$ .

**Lemma 1.31** (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.* Suppose there exists another occurrence of  $\mathbf{x}$  except those occurrences determined by  $\mathbf{x}^p$ , then it must start at the same position with a rotation of

$\mathbf{x}$ , which means  $\mathbf{x} = R_j(\mathbf{x})$  somewhere in the string. Thus  $\mathbf{x}$  is not primitive by Lemma 1.30, contradiction.  $\square$

string	$\mathbf{x}$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{p}$	comment
<u>cabca</u>	abc	a	c	1	$\mathbf{x}$ can't shift
ab <u>abbababbababb</u>	abbab	abb	ab	3	$\mathbf{x}$ can't shift

Table 1.20: Illustrations of Synchronization Principle holding for primitive string  $\mathbf{x}$ , the underlined letters are occurrences of  $\mathbf{x}$

string	$\mathbf{x}$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{p}$	comment
<u>babbabbabbabbab</u>	babbab	bab	bab	2	$\mathbf{x}$ occurs 4 times $> \mathbf{p}$
<u>bababababab</u>	abab	ab	b	2	$\mathbf{x}$ occurs 5 times $> \mathbf{p}$

Table 1.21: Illustrations of Synchronization Principle not holding for non-primitive string  $\mathbf{x}$ , the underlined letters are occurrences of  $\mathbf{x}$

**Lemma 1.32** (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  $\mathbf{x} \sim \mathbf{y}$ .*

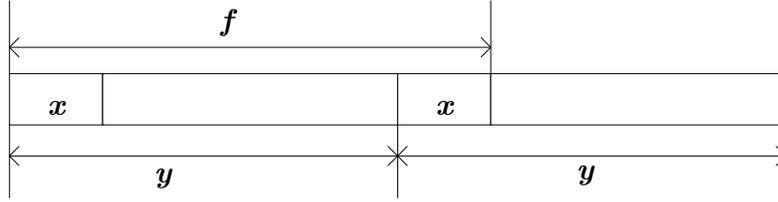
*Proof.* First consider the special case  $\mathbf{x}_1 = \mathbf{x}_2 = \mathbf{y}_1 = \mathbf{y}_2 = \varepsilon$ , where  $\mathbf{x}^p, \mathbf{y}^q$  have a common prefix  $\mathbf{f}$  of length  $|\mathbf{x}| + |\mathbf{y}|$ . We show that in this case  $\mathbf{x} = \mathbf{y}$ .

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 without loss of generality 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.

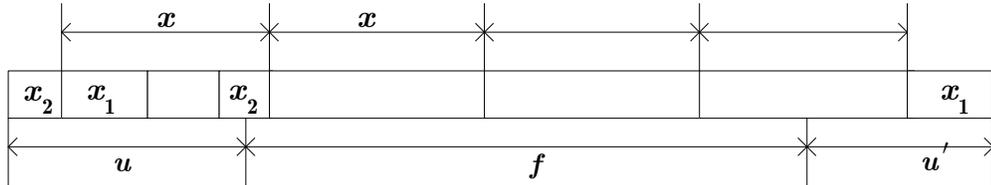
Hence there exists  $k \geq 1$  such that  $k|\mathbf{x}| < |\mathbf{y}|$  and  $(k+1)|\mathbf{x}| > |\mathbf{y}|$ . But since  $\mathbf{f} = \mathbf{y}\mathbf{x}$ , it follows that

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

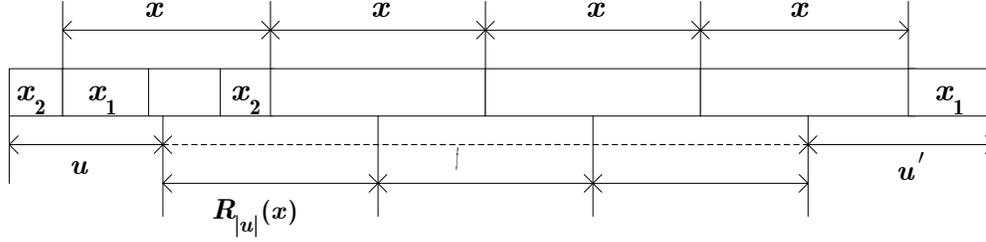
again by Lemma 1.30 contrary to the assumption that  $\mathbf{x}$  is primitive. We conclude that  $|\mathbf{x}| \not< |\mathbf{y}|$ , hence that  $|\mathbf{x}| = |\mathbf{y}|$  and  $\mathbf{x} = \mathbf{y}$ , as required.



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  $\mathbf{x}_2\mathbf{x}^{p-1}\mathbf{x}_1$ .



And so we can assume without loss of generality 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.  $\square$

string	$\mathbf{x}$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{p}$	length
<u>a</u> <u>ab</u> <u>ba</u> <u>ab</u> <u>ba</u> <u>ab</u>	<u>ab</u> <u>ba</u>	<u>ab</u>	<u>a</u>	3	4
string	$\mathbf{y}$	$\mathbf{y}_1$	$\mathbf{y}_2$	$\mathbf{p}$	length
<u>a</u> <u>bb</u> <u>a</u> <u>ab</u> <u>ba</u> <u>ab</u> <u>ba</u> <u>ab</u> <u>ba</u>	<u>ba</u> <u>ab</u>	<u>ba</u> <u>a</u>	<u>ab</u>	4	4

Table 1.22: Example of  $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$  and  $\mathbf{y}_2\mathbf{y}^q\mathbf{y}_1$  having a common factor of length  $|\mathbf{x}| + |\mathbf{y}|$ , then  $\mathbf{x} \sim \mathbf{y}$ . Underlined letters are the common factor.

string	$\mathbf{x}$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{p}$	length
<u>b</u> <u>ca</u> <u>b</u> <u>ca</u> <u>b</u> <u>ca</u>	<u>abc</u>	<u>a</u>	<u>bc</u>	2	3
string	$\mathbf{y}$	$\mathbf{y}_1$	$\mathbf{y}_2$	$\mathbf{p}$	length
<u>c</u> <u>ab</u> <u>ca</u> <u>bc</u> <u>ab</u> <u>ca</u> <u>bc</u>	<u>bca</u>	<u>bc</u>	<u>ca</u>	3	3

Table 1.23: Example of  $\mathbf{x}_2\mathbf{x}^p\mathbf{x}_1$  and  $\mathbf{y}_2\mathbf{y}^q\mathbf{y}_1$  having a common factor of length  $|\mathbf{x}| + |\mathbf{y}|$ , then  $\mathbf{x} \sim \mathbf{y}$ . Underlined letters are the common factor.

Note that Lemma 1.32 could be equivalently stated in a more general form:

**Lemma 1.33.** *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.

string	$\mathbf{x}$	$\bar{\mathbf{x}}$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{p}$	length
<u>bcabcabcabcabcabc</u>	abcabc	abc	ab	bcabc	2	6
string	$\mathbf{y}$	$\bar{\mathbf{y}}$	$\mathbf{y}_1$	$\mathbf{y}_2$	$\mathbf{p}$	length
<u>abcabcabcabc</u>	cabcab	cab	c	ab	2	6

Table 1.24: Example of the primitive root  $\bar{\mathbf{x}}$  of  $\mathbf{x}$  and the primitive root  $\bar{\mathbf{y}}$  of  $\mathbf{y}$  being conjugates. Underlined letters are the common factor.

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

**Lemma 1.34.** *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 1.32  $\mathbf{x} \sim \mathbf{y}$ . 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 1.34(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  $\mathbf{x} \neq \mathbf{y}$ .

## Chapter 2

# Canonical factorization of double squares

In this chapter based on [2], we discuss the most general form of canonical factorization of double squares. But first we need to defined precisely what is meant by a *double square*.

The term *double square* is due to Lam, [15] and he defined it as two rightmost occurrences of a square  $\mathbf{u}^2$  and a square  $\mathbf{v}^2$  starting at the same position of a string  $\mathbf{x}$ . Following on work of Fraenkel and Simpson [10], he tried to obtain a better bound for the number of distinct squares by bounding the number of double squares. In [8], such double squares are referred to as *FS-double squares*, to acknowledge the pioneering work of Fraenkel and Simpson in this field and to prevent a confusion with configurations of two squares starting at the same position. The bounding of FS-double squares

carried in [8] gave an improved upper bound of  $\frac{11}{6}|\mathbf{x}|$  for the number of distinct squares in  $\mathbf{x}$ .

From now on in this thesis, we will use the term *FS-double square* for two rightmost occurrences of squares starting at the same position, while the term *double square* for us is any two squares starting at the same position. We use the following notation:  $DS(\mathbf{u}, \mathbf{v})$  indicates that the squares  $\mathbf{u}^2$  and  $\mathbf{v}^2$  start at the same position and that  $|\mathbf{u}| < |\mathbf{v}|$ . For a double square  $DS(\mathbf{u}, \mathbf{v})$ , we say that the squares  $\mathbf{u}^2$  and  $\mathbf{v}^2$  are *proportional* if  $|\mathbf{u}| < |\mathbf{v}| < 2|\mathbf{u}| < 2|\mathbf{v}|$ . If the two squares of a double square are proportional, then the double square is called *balanced*.

Lam, [15], tried to provide a taxonomy of mutual relationships of FS-double squares, and for that purpose he described a factorization of an FS-double square into primitive components. In [8], some important properties of such factorization were discussed and the factorization was shown to be unique, i.e. it is shown that every FS-double square  $(\mathbf{u}, \mathbf{v})$  has a unique factorization  $\mathbf{u}_1, \mathbf{u}_2, e_1$  and  $e_2$  so that  $\mathbf{u}^2 = (\mathbf{u}_1\mathbf{u}_2)^2$  and  $\mathbf{v}^2 = \mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2\mathbf{u}_1^{e_2}$  where  $\mathbf{u}_1$  is primitive,  $\mathbf{u}_2$  is a non-trivial proper prefix of  $\mathbf{u}_1$ , and  $e_1, e_2$  are integers so that  $e_1 \geq e_2 \geq 1$ . The term *canonical factorization* was used there.

The main result of [2] is the fact that any balanced double square admits the canonical factorization. In a sense, it was a big surprise that no additional conditions are required. For the Three Squares Lemma, the squares were required to be primitive, or as Fraenkel and Simpson showed, the smallest

one must be primitive. For the New Periodicity Lemma, the regularity of the smallest square was required, even stronger property than primitiveness. The canonical factorization for FS-double squares was not surprising as being the rightmost occurrence is again stronger than being primitive.

The main result of this chapter is Lemma 2.1, *Two Squares Factorization Lemma*, which specifies the unique factorization imposed by the occurrence of two proportional squares at the same position in a string.

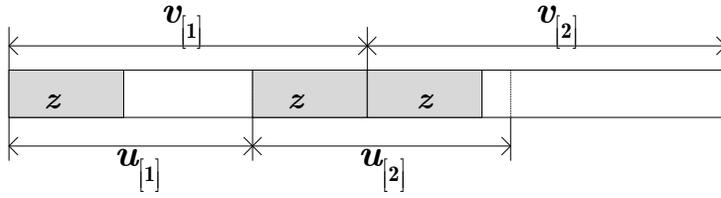
## 2.1 Two Squares Factorization

**Lemma 2.1** (Two Squares Factorization Lemma). *For a balanced 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  $\mathbf{u}$  is also primitive.*

*In both cases, the factorization is unique.*

*Proof.* 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  $\mathbf{u}$ .

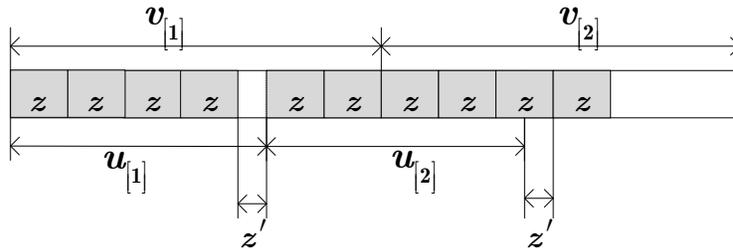


In general, there exists an integer  $k = \lfloor |u|/|z| \rfloor \geq 1$  such that  $u = z^k z'$  for some proper prefix  $z'$  of  $z$ . Let  $u_1$  be the primitive root of  $z$ , so that  $z = u_1^{e_2}$  for some integer  $e_2 \geq 1$ .

$$\begin{aligned} u &= z^k z' \\ &= u_1^{e_2 k} z' \\ &= u_1^{e_1} u_2 \end{aligned}$$

$$\begin{aligned} v &= uz \\ &= u_1^{e_1} u_2 u_1^{e_2} \end{aligned}$$

Therefore, for some  $e_1 \geq e_2 k$  and some prefix  $u_2$  of  $u_1$ ,  $u = u_1^{e_1} u_2$  and  $v = uz = u_1^{e_1} u_2 u_1^{e_2}$ , as required.



To prove the uniqueness of  $u_1$  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 1.34(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 1.34(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 without loss of generality 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

$$\mathbf{w}^{2k} = \mathbf{v}^2 = \mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2\mathbf{u}_1^{e_2}, \quad (2.1)$$

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 1.32 to conclude that  $\mathbf{w}$  and  $\mathbf{u}_1$  are conjugates, thus by (2.1) that  $\mathbf{w} = \mathbf{u}_1$ . But (2.1) 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  $\mathbf{u}_1$ , which means  $\mathbf{u}_2\bar{\mathbf{u}}_2 = \bar{\mathbf{u}}_2\mathbf{u}_2$  then  $\mathbf{u}_1$  is non-primitive by Lemma 1.30. This is a contradiction to 2.1 which requires  $\mathbf{u}_1$  to be primitive. Thus we conclude that  $\mathbf{v}$  is primitive.

Now suppose in addition that  $e_1 \geq 2$ , 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 2.1,  $\mathbf{w}^{2k}$  and  $\mathbf{u}_1^{e_1+e_2}$  have a common prefix  $\mathbf{u}_1^{e_1}\mathbf{u}_2$ . Note that  $|\mathbf{u}_1^{e_1}\mathbf{u}_2| \geq |\mathbf{w}| + |\mathbf{u}_1|$ , so that again applying Common Factor Lemma 1.32, 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 2.1(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  $\mathbf{v}^2 = abaabaabaababaabaabaab$ , 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  $\mathbf{u}$  is *not* primitive.

$\mathbf{u}^2$	$\mathbf{v}^2$	$\mathbf{u}_1$	$\mathbf{u}_2$
$abaabaaba \cdot abaabaaba$	$abaabaabaaba \cdot abaabaabaaba$	$aba$	$ \mathbf{u}_2  = 0$
$bbabba \cdot bbabba$	$bbabbabba \cdot bbabbabba$	$bba$	$ \mathbf{u}_2  = 0$

Table 2.1: Illustrations of Lemma 2.1 case (a), ‘.’ divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$ , in the first example,  $e_1 = 3$ ,  $e_2 = 1$ , and in the second example,  $e_1 = 2$ ,  $e_2 = 1$

$\mathbf{u}^2$	$\mathbf{v}^2$	$\mathbf{u}_1$	$\mathbf{u}_2$
$abcaabcaabc \cdot abcaabcaabc$	$abcaabcaabcabca \cdot abcaabcaabcabca$	$abca$	$abc$
$bcbcbcb \cdot bcbcbcb$	$bcbcbcbcbcb \cdot bcbcbcbcbcb$	$bc$	$b$

Table 2.2: Illustrations of Lemma 2.1 case (b), ‘.’ divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$ , in the first example,  $e_1 = 2$ ,  $e_2 = 1$ , and in the second example,  $e_1 = 3$ ,  $e_2 = 2$

$\mathbf{u}^2$	$\mathbf{v}^2$	$\mathbf{u}_1$	$\mathbf{u}_2$
$abccaabc \cdot abccaabc$	$abccaabcabcca \cdot abccaabcabcca$	$abcca$	$abc$
$aba \cdot aba$	$abaab \cdot abaab$	$ab$	$a$
$abaaba \cdot abaaba$	$abaabaabaab \cdot abaabaabaab$	$abaab$	$a$
$abaaabaa \cdot abaaabaa$	$abaaabaaabaaaba \cdot abaaabaaabaaaba$	$abaaaba$	$a$

Table 2.3: Example that the condition of Lemma 2.1 case (b) is not necessary : when  $e_1 = e_2 = 1$ ,  $\mathbf{u}$  can be either primitive or non-primitive, the first two examples are the case  $\mathbf{u}$  is not primitive, the last two examples are the case  $\mathbf{u}$  is primitive, ‘.’ divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$ , in all examples  $e_1 = e_2 = 1$

Lemma 2.1 gives credence to the following definition of terminology and notation:

**Definition 2.2.** For a balanced double square  $DS(\mathbf{u}, \mathbf{v})$ , we call the unique factorization guaranteed by Lemma 2.1, the canonical factorization of the double square 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  $\mathbf{u}_1 = \mathbf{u}_2\bar{\mathbf{u}}_2$ .

$\mathbf{u}^2$	$bbab \cdot bbab$
$\mathbf{v}^2$	$bbabbba \cdot bbabbba$
$\mathbf{u}_1$	$bbab$
$\mathbf{u}_2$	$b$
$\bar{\mathbf{u}}_2$	$bab$
$e_1$	$1$
$e_2$	$1$

Table 2.4: An example of canonical factorization, ‘.’ divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$

Lemma 2.1 also gives rise to a number of important observations:

**Observation 2.3.** In Lemma 2.1,  $|\mathbf{u}_2| > 0$  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 2.1,  $|\mathbf{u}_2| > 0$  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.

□

In the context of Observation 2.3(f), consider the double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$  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.

The following tables illustrate some of the cases discussed in this chapter.

$v$	$u$	$u_1$	$u_2$	$e_1$	$e_2$
$abaababa : v$ is primitive	$abaab$	$aba$	$ab :  u_2  > 0$	1	1
$ababab : v$ is non-primitive	$abab$	$ab$	$ u_2  = 0$	2	1

Table 2.5: Example of Observation 2.3 case (a)

$v$	$u$	$u_1$	$u_2$
$abaab$	$aba : u$ is primitive	$ab$	$a :  u_2  > 0$
$abcaabcaabcabca$	$abcaabcaabc : u$ is primitive	$abca$	$abc :  u_2  > 0$

Table 2.6: Example of Observation 2.3 case (b), in the first example,  $e_1 = e_2 = 1$ , and in the second example,  $e_1 = 2, e_2 = 1$ 

$v$	$u$	$u_1$	$u_2$	$e_1$	$e_2$
$abaabaabaab$	$abaaba$	$abaab$	$a$	1	1
$aaaaaaaa$	$aaaa$	$aaa$	$a$	1	1

Table 2.7: Example of the unnecessary of the condition in Observation 2.3 case (b):  $u$  is not primitive, but  $|u_2| > 0$ 

$v^2$	$u^2$	$u_1$	$u_2$	$e_1$	$e_2$
$\underline{baababaa} \cdot \underline{baababaa}$	$baaba \cdot baaba$	$baa$	$ba$	1	1
$\underline{abaab} \cdot \underline{abaab}$	$aba \cdot aba$	$ab$	$a$	1	1
$\underline{abbabbabb} \cdot \underline{abbabbabb}$	$abb \cdot abb$	$abb$	$ u_2  = 0$	2	1

Table 2.8: Example of Observation 2.3 case (c), in the third example, there exists other occurrence of  $u^2$  farther to the right in  $v^2$ , underlined letters are the right-most occurrence of  $u^2$  in  $v^2$ , '.' divides two roots of  $u^2$  and  $v^2$ 

$v^2$	$u^2$	$u_1$	$u_2$	$e_1$	$e_2$
$\underline{aaaaaaaa} \cdot \underline{aaaaaaaa}$	$aaaaa aaaaa$	$aaa$	$aa$	1	1

Table 2.9: Example of the unnecessary of the condition in Observation 2.3 case (c):  $u^2$  is not the rightmost in  $v^2$ , but  $|u_2| > 0$ , '.' divides two roots of  $u^2$  and  $v^2$

$\mathbf{v}$	$\mathbf{u}$	$\mathbf{u}_1$	$\mathbf{u}_2$	$e_1$	$e_2$
$abaaaba \cdot abaaaba$	$abaa \cdot abaa$	$aba$	$a$	1	1
$abcababc \cdot abcababc$	$abcab \cdot abcab$	$abc$	$ab$	1	1

Table 2.10: Example of Observation 2.3 case (d), '.' divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$

$\mathbf{v}$	$\mathbf{u}$	$\mathbf{u}_1$	$\mathbf{u}_2$	$e_1$	$e_2$
$abaabaaab \cdot abaabaaab$	<u><math>aa</math></u> $baaba \cdot aabaaba$	$aab$	$a$	2	1
$ababaaababa \cdot ababaaababa$	<u><math>ab</math></u> $abaa \cdot ababaa$	$ababa$	$a$	1	1

Table 2.11: Example of the unnecessary of the condition in Observation 2.3 case (d):  $\mathbf{u}^2$  is not regular, but  $|\mathbf{u}_2| > 0$ , underlined letters are the square prefix of  $\mathbf{u}^2$ , '.' divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$

$\mathbf{v}$	$\mathbf{u}$	$\mathbf{u}_1$	$\mathbf{u}_2$	$e_1$	$e_2$
$abcababc \cdot abcababc$	$abcab \cdot abcab$	$abc$	$ab$	1	1
$abbaabb \cdot abbaabb$	$abba \cdot abba$	$abb$	$a$	1	1

Table 2.12: Illustrations of Observation 2.3 case (f), '.' divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$

$\mathbf{v}$	$\mathbf{u}$	$\mathbf{u}_1$	$\mathbf{u}_2$	$e_1$	$e_2$
$abaabaabaab \cdot abaabaabaab$	<u><math>abaaba</math></u> $ \cdot abaaba$	$abaab$	$a$	1	1

Table 2.13: Example of the unnecessary of the condition in Observation 2.3 case (f):  $\mathbf{u}^2$  is not regular, but  $e_1 = e_2 = 1$  and  $\mathbf{u}_1$  is regular, underlined letters are the square prefix of  $\mathbf{u}^2$ , '.' divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$

## 2.2 Inversion Factors in a Balanced Double Square

Inversion factors were introduced in [8] and used in FS double squares, here we extend the concept to all double squares that admit special canonical factorization: consider a double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$  where  $|\mathbf{u}_2| > 0$ , and let  $\mathbf{u}_1 = \mathbf{u}_2\bar{\mathbf{u}}_2$ . Thus  $\mathbf{v}^2$  becomes

$$\begin{aligned}
\mathbf{v}^2 &= (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1} \mathbf{u}_2 (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2} \mathbf{u}_2 (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2} \\
&= (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1-1} (\mathbf{u}_2\bar{\mathbf{u}}_2) \mathbf{u}_2 (\mathbf{u}_2\bar{\mathbf{u}}_2) (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2-2} (\mathbf{u}_2\bar{\mathbf{u}}_2) \mathbf{u}_2 (\mathbf{u}_2\bar{\mathbf{u}}_2) (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2-1} \\
&= (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1-1} \mathbf{u}_2 (\bar{\mathbf{u}}_2 \mathbf{u}_2 \mathbf{u}_2 \bar{\mathbf{u}}_2) (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2-2} \mathbf{u}_2 (\bar{\mathbf{u}}_2 \mathbf{u}_2 \mathbf{u}_2 \bar{\mathbf{u}}_2) (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2-1} \\
&= (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1-1} \mathbf{u}_2 (\mathbb{IF}) (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2-2} \mathbf{u}_2 (\mathbb{IF}) (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2-1} \quad (2.2)
\end{aligned}$$

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

$\mathbf{v}^2$	$\mathbf{u}^2$	$\mathbf{u}_1$	$\mathbf{u}_2$	$\bar{\mathbf{u}}_2$	$\mathbb{IF}$
<u>ab</u> baa <u>ab</u> ba · <u>ab</u> baa <u>ab</u> ba	abbaa · abbaa	abba	a	bba	bbaaabba
<u>ab</u> ba <u>bb</u> ab <u>ab</u> · <u>ab</u> ba <u>bb</u> ab <u>ab</u>	abba <u>bb</u> ab · abba <u>bb</u> ab	abb	ab	b	bababb

Table 2.14: underlined letters emphasize the inversion factors, ‘·’ divides two roots of  $\mathbf{u}^2$  and  $\mathbf{v}^2$ , in the first example,  $e_1 = e_2 = 1$ , and in the second example,  $e_1 = 2$ ,  $e_2 = 1$

**Lemma 2.4.** *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  $\mathbb{IF}$  has exactly two occurrences in  $\mathbf{v}^2$  exactly a distance of  $|\mathbf{v}|$  apart as shown in (2.2).*

*Proof.* Let  $\mathbb{IF}_1$  denote the first inversion factor in a double square ,  $\mathbb{IF}_2$  denote the second inversion factor in a double square .

$$\mathbf{v}^2 = (\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1-1}\mathbf{u}_2(\mathbb{IF}_1)(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2-2}\mathbf{u}_2(\mathbb{IF}_2)(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_2-1}$$

The distance between  $\mathbb{IF}_1$  and  $\mathbb{IF}_2$  is  $|(\mathbf{u}_2\bar{\mathbf{u}}_2)^{e_1+e_2}| + |\mathbf{u}_2| = |\mathbf{u}_1^{e_1}\mathbf{u}_2\mathbf{u}_1^{e_2}| = |\mathbf{v}|$ . Suppose there exists an inversion factor occurring somewhere else than the positions of  $\mathbb{IF}_1$  and  $\mathbb{IF}_2$  in  $\mathbf{v}^2$ . According to the Synchronization principle Lemma 1.31, the subfactor  $\mathbf{u}_2\bar{\mathbf{u}}_2$  of the inversion factor must align with one of the three occurrences of  $\mathbf{u}_2\bar{\mathbf{u}}_2$  in  $\mathbf{v}^2$ , as  $\mathbf{u}_1 = \mathbf{u}_2\bar{\mathbf{u}}_2$  is primitive. On the other hand, another subfactor  $\bar{\mathbf{u}}_2\mathbf{u}_2$  always aligns with  $\mathbf{u}_2\bar{\mathbf{u}}_2$ . Thus  $\mathbf{u}_2\bar{\mathbf{u}}_2 = \bar{\mathbf{u}}_2\mathbf{u}_2$ , which means  $\mathbf{u}_1 = \mathbf{u}_2\bar{\mathbf{u}}_2$  is not primitive by the Lemma 1.30 , contradicting the primitiveness of  $\mathbf{u}_1$ .  $\square$

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 [8], 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 2.5** ([8]). *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  $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) \leq |\mathbf{u}_1| - 2$ .*

**Example**

$$\mathbf{u} = aababa, \mathbf{u}_1 = aabab, \mathbf{u}_2 = a, \bar{\mathbf{u}}_2 = abab$$

$$\mathbf{u}_2\bar{\mathbf{u}}_2 = aabab, \quad \bar{\mathbf{u}}_2\mathbf{u}_2 = ababa$$

$$\text{lcs} = 0, \text{lcp} = 1, \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 1 \leq |\mathbf{u}_1| - 2 = 3$$

In fact, in [8] 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 thesis):

**Lemma 2.6** ([8]). *Consider a double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$  with a non-empty  $\mathbf{u}_2$  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  $R_i(\mathbb{IF})$  or  $R_{-j}(\mathbb{IF})$  for some  $i \in 0..p$  or some  $j \in 0..s$ . Moreover, every  $R_i(\mathbb{IF})$  or  $R_{-j}(\mathbb{IF})$  appear exactly twice in  $\mathbf{v}^2$  exactly a distance  $|\mathbf{v}|$  apart for every  $i \in 0..p$  and every  $j \in 0..s$ .*

### Example

$$\mathbf{u} = aaabbaaabbaa, \quad \mathbf{u}_1 = aaabb, \quad \mathbf{u}_2 = aa, \quad \bar{\mathbf{u}}_2 = abb, \quad e_1 = 2, \quad e_2 = 2$$

$$\mathbf{u}_2\bar{\mathbf{u}}_2 = aaabb, \quad \bar{\mathbf{u}}_2\mathbf{u}_2 = abbaa$$

$$\mathbb{IF} = \bar{\mathbf{u}}_2\mathbf{u}_2\mathbf{u}_2\bar{\mathbf{u}}_2 = abbaaaaabb$$

$$s = \text{lcs} = 0, \quad p = \text{lcp} = 1, \quad R_1(\mathbb{IF}) = bbaaaaabba$$

$\mathbf{v}^2 = (aaabbaaabbaaaaabbaaabb)(aaabbaaabbaaaaabbaaabb)$ , where underlined letters are the two occurrences of  $R_1(\mathbb{IF})$ , the distance

$$\text{is } |aabbaaabbaaaaabba| = 22 = |\mathbf{v}|$$

## 2.3 Rare Factors in Balanced Double Squares

In this section, we discuss other factors with limited number of occurrences which we generically refer to as *rare factors* or *rarely occurring factors*. In particular, we introduce *right inversion subfactor* (RIS) and *left inversion subfactor* (LIS) in balanced double squares. These rare factors are half the size of the inversion factors, and, since they are shorter, they provide a stronger restriction on the positions and lengths of a third square. We will use them significantly in the next chapter for the New Periodicity Lemma.

In [19], Thierry discusses the core of the period interrupt, a very similar concept to the one we introduce here as RIS and LIS. For a double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$  with  $|\mathbf{u}_2| > 0$ , both RIS and LIS have only two occurrences in a significant portion of  $\mathbf{v}^2$ .

Lemma 2.7 will be used to prove rare occurrences of RIS and LIS in a balanced 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 2.7.** *Consider a string  $\mathbf{x} = \mathbf{x}[1..n]$ . If the substrings  $\mathbf{x}[i..i+k]$  and  $\mathbf{x}[i+1..i+1+k]$  are conjugates, then  $\mathbf{x}[i] = \mathbf{x}[i+k+1]$ .*

*Proof.* Since  $\mathbf{x}[i..i+k]$  and  $\mathbf{x}[i+1..i+1+k]$  are conjugates, the frequency of the alphabet symbols in both must be the same. Let  $\mathbf{x}[i] = a$ . Then

$\mathbf{x}[i..i+k]$  must have the same number of  $a$ 's as  $\mathbf{x}[i+1..i+1+k]$ , and so  $\mathbf{x}[i+1+k] = a$ .  $\square$

### Example

$$\begin{aligned}\mathbf{x}[1..n+1] &= a \underbrace{babbab}a \\ \mathbf{x}[1..n] &= a \underbrace{babbab} \\ \mathbf{x}[2..n+1] &= \underbrace{babbab}a\end{aligned}$$

Note that Lemma 2.7 does not hold if  $\mathbf{u}$  and  $\mathbf{u}'$  overlap less. For example:

### Example

$$\begin{aligned}\mathbf{x}[1..n+1] &= ab \underbrace{ababa}ba \\ \mathbf{x}[1..n-1] &= ab \underbrace{ababa} \\ \mathbf{x}[3..n+1] &= \underbrace{ababa}ba\end{aligned}$$

neither  $\mathbf{x}[1..n-1]$  can be cyclically shifted to the right nor  $\mathbf{x}[3..n+1]$  to the left, yet  $\mathbf{x}[1..n-1]$  and  $\mathbf{x}[3..n+1]$  are conjugates.

**Note:** From the definition below to the end of this chapter, we are going to index from 0 to keep consistent with the publication *The New Periodicity Lemma Revisited* [3].

**Definition 2.8.** Consider a double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$  with nonempty  $\mathbf{u}_2$ . The right inversion subfactors (or RIS) are two substrings

$\mathbf{x}[i..j]$  and  $\mathbf{x}[i + |\mathbf{v}|, j + |\mathbf{v}|]$ , both of length  $|\mathbf{u}_1|$ , where

$$i = (e_1 - 1)|\mathbf{u}_1| + |\mathbf{u}_2| + 1 + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$$

$$j = e_1|\mathbf{u}_1| + |\mathbf{u}_2| + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$$

The left inversion subfactors (or LIS) are two substrings  $\mathbf{x}[i..j]$  and  $\mathbf{x}[i + |\mathbf{v}|, j + |\mathbf{v}|]$ , both of length  $|\mathbf{u}_1|$ , where

$$i = e_1|\mathbf{u}_1| + |\mathbf{u}_2| - 1 - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$$

$$j = (e_1 + 1)|\mathbf{u}_1| + |\mathbf{u}_2| - 2 - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$$

Here are two examples of RIS and LIS :

### Example 1

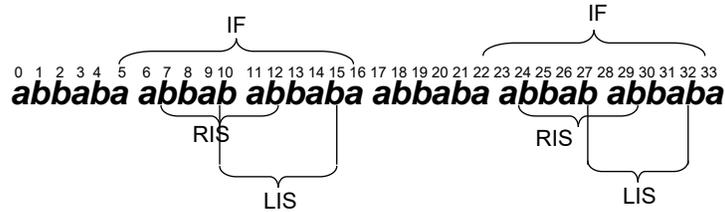
$$\mathbf{u}_2 = \text{abbab}, \bar{\mathbf{u}}_2 = a, \mathbf{u}_1 = \text{abbaba}, e_1 = e_2 = 1$$

$$\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 1, \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 0$$

$$\mathbf{u} = \mathbf{u}_1\mathbf{u}_2 = \text{abbabaabbab}, \mathbf{v} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1 = \text{abbabaabbababbaba}$$

$$\text{IF} = \bar{\mathbf{u}}_2\mathbf{u}_2\bar{\mathbf{u}}_2 = a \text{ abbab abbab } a$$

$$\text{RIS} = \text{bbabab} \text{ and LIS} = \text{babbab}$$



RIS:

$$i = (e_1 - 1)|\mathbf{u}_1| + |\mathbf{u}_2| + 1 + \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 0 + 5 + 1 + 1 = 7 \text{ and}$$



$$15 + 19 = 34$$

LIS:

$$i = e_1|\mathbf{u}_1| + |\mathbf{u}_2| - 1 - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 2 \times 5 + 4 - 1 - 1 = 12 \text{ and}$$

$$i + |\mathbf{v}| = 12 + 19 = 31$$

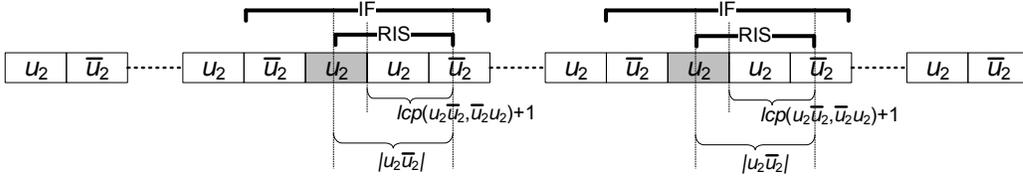
$$j = (e_1 + 1)|\mathbf{u}_1| + |\mathbf{u}_2| - 2 - \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 3 \times 5 + 4 - 2 - 1 = 16$$

$$\text{and } j + |\mathbf{v}| = 16 + 19 = 35$$

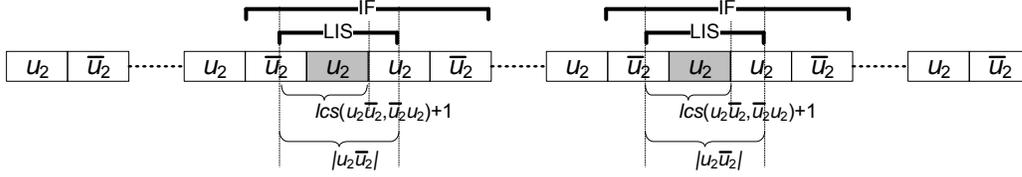
It is useful to view  $\mathbb{RIS}$  in the following way: let  $\mathbf{x}[i..j]$  be the maximum right cyclic shift of the left half of the first  $\mathbb{IF}$ , respectively of the left half of the second  $\mathbb{IF}$ . Then  $\mathbb{RIS}$  is  $\mathbf{x}[i+1, j+1]$ . According to Lemma 2.7, if  $\mathbf{x}[i, j]$  and  $\mathbf{x}[i+1, j+1]$  were conjugate,  $\mathbf{x}[i+1, j+1]$  would have to be a right cyclic shift of  $\mathbf{x}[i, j]$ , which is a contradiction. Thus  $\mathbf{x}[i+1, j+1]$  cannot be conjugate with  $\mathbf{x}[i, j]$  and hence a copy of it can occur neither in  $\mathbf{u}_1^{e_1+e_2}$ , nor in  $\mathbf{u}_1^{e_2}$ .

Similarly for  $\mathbb{LIS}$ : let  $\mathbf{x}[i, j]$  be the maximum left cyclic shift of the right half of the first  $\mathbb{IF}$ , respectively of the right half of the second  $\mathbb{IF}$ . Then  $\mathbb{LIS}$  is  $\mathbf{x}[i-1, j-1]$ . Again, according to Lemma 2.7, if  $\mathbf{x}[i, j]$  and  $\mathbf{x}[i-1, j-1]$  were conjugate,  $\mathbf{x}[i-1, j-1]$  would have to be a left cyclic shift of  $\mathbf{x}[i, j]$ , which is a contradiction. Thus  $\mathbf{x}[i-1, j-1]$  cannot be conjugate with  $\mathbf{x}[i, j]$  and hence a copy of it can occur neither in  $\mathbf{u}_1^{e_1}$ , nor in  $\mathbf{u}_1^{e_1+e_2}$ .

For a better understanding, we illustrate the two natural occurrences of  $\mathbb{RIS}$  in the following diagram:



and the two natural occurrences of LIS as follows:



The purpose of presenting RIS and LIS is to have some structures similar to IF's but shorter, as shorter factors provide a stronger constraint about where the third square can occur in a balanced double square. This will be used in the proof of a reformulated New Periodicity Lemma in Chapter 3.

Let us explain what we mean by “stronger constraint”: when using just IF's to prove where we can have a third square in a balanced double square  $DS(\mathbf{u}, \mathbf{v})$ , the argument follows this line: suppose we have a third square  $\mathbf{w}^2$  such that  $\mathbf{w}_{[1]}$  contains the first IF, then necessarily  $\mathbf{w}_{[2]}$  must contain the second IF (and so  $|\mathbf{w}| = |\mathbf{v}|$ ), or  $\mathbf{w}_{[1]}$  must contain a part of the second IF (and so  $|\mathbf{w}| > |\mathbf{v}|$ ); i.e. it forces the square  $\mathbf{w}^2$  to be big and hence a smaller square cannot exist. However, this argument fails when  $\mathbf{w}_{[1]}$  just contains a part of the first IF. Nevertheless,  $\mathbf{w}_{[1]}$  may still contain the first RIS or LIS and it would force the size of  $\mathbf{w}^2$  in the same way.

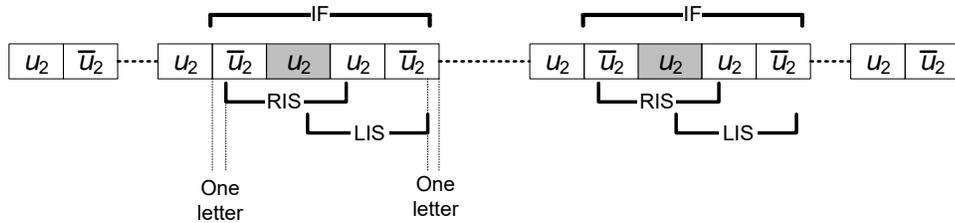
In the last part of this chapter, we will analyze the mutual positions of RIS and LIS by describing the extreme cases when RIS and LIS are almost

disjoint and when they completely overlap. Let  $\mathbf{R}_1$  denote the maximal right cyclic shift of the left half of  $\mathbb{IF}_1$ ,  $\mathbf{R}_2$  denote the maximal right cyclic shift of the left half of  $\mathbb{IF}_2$ , while  $\mathbf{L}_1$  denote the maximal left cyclic shift of the right half of  $\mathbb{IF}_1$ , and  $\mathbf{L}_2$  denote the maximal left cyclic shift of the right half of  $\mathbb{IF}_2$ . We already established that the maximal right cyclic shift of the left half of  $\mathbb{IF}$  is determined by  $\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$  and the maximal left cyclic shift of the right half of  $\mathbb{IF}$  is determined by  $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$ . Lemma 2.5 bounds the size of the shifts as  $\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) \leq |\mathbf{u}_1| - 2$ .

**Note:** Since we get the first pair of RIS and LIS in the same way as the second pair of RIS and LIS, we only discuss how we get the first pair in the following cases.

**Case 1:**  $\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 0$

This is the case when RIS and LIS are almost disjoint. In this case,  $\mathbf{R}_1$  is exactly the left half of  $\mathbb{IF}_1$ , and if we shift  $\mathbf{R}_1$  one position to the right, then we get RIS, and  $\mathbf{L}_1$  is exactly the right half of  $\mathbb{IF}_2$ , and if we shift  $\mathbf{L}_1$  one position to the left then, we get LIS.



**Example 1**

$$\mathbf{u}_2 = abc, \bar{\mathbf{u}}_2 = cb, \mathbf{u}_1 = abccb, e_1 = 2, e_2 = 1$$

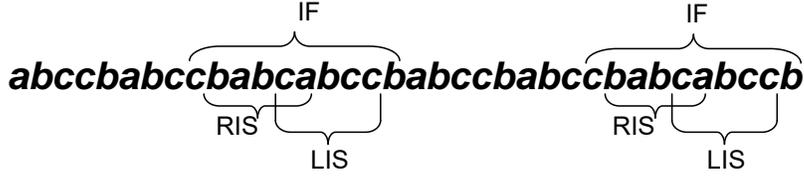
$$\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcp}(abccb, cbabc) = 0$$

$$\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcs}(abccb, cbabc) = 0$$

$$\mathbf{u} = \mathbf{u}_1^2\mathbf{u}_2 = abccbabc, \mathbf{v} = \mathbf{u}_1^2\mathbf{u}_2\mathbf{u}_1 = abccbabcabccb$$

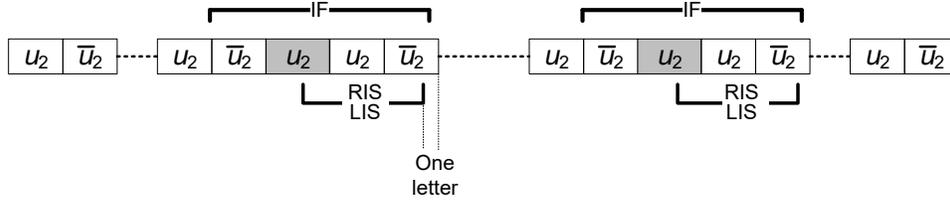
$$\mathbb{IF} = \bar{\mathbf{u}}_2\mathbf{u}_2\bar{\mathbf{u}}_2 = cbabc$$

$$\text{RIS} = babca, \text{LIS} = cabcc$$



$$\text{Case 2: } \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = |\mathbf{u}_1| - 2, \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = 0$$

This is the case when RIS and LIS completely overlap. In this case,  $\mathbf{L}_1$  is exactly the right half of  $\mathbb{IF}_1$ , and if we shift  $\mathbf{L}_1$  one position to the left, then we get LIS, and RIS is at the furthest position to the right as it can be.

**Example 2**

$$\mathbf{u}_2 = ab, \bar{\mathbf{u}}_2 = aba, \mathbf{u}_1 = ababa, e_1 = 2, e_2 = 1$$

$$\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcp}(ababa, abaab) = 3$$

$$\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcs}(ababa, abaab) = 0$$

$$\mathbf{u} = \mathbf{u}_1^2\mathbf{u}_2 = ababaababaab, \mathbf{v} = \mathbf{u}_1^2\mathbf{u}_2\mathbf{u}_1 = ababaababaabababa$$



**Case 4:**  $\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) \neq 0$ ,  $\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) \neq 0$ ,  $\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) + \text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = |\mathbf{u}_1| - 2$

This, too, is a case when RIS and LIS completely overlap.

**Example 4**

$$\mathbf{u}_2 = abb, \bar{\mathbf{u}}_2 = ab, \mathbf{u}_1 = abbab, e_1 = 2, e_2 = 2$$

$$\text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcp}(abbab, ababb) = 2$$

$$\text{lcs}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2) = \text{lcs}(abbab, ababb) = 1$$

$$\mathbf{u} = \mathbf{u}_1^2\mathbf{u}_2 = abbababbababb, \mathbf{v} = \mathbf{u}_1^2\mathbf{u}_2\mathbf{u}_1^2 = abbababbababbababbab$$

$$\text{IF} = \bar{\mathbf{u}}_2\mathbf{u}_2\bar{\mathbf{u}}_2 = ababbabbab$$

$$\text{RIS} = \text{LIS} = bbabb$$

$$\mathbf{u} \quad \underbrace{\hspace{10em}}_{\text{IF}} \quad \underbrace{\hspace{10em}}_{\text{RIS=LIS}} \quad \underbrace{\hspace{10em}}_{\text{IF}} \quad \underbrace{\hspace{10em}}_{\text{RIS=LIS}}$$

## Chapter 3

# Application of canonical factorization to New Periodicity Lemma

In this chapter based on the paper [3], we discuss an application of the canonical factorization of balanced double squares to the New Periodicity Lemma, [9], which shows 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 is complex, breaking down into 14 subcases, and requires very strong assumptions, namely that the shorter of the two squares be regular and the bigger one primitive.

Before proving a strengthened version of Lemma 3.1, we state the original

New Periodicity Lemma:

**Lemma 3.1** ([9], New Periodicity Lemma). *Let a regular  $\mathbf{u}^2$  and a primitively rooted  $\mathbf{v}^2$  be prefixes of a string  $\mathbf{x}$  so that  $|\mathbf{u}| < |\mathbf{v}| < 2|\mathbf{u}|$ . Then for all integers  $k$  and  $w$  such that  $0 \leq k < |\mathbf{v}| - |\mathbf{u}| < w < |\mathbf{v}|$ ,  $w \neq |\mathbf{u}|$ ,  $\mathbf{x}[k+1..k+2w]$  is not a square.*

### Example 1

$\mathbf{u} = aba$ , so  $\mathbf{u}^2$  is regular

$\mathbf{v} = abaab$  is primitive, so  $\mathbf{v}^2$  is primitively rooted.

$|\mathbf{u}| = 3$ ,  $|\mathbf{v}| = 5$ ,  $|\mathbf{v}| - |\mathbf{u}| = 2$

$\mathbf{x} = abaababaab \dots$

$0 \leq k < 2 < w < 5$

e.g. when  $k = 0$  and  $w = 4$ , then  $\mathbf{x}[1..8] = abaa|baba$  is not a square, or

e.g. when  $k = 1$  and  $w = 3$ , then  $\mathbf{x}[2..7] = baa|bab$  is not a square, and

so on for all possible values of  $k$  and all possible values of  $w$ .

### Example 2

$\mathbf{u} = abbaabb$  and so  $\mathbf{u}^2$  is regular

$\mathbf{v} = abbaabbabba$  is primitive and so  $\mathbf{v}^2$  is primitively rooted

$|\mathbf{u}| = 7$ ,  $|\mathbf{v}| = 11$ ,  $|\mathbf{v}| - |\mathbf{u}| = 4$

$\mathbf{x} = abbaabbabbaabbaabbabba \dots$

$0 \leq k < 4 < w < 11$

e.g. when  $k = 0$  and  $w = 5$ , then  $\mathbf{x}[1..10] = abbaa|bbabb$  is not a square,

or e.g. when  $k = 0$  and  $w = 6$ , then  $\mathbf{x}[1..12] = abbaab|babbaa$  is not a

square, or e.g. when  $k = 1$  and  $w = 5$ , then  $\mathbf{x}[2..11] = bbaab|babba$  is not a square, or e.g. when  $k = 1$  and  $w = 8$ , then  $\mathbf{x}[2..17] = bbaabbab|baabbaab$  is not a square, and so on for all possible values of  $k$  and all possible values of  $w$ .

### Example 3

$\mathbf{u} = abcaabc$  and so  $\mathbf{u}^2$  is regular

$\mathbf{v} = abcaabcabca$  is primitive and so  $\mathbf{v}^2$  is primitively rooted

$|u| = 7, |v| = 11, |v| - |u| = 4$

$\mathbf{x} = abcaabcabcaabcaabcabca \dots$

$0 \leq k < 4 < w < 11$

e.g. when  $k = 0$  and  $w = 5$ , then  $\mathbf{x}[1..10] = abcaa|bcabc$  is not a square, or e.g. when  $k = 0$  and  $w = 6$ , then  $\mathbf{x}[1..12] = abcaab|cabcaa$  is not a square, or e.g. when  $k = 1$  and  $w = 5$ , then  $\mathbf{x}[2..11] = bcaab|cabca$  is not a square, or e.g. when  $k = 1$  and  $w = 8$ , then  $\mathbf{x}[2..17] = bcaabcab|caabcaab$  is not a square, and so on for all possible values of  $k$  and all possible values of  $w$ .

In fact, the requirement that  $\mathbf{u}^2$  be regular seems unnecessary. Here we list a few examples that satisfy the conclusion of Lemma 3.1 yet violate the regularity condition:

### Example 4

$\mathbf{u} = abbababbababbababb$ ,  $\mathbf{u}^2$  is not regular as for instance  $(abab)^2$  is a prefix of  $\mathbf{u}$

$\mathbf{v} = abbababbababbababbababbab$  is primitive



e.g. when  $k = 3$  and  $w = 6$ , then  $\mathbf{x}[4..15] = aaabaa|baaaba$  is not a square, or e.g. when  $k = 3$  and  $w = 7$ , then  $\mathbf{x}[4..17] = aaabaab|aaabaaa$  is not a square, and so on for all possible values of  $k$  and all possible values of  $w$ .

The Examples 4, 5, and 6 shows that the regularity of  $\mathbf{u}^2$  is maybe an unnecessary assumption. Since  $\mathbf{u}$  and  $\mathbf{v}$  are proportional squares, they form a balanced double square. Thus, the required assumption of primitiveness of  $\mathbf{v}$  is redundant by Observation 2.3: the fact that  $\mathbf{u}^2$  is regular already forces the primitiveness of  $\mathbf{v}$ . Also note that the regularity of  $\mathbf{u}^2$  necessarily implies that in the canonical factorization of a double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$  the exponents are all equal 1, i.e.  $e_1 = e_2 = 1$ . Let us see what Lemma 3.1 says in the terms of the canonical factorization: so we have a double square  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, 1, 1)$ . As always, 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 gives

$$\mathbf{v}^2 = \underbrace{\overbrace{(\mathbf{u}_2\bar{\mathbf{u}}_2)}^{\mathbf{u}_1} \mathbf{u}_2 (\mathbf{u}_2\bar{\mathbf{u}}_2)}_{\mathbb{IF}_1} \underbrace{(\mathbf{u}_2\bar{\mathbf{u}}_2) \mathbf{u}_2 (\mathbf{u}_2\bar{\mathbf{u}}_2)}_{\mathbb{IF}_2}$$

Thus,  $0 \leq k < |\mathbf{v}| - |\mathbf{u}| < w < |\mathbf{v}|$  gives  $0 \leq k < |\mathbf{u}_1| < w < |\mathbf{v}|$ . Therefore, Lemma 3.1 says that there is no square  $\mathbf{w}^2$  that starts in the first  $\mathbf{u}_1$  and such that  $|\mathbf{u}_1| < |\mathbf{w}| < |\mathbf{v}|$  and  $|\mathbf{w}| \neq |\mathbf{u}|$ .

Now, 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 just from the properties of the canonical factorization. If for instance  $\mathbf{w}$  starts in the first  $\mathbf{u}_2$  and ends in

the fourth  $\mathbf{u}_2$ , then  $\mathbf{w}_{[1]}$  contains fully the  $\mathbb{IF}_1$ , so  $\mathbf{w}_{[2]}$  has to contain  $\mathbb{IF}_2$ , and so  $|\mathbf{w}| \geq |\mathbf{v}|$ , a contradiction. If  $\mathbf{w}$  ends in the second  $\bar{\mathbf{u}}_2$  we cannot argue using  $\mathbb{IF}$ , but still knowing that  $\mathbf{u}_2\bar{\mathbf{u}}_2$  is primitive and also all its rotations are primitive, using the Synchronization principle Lemma 1.31 can be applied to obtain a contradiction.

Almost all possible cases for  $\mathbf{w}^2$  except two can be easily shown impossible using only the properties of the canonical factorization. Thus, it was clear that the canonical factorization could not only provide a significantly simplified proof of Lemma 3.1, but also could allow to significantly reduce the assumption that  $\mathbf{u}$  must be regular. This was the motivation for the research culminating in Theorem 3.2. To increase the power of the use of inversion factors for the use in the proof of Theorem 3.2, the right and left inversion subfactors RIS and LIS were studied, see Chapter 2. The main result in this chapter is the following Theorem 3.2.

**Theorem 3.2.** *Consider a balanced double square  $DS(\mathbf{u}, \mathbf{v})$  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 substring 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 prove the theorem, let us discuss how it relates to the original Lemma 3.1. As mentioned above, for a very specific double square, the

lemma 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 3.2 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 the lemma 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 balanced double squares without any additional conditions or constraints**. Now we can proceed with the proof of Theorem 3.2.

*Proof.* If  $|w| \geq |v|$ , then  $w = v$  since  $w^2$  is a substring of  $v^2$ , and thus case (a) holds. Hence, for the remainder of the proof we can assume that  $|w| < |v|$ . Since double square  $DS(\mathbf{u}, \mathbf{v})$  is balanced, it admits the canonical factorization:  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ , by Lemma 2.1. Then  $\mathbf{u}' = \mathbf{u}_1^{e_2}$  and the primitive root of  $\mathbf{u}'$  is  $\mathbf{u}_1$ .

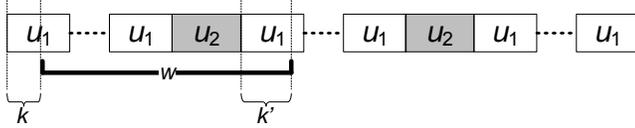
We first deal with the case  $|\mathbf{u}_2| = 0$ . By Lemma 2.1,  $\mathbf{u}_1$  is the primitive root of  $\mathbf{u} = \mathbf{u}_1^{e_1}$  and of  $\mathbf{v} = \mathbf{u}_1^{e_1+e_2}$  with  $e_1 > e_2 \geq 1$ . If case (b) does not hold, we must have  $|\mathbf{w}| \geq |\mathbf{u}| = |\mathbf{u}_1^{e_1}| > |\mathbf{u}_1|$ . Thus  $\mathbf{w}^2$  and  $\mathbf{u}_1^{2e_1+2e_2}$  have a common factor of size  $|\mathbf{u}_1|+|\mathbf{w}|$ , so that by Common factor Lemma 1.32, 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  $|u| \leq |w| < |v|$ .

Let us assume that there is a square  $\mathbf{w}^2$  starting in  $\mathbf{u}_1^{e_1}$  such that  $|\mathbf{w}| > |\mathbf{u}|$ . Since for  $|\mathbf{w}| = |\mathbf{u}|$ ,  $\mathbf{w}$  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  $|\mathbf{w}| > |\mathbf{u}|$ . 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)  $\mathbf{w}_{[1]}$  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  $|\mathbf{w}| > |\mathbf{u}| = |\mathbf{u}_1^{e_1}| + |\mathbf{u}_2|$ ,  $k < k'$ .

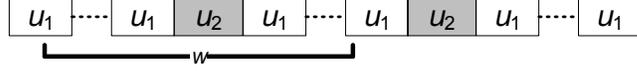
- (i)  $k' \leq \text{lcp}(\mathbf{u}_2\bar{\mathbf{u}}_2, \bar{\mathbf{u}}_2\mathbf{u}_2)$

Then  $\mathbf{w}_{[1]}$  has as a prefix a  $k$ -th rotation of  $\mathbf{u}_1$  and  $\mathbf{w}_{[2]}$  has as a prefix a  $k'$ -th rotation of  $\mathbf{u}_1$ . By the Synchronization principle Lemma 1.31,  $k = k'$ , a contradiction. This case is not possible.

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

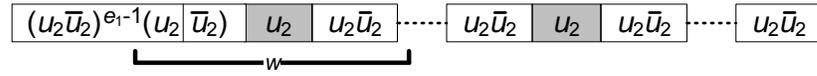
Here  $\mathbf{w}_{[1]}$  contains the first RIS, and so  $\mathbf{w}_{[2]}$  must contain an occurrence of RIS. Since  $\mathbf{w}_{[2]}$  is a factor in  $\mathbf{u}_1^{e_1+e_2}\mathbf{u}_2$ , therefore by Lemma 2.8  $\mathbf{w}_{[2]}$  must contain the second RIS and so  $|\mathbf{w}| \geq |\mathbf{v}|$ , a contradiction. This case is not possible.

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



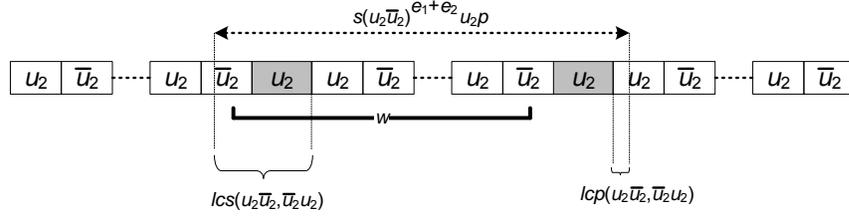
The same argument as in (1)(ii) gives  $|w| \geq |v|$ , a contradiction. This case is not possible.

- (3)  $w_{[1]}$  starts in  $u_1^{e_1-1}u_2$  but not in the first  $u_1$ .



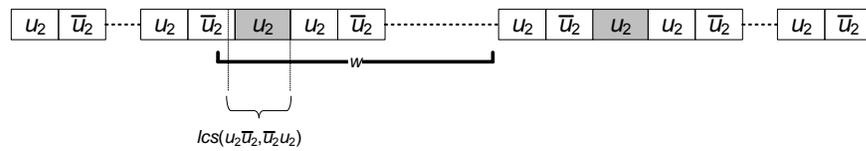
Then  $e_1 > 1$  and since  $|w| > |u| = |u_1^{e_1}u_2|$ ,  $w_{[1]}$  ends past the first  $u_1$  of  $u_1^{e_1+e_2}$ . Therefore,  $w_{[1]}$  contains RIS and so  $|w| \geq |v|$ , i.e. this case is not possible.

- (4)  $w_{[1]}$  starts in the suffix of  $\bar{u}_2u_2$  of  $u$  whose length  $\leq \text{lcs}(u_2\bar{u}_2, \bar{u}_2u_2)$ .



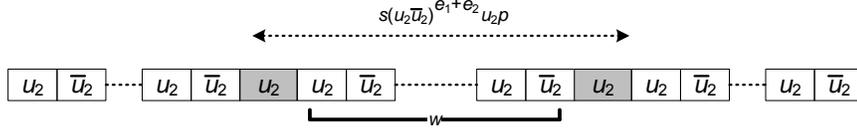
Here  $w_{[1]}$  is a factor in  $su_1^{e_1+e_2}u_2p$ , where  $s$  is the maximal common suffix and  $p$  the maximal common prefix of  $u_2\bar{u}_2$  and  $\bar{u}_2u_2$ . Thus  $w^2$  and  $su_1^{e_1+e_2}u_2p$  have a common factor of size  $|u_1+w|$  and by the Common Factor Lemma 1.32, the primitive root of  $w$  is a conjugate of  $u_1$ , i.e. case (c) holds.

- (5)  $w_{[1]}$  starts in the suffix of  $\bar{u}_2u_2$  of  $u$  whose length  $> \text{lcs}(u_2\bar{u}_2, \bar{u}_2u_2)$ .



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

(6)  $\mathbf{w}_{[1]}$  starts past the first  $\mathbf{u}$ .



The same argument as in (4) shows that the primitive root of  $\mathbf{w}$  is a conjugate of  $\mathbf{u}_1$  and so case (c) holds.

□

In the following we present a few examples of the Theorem 3.2 illustrating various situations encountered in the proof.

### Example 1

This is an example of a square  $\mathbf{w}^2$  such that  $|\mathbf{w}| < |\mathbf{u}_1|$

$$\mathbf{u}_2 = abaa, \bar{\mathbf{u}}_2 = baa, \mathbf{u}_1 = abaabaa, e_1 = 1, e_2 = 1$$

$$\mathbf{u} = \mathbf{u}_1 \mathbf{u}_2 = abaabaaabaa$$

$$\mathbf{v} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 = abaabaaabaabaa$$

$$\mathbf{w}^2 = abaaba$$

$$\mathbf{v}^2 = \underbrace{abaaba}_{\mathbf{w}^2} aabaaabaabaaabaabaaabaabaa$$

### Example 2

This is an example of a square  $\mathbf{w}^2$  such that  $|\mathbf{w}| < |\mathbf{u}|$

$$\mathbf{u}_2 = abb, \bar{\mathbf{u}}_2 = ab, \mathbf{u}_1 = abbab, e_1 = 2, e_2 = 1$$

$$\mathbf{u} = \mathbf{u}_1^2 \mathbf{u}_2 = abbababbababb$$

$$\mathbf{v} = \mathbf{u}_1^2 \mathbf{u}_2 \mathbf{u}_1 = \text{abbababbababbabbab}$$

$$\mathbf{w}^2 = \text{abbababbab}$$

$$\mathbf{v}^2 = \underbrace{\text{abbababbab}}_{\mathbf{w}^2} \text{abbabbababbababbababbabbab}$$

**Example 3**

This is an example of a square  $\mathbf{w}^2$  such that  $|\mathbf{w}| < |\mathbf{u}|$

$$\mathbf{u}_2 = \text{abb}, \bar{\mathbf{u}}_2 = \text{ab}, \mathbf{u}_1 = \text{abbab}, e_1 = 2, e_2 = 2$$

$$\mathbf{u} = \mathbf{u}_1^2 \mathbf{u}_2 = \text{abbababbababb}$$

$$\mathbf{v} = \mathbf{u}_1^2 \mathbf{u}_2 \mathbf{u}_1^2 = \text{abbababbababbabbababbab}$$

$$\mathbf{w}^2 = \text{babbababbababbababba}$$

$$\mathbf{v}^2 = \text{abbababbabab} \underbrace{\text{babbababba babbababba}}_{\mathbf{w}^2} \text{babbabbababbab}$$

**Example 4**

This is an example of a square  $\mathbf{w}^2$  such that  $|\mathbf{w}| < |\mathbf{u}|$

$$\mathbf{u}_2 = \text{abaa}, \bar{\mathbf{u}}_2 = \text{a}, \mathbf{u}_1 = \text{abaaa}, e_1 = 2, e_2 = 1$$

$$\mathbf{u} = \mathbf{u}_1^2 \mathbf{u}_2 = \text{abaaaabaaaabaa}$$

$$\mathbf{v} = \mathbf{u}_1^2 \mathbf{u}_2 \mathbf{u}_1 = \text{abaaaabaaaabaaabaaa}$$

$$\mathbf{w}^2 = \text{abaaaabaaaabaaaabaa}$$

$$\mathbf{v}^2 = \text{abaaaabaaaaba} \underbrace{\text{abaaaabaa abaaaabaa}}_{\mathbf{w}^2} \text{abaaa}$$

**Example 5**

This is an example of a square  $\mathbf{w}^2$  such that  $|\mathbf{w}| < |\mathbf{u}|$

$$\mathbf{u}_2 = \text{ab}, \bar{\mathbf{u}}_2 = \text{b}, \mathbf{u}_1 = \text{abb}, e_1 = 2, e_2 = 2$$

$$\mathbf{u} = \mathbf{u}_1^2 \mathbf{u}_2 = \text{abbabbab}$$

$$\mathbf{v} = \mathbf{u}_1^2 \mathbf{u}_2 \mathbf{u}_1^2 = \text{abbabbababbabb}$$

$$\mathbf{w}^2 = \text{bbabbabbabba}$$

$$\mathbf{v}^2 = \text{abbabbaba} \underbrace{\text{bbabbabbabba}}_{\mathbf{w}^2} \text{babbabb}$$

### Example 6

This is an example of a square  $\mathbf{w}^2$  such that  $|\mathbf{w}| < |\mathbf{u}|$

$$\mathbf{u}_2 = \text{abc}, \bar{\mathbf{u}}_2 = \text{bc}, \mathbf{u}_1 = \text{abcbc}, e_1 = 3, e_2 = 1$$

$$\mathbf{u} = \mathbf{u}_1^3 \mathbf{u}_2 = \text{abc bcaabc bcaabc bcaabc}$$

$$\mathbf{v} = \mathbf{u}_1^3 \mathbf{u}_2 \mathbf{u}_1 = \text{abc bcaabc bcaabc bcaabc bcaabc}$$

$$\mathbf{w}^2 = \text{bcbcaabc bcaabc bcaabc bcaa}$$

$$\mathbf{v}^2 = \text{abc bcaabc bcaabc bcaabc bcaa} \underbrace{\text{bcbcaabc bcaa bcbcaabc bcaa}}_{\mathbf{w}^2} \text{bcaabc bc}$$

## Chapter 4

# Application of canonical factorization to Three Squares Lemma

In this chapter based on the paper [1], we discuss an application of the canonical factorization of double squares to the Three Squares Lemma by Crochemore and Rytter introduced in 1995. It is a structural lemma on three squares starting at the same position. This influential lemma has been used by many researchers in the field of periodicities in strings. In particular, Fraenkel and Simpson used it in 1998 to obtain a universal upper bound for the maximum number of distinct squares occurring in a string. We present a generalization of Three Squares Lemma by exploiting the canonical factorization of balanced double squares discussed in Chapter 2. Let us first state

the original lemma.

**Lemma 4.1** ([7], Three Squares Lemma). *Let  $\mathbf{u}^2 \neq \mathbf{v}^2$  be proper prefixes of  $\mathbf{w}^2$  and let  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  be primitive, then  $|\mathbf{u}| + |\mathbf{v}| < |\mathbf{w}|$ .*

Lemma 4.1 has been used by many researchers including Kolpakov and Kucherov [13], Stoye and Gusfield [18], Fan, Puglisi, Smyth, and Turpin [9], Simpson [16]. Lemma 4.1 was essential for the 1998 result by Fraenkel and Simpson [10] giving a universal upper bound of  $2n$  for the number of distinct squares in a string of length  $n$ . Note that for the problem of distinct squares, every type of square is only counted once, i.e. the types, rather than the occurrences, are counted. For illustration, *aabaab* contains the following three underlined squares aabaab, aabaab and aabaab while the number of distinct squares is 2: *aa* and *aabaab*. Ilie [12] provided in 2005 an alternate proof of the main theorem of [10] not directly using Lemma 4.1. Noticing that the proof of Lemma 4.1 by Crochemore and Rytter only requires the primitiveness of the shortest square, Fraenkel and Simpson [10] proposed the following strengthening referred to as three-prefix-square Lemma in [6] where additional context and references can be found.

**Lemma 4.2** ([10], Three Squares Lemma – Fraenkel and Simpson variant). *Let  $\mathbf{u}^2 \neq \mathbf{v}^2$  be proper prefixes of  $\mathbf{w}^2$  and let the shorter of the two strings  $\mathbf{u}$  and  $\mathbf{v}$  be primitive, then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$ .*

**Example 1**

$$\mathbf{w}^2 = aabaaababaabaaabab$$

$$\mathbf{w} = aabaaabab, \mathbf{v} = aaba, \mathbf{u} = a$$

$$|\mathbf{w}| = 9, |\mathbf{v}| = 4, |\mathbf{u}| = 1, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

**Example 2**

$$\mathbf{w}^2 = abbabbaabbabbaaabbabbaabbabbaa$$

$$\mathbf{w} = abbabbaabbabbaa, \mathbf{v} = abbabba, \mathbf{u} = abb$$

$$|\mathbf{w}| = 15, |\mathbf{v}| = 7, |\mathbf{u}| = 3, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

**Example 3**

$$\mathbf{w}^2 = abaabaabaaba$$

$$\mathbf{w} = abaabaabaaba, \mathbf{v} = abaaba, \mathbf{u} = aba$$

$$|\mathbf{w}| = 12, |\mathbf{v}| = 6, |\mathbf{u}| = 3, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

**Example 4**

$$\mathbf{w}^2 = abababababab$$

$$\mathbf{w} = ababab, \mathbf{v} = abab, \mathbf{u} = ab$$

$$|\mathbf{w}| = 6, |\mathbf{v}| = 4, |\mathbf{u}| = 2, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| = |\mathbf{w}|$$

Fraenkel and Simpson illustrated the necessity of the primitiveness for the shortest square with the following example:  $\mathbf{u} = a^2$ ,  $\mathbf{v} = a^4$ , and  $\mathbf{w} = a^5$ . We present a further strengthening based on the recently investigated structural properties of two squares starting at the same position, see [2, 8].

**Fraenkel and Simpson counterexample**

$$\mathbf{w}^2 = aaaaaaaaaa$$

$$\mathbf{w} = aaaaa, \mathbf{v} = aaaa, \mathbf{u} = aa$$

$$|\mathbf{w}| = 5, |\mathbf{v}| = 4, |\mathbf{u}| = 2, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| > |\mathbf{w}|$$

In the paper [1], our strengthened version of the Three Square Lemma was formulated in the following way

**Lemma 4.3.** *Let  $\mathbf{u}^2 \neq \mathbf{v}^2$  be proper prefixes of  $\mathbf{w}^2$ , then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$  unless  $\mathbf{u}, \mathbf{v}$ , and  $\mathbf{w}$  have the same primitive root.*

After the publication of [1], a fellow AdvOL Ph.D. student Adrien Thierry came with a counterexample for the lemma:

$\mathbf{u} = (ab)^2$ ,  $\mathbf{v} = (ab)^3a$ , and  $\mathbf{w} = (ab)^3a(ab)$ , clearly  $|\mathbf{u}| + |\mathbf{v}| = 4 + 7 = 11$  while  $|\mathbf{w}| = 9$ , yet  $\mathbf{w}\mathbf{w} = \overbrace{abab\ \mathbf{u}\ \abaabababab}^{\mathbf{v}\mathbf{v}}baab$  and  $\mathbf{u}, \mathbf{v}$ , and  $\mathbf{w}$  do not have the same primitive root.

Thierry's counterexample pointed out a little oversight in the proof of Lemma 4.3. The proof was corrected and it weakened the conclusion of the lemma slightly. Here we present the corrected Lemma 4.4 and its proof. The corrigendum was published soon after the error was detected, see [1].

**Lemma 4.4.** *Let  $\mathbf{u}^2 \neq \mathbf{v}^2$  be proper prefixes of  $\mathbf{w}^2$ , then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$  unless  $\mathbf{u} = \mathbf{v}_1^t$ ,  $\mathbf{v} = \mathbf{v}_1^{p_1}\mathbf{v}_2$ , and  $\mathbf{w} = \mathbf{v}_1^{p_1}\mathbf{v}_2\mathbf{v}_1^{p_2}$  where  $\mathbf{v}_1$  is primitive,  $\mathbf{v}_2$  a proper possibly empty prefix of  $\mathbf{v}_1$ ,  $t > p_2$ , and  $p_1 \geq p_2 \geq 1$ .*

Lemma 4.4 shows that the squares  $\mathbf{u}^2$ ,  $\mathbf{v}^2$ , and  $\mathbf{w}^2$  violating  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$  consist of two types; one corresponding to the counterexample given by Fraenkel and Simpson (when  $\mathbf{v}_2 = \varepsilon$ ) and one corresponding to the coun-

terexample of Thierry (when  $\mathbf{v}_2 \neq \varepsilon$ ). Bellow is essentially Fraenkel-Simpsons counterexample:

$$\mathbf{w}^2 = abababababababab$$

$$\mathbf{w} = abababab, \mathbf{v} = ababab, \mathbf{u} = abab,$$

$$|\mathbf{w}| = 8, |\mathbf{v}| = 6, |\mathbf{u}| = 4, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| > |\mathbf{w}|$$

Corollary 4.5 illustrates that Lemma 4.4 is a true generalization of Lemma 4.2.

**Corollary 4.5.** *Let  $\mathbf{u}^2$  be a proper prefix of  $\mathbf{v}^2$  that is a proper prefixes of  $\mathbf{w}^2$  and let  $\mathbf{u}$  be primitive, then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$ . Moreover, if  $|\mathbf{u}| < |\mathbf{v}| < 2|\mathbf{u}|$  and either  $\mathbf{v}$  or  $\mathbf{w}$  is primitive, then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$ .*

*Proof.* Let us assume by contradiction that  $|\mathbf{u}| + |\mathbf{v}| > |\mathbf{w}|$ . Then by Lemma 4.4,  $\mathbf{u} = \mathbf{v}_1^t$ ,  $\mathbf{v} = \mathbf{v}_1^{p_1}\mathbf{v}_2$  and  $\mathbf{w} = \mathbf{v}_1^{p_1}\mathbf{v}_2\mathbf{v}_1^{p_2}$  for a primitive  $\mathbf{v}_1$ , a proper possibly empty prefix  $\mathbf{v}_2$  of  $\mathbf{v}_1$ , and  $t > p_2$ ,  $p_1 \geq p_2 \geq 1$ . If  $\mathbf{u}$  is primitive,  $t = 1$  and so  $t > p_2 \geq 1$  is a contradiction. If  $|\mathbf{v}| < 2|\mathbf{u}|$ , then  $\mathbf{v}_1^{p_1}\mathbf{v}_2$  is a prefix of  $\mathbf{v}_1^{2t}$ , which can only be true when  $\mathbf{v}_2$  is empty due to Synchronization Principle Lemma 1.31. If  $\mathbf{v}$  is primitive, then  $p_1 = 1$  and so  $p_2 = 1$  and so  $\mathbf{u} = \mathbf{v}_1^t$ ,  $t > 1$  and  $\mathbf{v} = \mathbf{v}_1$  and  $\mathbf{w} = \mathbf{v}_1^2$ , and so  $|\mathbf{u}| \geq |\mathbf{w}|$ , a contradiction. If  $\mathbf{w}$  is primitive, then  $\mathbf{w} = \mathbf{v}_1$ , and so  $|\mathbf{w}| = |\mathbf{v}|$ , a contradiction.  $\square$

**Example 1 : when  $\mathbf{u}$  is primitive**

$$\mathbf{w}^2 = abababababababab$$

$$\mathbf{w} = abababab, \mathbf{v} = abab, \mathbf{u} = ab$$

$$|\mathbf{w}| = 8, |\mathbf{v}| = 4, |\mathbf{u}| = 2, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

**Example 2 : when  $|\mathbf{u}| < |\mathbf{v}| < 2|\mathbf{u}|$  and  $\mathbf{v}$  is primitive**

$$\mathbf{w}^2 = aabaabaabaab$$

$$\mathbf{w} = aabaab, \mathbf{v} = aab, \mathbf{u} = aa$$

$$|\mathbf{w}| = 6, |\mathbf{v}| = 3, |\mathbf{u}| = 2, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

**Example 3 : when  $|\mathbf{u}| < |\mathbf{v}| < 2|\mathbf{u}|$  and  $\mathbf{w}$  is primitive**

$$\mathbf{w}^2 = aaaaaabaaaaaab$$

$$\mathbf{w} = aaaaaab, \mathbf{v} = aaa, \mathbf{u} = aa$$

$$|\mathbf{w}| = 7, |\mathbf{v}| = 3, |\mathbf{u}| = 2, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

**Example 4 : when  $\mathbf{v}$  is primitive, this example shows that conditions in Corollary 4.5 are not necessary to conclude that  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$ , since  $\mathbf{u}$  is not primitive and  $|\mathbf{v}| > 2|\mathbf{u}|$ .**

$$\mathbf{w}^2 = aaaabaaaabaaaabaaaab$$

$$\mathbf{w} = aaaabaaaab, \mathbf{v} = aaaab, \mathbf{u} = aa$$

$$|\mathbf{w}| = 10, |\mathbf{v}| = 5, |\mathbf{u}| = 2, \text{ thus, } |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$$

#### Proof of Lemma 4.4

*Proof.* Let  $\mathbf{u} \neq \mathbf{v}$ , and  $\mathbf{u}^2$  and  $\mathbf{v}^2$  be both proper prefixes of  $\mathbf{w}^2$ . Lemma 4.4 states that

$$\{ \mathbf{u} = \mathbf{v}_1^t, \mathbf{v} = \mathbf{v}_1^{p_1} \mathbf{v}_2, \mathbf{w} = \mathbf{v}_1^{p_1} \mathbf{v}_2 \mathbf{v}_1^{p_2}, t > p_2, p_1 \geq p_2 \geq 1 \} \text{ or } \{ |\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}| \}. \quad (S)$$

Without loss of generality, we can assume that  $|\mathbf{u}| < |\mathbf{v}|$ .

If  $2|\mathbf{v}| \leq |\mathbf{w}|$ , then  $|\mathbf{v}| + |\mathbf{u}| < |\mathbf{w}|$  as  $|\mathbf{u}| < |\mathbf{v}|$ , and thus (S) holds. Therefore,

we can assume that  $|\mathbf{w}| < 2|\mathbf{v}|$ ; that is,  $(\mathbf{v}, \mathbf{w})$  is a balanced double square and thus admits a canonical factorization  $(\mathbf{v}, \mathbf{w} : \mathbf{v}_1, \mathbf{v}_2, p_1, p_2)$  by Lemma 2.1. We consider the following cases.

1. Case when  $\mathbf{u}$  and  $\mathbf{v}$  are not proportional, i.e.  $2|\mathbf{u}| \leq |\mathbf{v}|$ .

If  $|\mathbf{u}| < |\mathbf{v}_1|$ , then

$$\begin{aligned}
 |\mathbf{u}| + |\mathbf{v}| &= |\mathbf{u}| + |\mathbf{v}_1^{p_1} \mathbf{v}_2| \\
 &< |\mathbf{v}_1| + |\mathbf{v}_1^{p_1} \mathbf{v}_2| \\
 &= |\mathbf{v}_1^{p_1+1} \mathbf{v}_2| \\
 &\leq |\mathbf{v}_1^{p_1+p_2} \mathbf{v}_2| \\
 &= |\mathbf{w}|
 \end{aligned}$$

If  $|\mathbf{u}| \geq |\mathbf{v}_1|$ , since  $\mathbf{u}^2$  is a prefix of  $\mathbf{v} = \mathbf{v}_1^{p_1} \mathbf{v}_2$ , then  $\mathbf{u}^2$  and  $\mathbf{v}_1^{p_1} \mathbf{v}_2$  have a common factor of length  $|\mathbf{u}| + |\mathbf{v}_1|$ , and by Lemma 1.32,  $\mathbf{u}$  and  $\mathbf{v}_1$  have the same primitive root, and so  $\mathbf{v}_1$  is the primitive root of  $\mathbf{v}_1$ .

Thus  $\mathbf{u} = \mathbf{v}_1^t$  for some  $t \geq 1$ ,  $\mathbf{v} = \mathbf{v}_1^{p_1} \mathbf{v}_2$ , and  $\mathbf{w} = \mathbf{v}_1^{p_1} \mathbf{v}_2 \mathbf{v}_1^{p_2}$ .

Now we check the relation between  $|\mathbf{v}_1^t + \mathbf{v}_1^{p_1} \mathbf{v}_2|$  and  $|\mathbf{v}_1^{p_1} \mathbf{v}_2 \mathbf{v}_1^{p_2}|$ , which actually is the relation between  $|\mathbf{v}_1^t|$  and  $|\mathbf{v}_1^{p_2}|$ .

If  $t \leq p_2$ , then  $|\mathbf{u}| + |\mathbf{v}| \leq |\mathbf{w}|$ , but if  $t > p_2$ , then  $|\mathbf{u}| + |\mathbf{v}| > |\mathbf{w}|$ . In both cases,  $(S)$  holds.

2. Case when  $\mathbf{u}$  and  $\mathbf{v}$  are proportional, i.e.  $|\mathbf{v}| < 2|\mathbf{u}|$ .

Then  $DS(\mathbf{u}, \mathbf{v})$  is a balanced double square and thus admits by the Two Square Factorization Lemma 2.1 a canonical factorization  $DS(\mathbf{u}, \mathbf{v} : \mathbf{u}_1, \mathbf{u}_2, e_1, e_2)$ . As a result, we have:

$$\mathbf{u} = \mathbf{u}_1^{e_1} \mathbf{u}_2$$

$$\mathbf{v} = \mathbf{u}_1^{e_1} \mathbf{u}_2 \mathbf{u}_1^{e_2}$$

$$\mathbf{v}^2 = \mathbf{u}_1^{e_1} \mathbf{u}_2 \mathbf{u}_1^{e_1+e_2} \mathbf{u}_2 \mathbf{u}_1^{e_2}$$

and

$$\mathbf{v} = \mathbf{v}_1^{p_1} \mathbf{v}_2$$

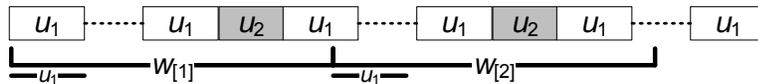
$$\mathbf{w} = \mathbf{v}_1^{p_1} \mathbf{v}_2 \mathbf{v}_1^{p_2}$$

$$\mathbf{w}^2 = \mathbf{v}_1^{p_1} \mathbf{v}_2 \mathbf{v}_1^{p_1+p_2} \mathbf{v}_2 \mathbf{v}_1^{p_2}$$

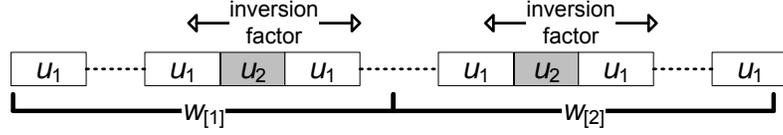
- (i) Case when  $|\mathbf{u}_2| = 0$ . Then  $e_1 > e_2$ ,  $\mathbf{u} = \mathbf{u}_1^{e_1}$ , and  $\mathbf{v} = \mathbf{u}_1^{e_1+e_2}$ .

Let us assume that  $|\mathbf{w}| < |\mathbf{u}| + |\mathbf{v}| = (2e_1 + e_2)|\mathbf{u}_1|$ . Then  $\mathbf{w}^2$  and  $\mathbf{u}_1^{2e_1+2e_2}$  have a common factor of length  $|\mathbf{w}| + |\mathbf{u}_1|$ , and by the Common Factor Lemma 1.32 the primitive root of  $\mathbf{w}$  is a conjugate of  $\mathbf{u}_1$ , i.e. equals  $\mathbf{u}_1$ . Thus,  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  all have the same primitive root, and thus (S) holds.

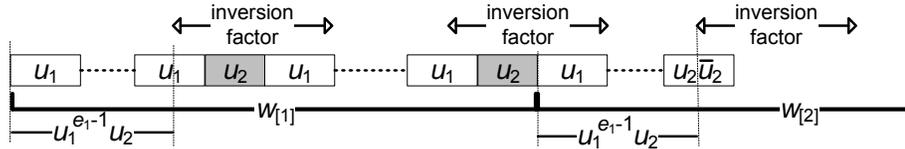
- (ii) Case when  $|\mathbf{u}_2| > 0$ . Let  $\mathbf{w}_{[1]}$  refer to the first occurrence of  $\mathbf{w}$  and  $\mathbf{w}_{[2]}$  to the second. First, we have to show that  $\mathbf{w}_{[1]}$  does not end in the first  $\mathbf{u}_1$  of  $\mathbf{u}_1^{e_1+e_2}$ . If it did, then it would contradict the Synchronization Principle Lemma 1.31 as  $\mathbf{w}_{[1]}$  and hence  $\mathbf{w}_{[2]}$  has the primitive  $\mathbf{u}_1$  as a prefix as indicated by the following diagram:



Thus,  $w_{[1]}$  must end somewhere past the first  $u_1$  of  $u_1^{e_1+e_2}$ :

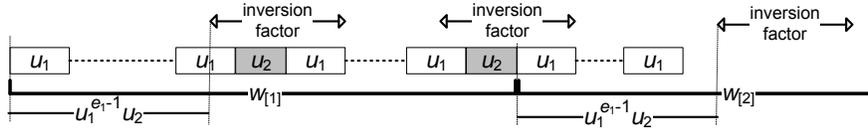


As a consequence,  $ws_{[1]}$  contains the first inversion factor of  $u_1^{e_1}u_2u_1^{e_1+e_2}u_2u_1^{e_2}$  exactly at a distance of  $|u_1^{e_1-1}u_2|$  from the beginning. It follows that  $w_{[2]}$  must contain an occurrence of the inversion factor at exactly the same distance from the beginning of  $w_{[2]}$ . If it were the second inversion factor, then the length of  $w$  would be exactly  $|v|$  as it is the distance between the two occurrences of the inversion factor in  $w^2$ , a contradiction. Thus, it must be an occurrence of the inversion factor past the second one. The first possible start of another occurrence of the inversion factor is the suffix  $\bar{u}_2$  of  $u_1^{e_1}u_2u_1^{e_1+e_2}u_2u_1^{e_2}$ . If  $e_1 = e_2$ , then it is the case that  $w_{[1]} = w_{[2]} = u_1^{e_1-1}(\bar{u}_2u_2u_2\bar{u}_2)u_1^{e_1+e_2-1}u_2$  (see below) and then

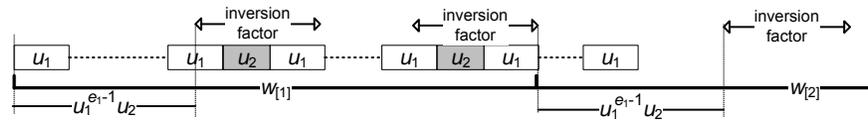


$|w| = |u| + |v|$  as  $|u| = |u_1^{e_1}u_2|$  and  $|v| = |u_1^{e_1+e_2}u_2|$ , thus (S) holds. If  $e_1 > e_2$ , then the occurrence of the inversion factor in  $w_{[2]}$  must be at a distance  $|u_1^{e_1-1}u_2|$  from the beginning of  $w_{[2]}$ . By the Synchronization Principle Lemma 1.31, the prefix  $u_1^{e_1-1}u_2$  of

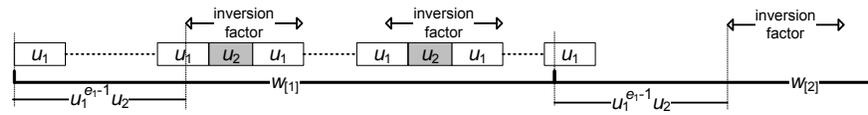
$w_{[2]}$  must align with  $u_1^{e_2}$  or start in the last  $u_1$  of  $u_1^{e_2}$ , and so  $w_{[1]}$  must have  $u_1^{e_1} u_2 u_1^{e_1+e_2} u_2$  as a prefix, again yielding  $|w| \geq |u|+|v|$  (see below), thus  $(S)$  holds.



OR



OR



This completes the proof. □

# Chapter 5

## Conclusion and Future work

We presented a unique canonical factorization of a balanced double square consisting of two squares starting in the same position and of comparable lengths. The interesting part is that this unique factorization is guaranteed without any additional conditions. We described and analyzed three kinds of factors occurring rarely in the longer square: inversion factors, right inversion subfactors and left inversion subfactors. We discussed how such rare factors can be used to constraint the sizes of any possible third square starting in proximity of the start of the double square. The right and left inversion subfactors provide stronger constraint in comparison to the inversion factors due to their length that is half of the length of the inversion factors. The right and left inversion subfactors were discovered and analyzed exactly for this stronger constraint for the purpose of proving the New Periodicity Lemma.

A stronger version of the New Periodicity Lemma was formulated with

a significantly weaker assumption using the canonical factorization and the right and left inversion subfactors. The range of possible starts for "forbidden" squares was enlarged without any additional constraints. We also gave a much simpler proof of the generalized New Periodicity Lemma in comparison to the proof of the original New Periodicity Lemma.

The canonical factorization and the inversion factors were applied to prove a stronger version of the Three Squares Lemma. In order to guarantee the conclusion of the original Three Squares Lemma, our generalization only needs one of the three squares to be primitive-rooted when the constraints of length are satisfied, while the original lemma requires the shortest square's root to be primitive.

The future work and research directions will focus on strengthening the ability of the rare factors to constrain the possible positions where the third square can occur in a balanced double square. This will require discovering some shorter factors that occur rarely in the double square. With shorter rare factors, both the New Periodicity Lemma and the Three Squares Lemma could be generalized and extended further. Since the rare factors were essential to proving the best-to-date upper bound of  $11n/6$  for the number of distinct squares in a string of length  $n$ , it is quite promising that shorter rare factors will further improve this upper bound.

# Bibliography

- [1] H. Bai, A. Deza, and F. Franek. On a lemma of Crochemore and Rytter. *Journal of Discrete Algorithms*, 34:18–22, 2015. (Corrigendum, *Journal of Discrete Algorithms*, 38-41:50-51, 2016).
- [2] H. Bai, F. Franek, and W. Smyth. Two squares canonical factorization. In *Proceedings of Prague Stringology Conference PSC 2014*, pages 52–58, 2014.
- [3] H. Bai, F. Franek, and W. Smyth. The new periodicity lemma revisited. *Discrete Applied Mathematics*, 212(30):30–36, 2016.
- [4] W. Bland and W. Smyth. Three overlapping squares. *Theor. Comput. Sci.*, 596(C):23–40, 2015.
- [5] K. Chen, R. Fox, and R. Lyndon. Free differential calculus, IV. The quotient groups of the lower central series. *The Annals of Mathematics*, 68(1):81–95, 1958.

- [6] M. Crochemore, C. Hancart, and T. Lecroq. *Algorithms on Strings*. Cambridge University Press, New York City, New York, 2007.
- [7] M. Crochemore and W. Rytter. Squares, cubes, and time-space efficient string searching. *Algorithmica*, 13(5):405–425, 1995.
- [8] A. Deza, F. Franek, and A. Thierry. How many double squares can a string contain? *Discrete Applied Mathematics*, 180(10):52–69, 2015.
- [9] K. Fan, S. Puglisi, W. Smyth, and A. Turpin. A new periodicity lemma. *SIAM Journal on Discrete Mathematics*, 20(3):656–668, 2006.
- [10] A. Fraenkel and J. Simpson. How many squares can a string contain? *Journal of Combinatorial Theory, Series A*, 82(1):112–120, 1998.
- [11] F. Franek, R. Fuller, J. Simpson, and W. Smyth. More results on overlapping squares. *Journal of Discrete Algorithms*, 17:2–8, 2012.
- [12] L. Ilie. A simple proof that a word of length  $n$  has at most  $2n$  distinct squares. *Journal of Combinatorial Theory, Series A*, 112(1):163–164, 2005.
- [13] R. Kolpakov and G. Kucherov. Finding maximal repetitions in a word in linear time. In *Proceedings of 40th Annual Symposium on Foundations of Computer Science*, pages 596–604, 1999.
- [14] E. Kopylova and W. Smyth. The three squares lemma revisited. *Journal of Discrete Algorithms*, 11:3–14, 2012.

- [15] N. Lam. On the number of squares in a string. Technical report, AdvOL-Report No. 2013/02, McMaster University, 2013.
- [16] J. Simpson. Intersecting periodic words. *Theoretical Computer Science*, 374(20):58–65, 2007.
- [17] W. Smyth. *Computing Patterns in Strings*. Addison-Wesley, 2002.
- [18] J. Stoye and D. Gusfield. Simple and flexible detection of contiguous repeats using a suffix tree. *Theoretical Computer Science*, 270(6):843–856, 2002.
- [19] A. Thierry. Combinatorics of the interrupted period. In *Proceedings of Prague Stringology Conference PSC 2015*, pages 17–21, 2015.