

# Completing the New Periodicity Lemma

Widmer Bland

A thesis submitted to the School of Graduate Studies in partial fulfilment of the requirements for the degree of Master of Science

Master of Science  
Computer Science  
2014

Department of Computing & Software  
McMaster University  
Hamilton, Ontario, Canada

Title:                    Completing the New Periodicity Lemma

Author:                 Widmer Bland

Supervisor:            Dr W.F. Smyth

Number of pages:    iv, 45

### Abstract

The “Three Squares Lemma” (Crochemore and Rytter 1995) famously explored the consequences of supposing that three squares occur at the same position in a string. Essentially, it showed that this phenomenon could not occur unless the longest of the three squares was at least the sum of the lengths of the other two. More recently, several papers (Fan et al. 2006; Franek, Fuller, et al. 2012; Kopylova and Smyth 2012; Simpson 2007) have greatly extended this result to a “New Periodicity Lemma” (NPL) by supposing that only two of the squares occur at the same position, with a third occurring in a neighbourhood to the right. The proof of the NPL involves fourteen subcases, twelve of which have been proven over the last seven years. In this thesis, we prove the final two remaining.

### **Acknowledgements**

Bill Smyth imparted his enthusiasm for strings, and provided direction and encouragement that enabled my work. I feel greatly fortunate to have had his mentorship.

Michael Soltys offered helpful comments on this thesis and was kindly supportive on many occasions over the last few years.

Thanks.

# Contents

|   |           |
|---|-----------|
| List of Figures   | iv        |
| <b>1 Introduction</b>                                       | <b>1</b>  |
| <b>2 Background</b>   | <b>3</b>  |
| 2.1 Strings . . . . .                                       | 3         |
| 2.2 String data structures . . . . .                        | 3         |
| 2.3 Periodicity . . . . .                                   | 6         |
| 2.4 Repetitions and runs . . . . .                          | 8         |
| 2.4.1 Algorithms for finding repetitions and runs . . . . . | 9         |
| 2.4.2 The combinatorics of runs . . . . .                   | 10        |
| 2.5 Three overlapping squares . . . . .                     | 11        |
| 2.5.1 The New Periodicity Lemma . . . . .                   | 11        |
| 2.5.2 The general case characterized . . . . .              | 14        |
| <b>3 Completing the New Periodicity Lemma</b>               | <b>16</b> |
| 3.1 Subcase 3 . . . . .                                     | 16        |
| 3.2 Subcase 7 . . . . .                                     | 26        |
| <b>4 On to the General Case</b>                             | <b>38</b> |
| <b>Bibliography</b>   | <b>40</b> |

# List of Figures

|      |   |    |
|------|---|----|
| 2.1  | $ST_x$ for $x = abaababa$ . . . . .   | 4  |
| 2.2  | $SA_x$ and $LCP_x$ for $x = abaababa$ . . . . .                                     | 5  |
| 2.3  | Three overlapping squares, as postulated in Lemma 13. . . . .                       | 12 |
| 2.4  | The 14 subcases . . . . .   | 13 |
| 2.5  | Structure of $x$ for the 14 subcases . . . . .                                      | 14 |
| 3.1  | String $u$ in Subcase 3 . . . . .   | 17 |
| 3.2  | String $u$ in Subcase 3 when (3.5) holds . . . . .                                  | 18 |
| 3.3  | Subcase 3 when (C1) holds . . . . .   | 19 |
| 3.4  | Subcase 3 when (C2) holds . . . . .   | 20 |
| 3.5  | Subcase 3 when (C2) holds and $z < g$ . . . . .                                     | 21 |
| 3.6  | Subcase 3 when (C2) holds, $z < g$ , $g_1 < g_2$ , and $g' \leq z$ . . . . .        | 22 |
| 3.7  | Subcase 3 when (C2) holds, $z < g$ , $g_1 > g_2$ , and $g' \leq \ell$ . . . . .     | 23 |
| 3.8  | Subcase 3 when (C2) holds, $z < g$ , $g_1 > g_2$ , and $\ell < g' \leq z$ . . . . . | 24 |
| 3.9  | Subcase 3 when (C3) holds . . . . .   | 24 |
| 3.10 | String $uu_1$ in Subcase 7 . . . . .  | 27 |
| 3.11 | Subcase 7 when (C1) holds and $g > 0$ . . . . .                                     | 29 |
| 3.12 | Subcase 7 when (C1) holds and $g < 0$ . . . . .                                     | 31 |
| 3.13 | Subcase 7 when (C2) holds . . . . .   | 31 |
| 3.14 | Subcase 7 when (C3) holds . . . . .   | 32 |
| 3.15 | Subcase 7 when (C3) holds, $g_1 < g_2$ , and $g \leq z$ . . . . .                   | 33 |
| 3.16 | Subcase 7 when (C3) holds, $g_1 > g_2$ and $g \leq z$ . . . . .                     | 34 |
| 3.17 | Subcase 7 when (C4) holds . . . . .   | 35 |
| 4.1  | Case [db] . . . . .   | 38 |

# Chapter 1

## Introduction

There has for several years been considerable interest in the limitations that may exist on periodicity in strings. The “Three Squares Lemma” (Crochemore and Rytter 1995) showed that three squares could exist at the same position in a string only if the longest of the three was at least the sum of the lengths of the other two. A sequence of papers (Fan et al. 2006; Franek, Fuller, et al. 2012; Kopylova and Smyth 2012; Simpson 2007) greatly generalized this result by considering two squares  $\mathbf{u}^2$  and  $\mathbf{v}^2$  at the same position, with however the third square  $\mathbf{w}^2$  offset a distance  $k \geq 0$  to the right. First stated and proved as the “New Periodicity Lemma” (NPL) in Fan et al. (2006), the main theorem has since been made more specific: the existence of three neighbouring squares in certain well-defined configurations has been shown to cause a breakdown into repetitions of small period. The statement of the NPL includes 14 subcases, with 12 previously proven. This thesis contributes proofs of the two that remain.

Interest has been added to this research by a parallel development over the last dozen years or so: the attempt to specify sharp bounds on the number of maximal periodicities (“runs”) that can occur in any string of given length  $n$ . Kolpakov and Kucherov (2000) showed that the maximum number of runs—usually denoted  $\rho(n)$ —is linear in  $n$ , and moreover they described a linear-time algorithm to compute all the runs in any given string. But their proof of linearity was nonconstructive: the maximum number of runs was shown to be  $\mathcal{O}(n)$  but no constant of proportionality was specified. Subsequent research has shown that  $\rho(n)$  is at least  $0.9445757n$  (Kusano et al. 2013; Simpson 2010) and asymptotically at most  $1.029n$  (Crochemore et al. 2011), or, in other words, more or less the string length  $n$ .

What links these two streams of research is a simple observation:

If the maximum number of runs over all strings of length  $n$  is itself approximately  $n$ , then on average there will be about one run starting at each position. Thus, if two runs start at some position, there must be some other position, probably nearby, at which no run can start — “probably nearby” because the interference of overlapping squares typically precludes periodic behaviour at one or more positions within the range of the double periodicity. More generally, determining combinatorial constraints on the occurrence of overlapping squares (runs) may lead to a better characterization of  $\rho(n)$ .

There is a third avenue of research that relates closely to overlapping squares: the computation of all the runs/repetitions in a given string. At present the only way that this can be done involves global data structures (suffix array, longest common prefix array, Lempel-Ziv factorization) that need to be computed in an extended preprocessing phase. This seems uneconomical considering that runs are generally a local phenomenon, and it has been shown (Puglisi and Simpson 2008) that the expected number of runs in a string is much less than the string’s length (i.e. runs tend to occur sparsely). Another “local” problem, string searching, has local solutions: there exist many string searching algorithms that only preprocess the query, not the full text to be searched. Thus the current global approach to computing runs might be unnecessary but for the absence of a detailed understanding of the combinatorics of overlapping runs. A local approach would be desirable for space efficiency, particularly as string data, such as biological sequences, gets larger.

In Chapter 2, we review terminology, notation and the relevant background. Then, in Chapter 3, we prove Subcases 3 and 7 of the NPL. We conclude in Chapter 4 with a discussion of future research directions, namely, the general case of three overlapping squares (no two constrained to begin at the same position).



# Chapter 2

## Background

### 2.1 Strings

We begin with some basic terminology<sup>1</sup>. A *string* is a finite sequence of symbols (*letters*) drawn from some (possibly infinite) set  $\Sigma$  called the *alphabet*. The alphabet *size* is  $\sigma = |\Sigma|$ . To reduce notational clutter, we write a string  $\mathbf{x}$  in mathbold and its length  $x = |\mathbf{x}|$  in plain math font. We represent a string  $\mathbf{x}$  as an array  $\mathbf{x}[1..x]$  for  $x \geq 0$ . For  $x = 0$ ,  $\mathbf{x} = \varepsilon$ , the *empty string*. If  $\mathbf{x} = \mathbf{uvw}$ , then  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{w}$  are *substrings* of  $\mathbf{x}$ , and furthermore,  $\mathbf{u}$  is a *prefix* and  $\mathbf{w}$  is a *suffix* of  $\mathbf{x}$ . If  $\mathbf{vw} \neq \varepsilon$ , then  $\mathbf{u}$  is a *proper prefix*. Similarly, if  $\mathbf{uv} \neq \varepsilon$ ,  $\mathbf{w}$  is a *proper suffix*. If  $\mathbf{x} = \mathbf{uv} = \mathbf{wu}$  for  $u < x$ , then  $\mathbf{u}$  is a *border* of  $\mathbf{x}$ , that is, a proper prefix that equals a proper suffix. (Note that every nonempty string has an empty border.) If  $\mathbf{x} = \mathbf{uv}$ ,  $0 \leq u < x$ , then  $\mathbf{vu}$  is said to be the  $u^{\text{th}}$  *rotation* of  $\mathbf{x}$ , written  $R_u(\mathbf{x})$ .

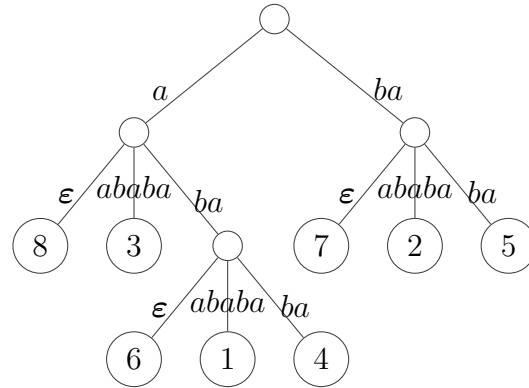
### 2.2 String data structures

In this section, we briefly describe a few commonly-used string data structures that, as explained in Section 2.4.1, help motivate our work:

- suffix tree
- suffix array
- longest common prefix array

---

<sup>1</sup>Our usage generally follows Smyth (2003).

Figure 2.1:  $ST_x$  for  $x = abaababa$ 

- LZ factorization

The *suffix tree* is an instance of a more general data structure, the trie. Used to index a set of strings, a trie (also known as a radix tree or prefix tree) is a tree in which the root is associated with the empty string, leaves are associated with strings in the set, and each internal node is associated with the longest common prefix of all the strings associated with its descendants (Fredkin 1960). The suffix tree  $ST_x$  of a string  $x$  is a trie containing all the suffixes of  $x$ . The leaves of a suffix tree are labeled with the starting position of the associated suffix, and a node’s children are ordered lexicographically. Figure 2.1 shows the suffix tree of  $x = abaababa$ .

Weiner (1973) introduced the suffix tree, along with a  $\mathcal{O}(n \log \sigma)$  time construction algorithm. A few years later, McCreight (1976) presented another  $\mathcal{O}(n \log \sigma)$  construction algorithm that is significantly more time- and space-efficient in practice (Giegerich and Kurtz 1997). Nearly twenty years passed before a third algorithm appeared (Ukkonen 1995). Ukkonen’s algorithm is on-line and simpler than its predecessors, but in practice marginally slower than McCreight’s (Giegerich and Kurtz 1997). Though not practical for large strings, a  $\mathcal{O}(n)$  time (independent of  $\sigma$ ) algorithm has also been conceived (Farach 1997).

Suffix trees enable many linear-time string algorithms (Gusfield 1997, Chapter 7) and have been memorably lauded for their “myriad virtues” (Apostolico 1985). Nevertheless, large (though linear) space requirements limit their usefulness. A space-efficient implementation might require  $10n$  bytes on average, and twice as much in the worst case (Kurtz 1999). Because

|                             |     |     |     |     |     |     |     |     |
|-----------------------------|-----|-----|-----|-----|-----|-----|-----|-----|
|                             | 1   | 2   | 3   | 4   | 5   | 6   | 7   | 8   |
| $\mathbf{x} =$              | $a$ | $b$ | $a$ | $a$ | $b$ | $a$ | $b$ | $a$ |
| $\text{SA}_{\mathbf{x}} =$  | 8   | 3   | 6   | 1   | 4   | 7   | 2   | 5   |
| $\text{LCP}_{\mathbf{x}} =$ | 0   | 1   | 1   | 3   | 3   | 0   | 2   | 2   |

Figure 2.2:  $\text{SA}_{\mathbf{x}}$  and  $\text{LCP}_{\mathbf{x}}$  for  $\mathbf{x} = abaababa$ 

of this significant drawback, the suffix tree has been mostly superseded by the more economical suffix array (Abouelhoda et al. 2004).

The *suffix array*  $\text{SA}_{\mathbf{x}}$  of a string  $\mathbf{x}$  is a sorted list of the suffixes of  $\mathbf{x}$ , succinctly encoded: for  $i, j \in [1..x]$ ,  $\text{SA}_{\mathbf{x}}[j] = i$ , where  $i$  is the starting position of the  $j^{\text{th}}$  suffix of  $\mathbf{x}$  in lexicographic order (Manber and Myers 1993). Figure 2.2 shows the suffix array of  $\mathbf{x} = abaababa$ . (Comparing Figures 2.2 and 2.1, note that performing a depth-first traversal of a suffix tree while printing leaf node labels produces the corresponding suffix array.) Whereas tree traversal is used to query suffix trees, suffix arrays can be queried by binary search. The original suffix array construction algorithms were  $\mathcal{O}(n \log n)$ -time, but linear algorithms were later developed (Kärkkäinen and Sanders 2003; Kim et al. 2003; Ko and Aluru 2003; Nong et al. 2009). However, according to a survey of suffix array construction algorithms (Puglisi, Smyth, et al. 2007), the fastest algorithm in practice (Maniscalco and Puglisi 2006) actually has worst-case  $\mathcal{O}(n^2 \log n)$  time complexity. The most space-efficient algorithms in practice (Maniscalco and Puglisi 2008, 2006; Manzini and Ferragina 2004) use as little as  $5n$  bytes on average for a string of length  $n$ .

The suffix array is enhanced by other data structures, most notably the *longest common prefix* array, in which  $\text{LCP}_{\mathbf{x}}[1] = 0$  and for  $j \in [2..x]$ ,  $\text{LCP}_{\mathbf{x}}[j]$  is the length of the longest prefix between the  $j^{\text{th}}$  and  $(j-1)^{\text{th}}$  suffixes in  $\text{SA}_{\mathbf{x}}$ . Figure 2.2 shows an example LCP array.

The Lempel-Ziv (LZ) factorization (Lempel and Ziv 1976) partitions a string into substrings:  $\mathbf{x} = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_k$ , where for all  $j \in [1..k]$ ,  $\mathbf{w}_j$  is

1. a single character that does not occur in  $\mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_{j-1}$ ; or
2. the longest substring that occurs twice in  $\mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_j$ .

In case 2, the first of the two occurrences may overlap with  $\mathbf{w}_j$ . Note that  $\mathbf{w}_1 = \mathbf{x}[1]$ . As an example, the LZ factorization of  $\mathbf{x} = abaababa$  is  $\mathbf{x} = (a)(b)(a)(aba)(ba)$ . Al-Hafeedh et al. (2012) compares the many LZ

factorization algorithms, some of which are linear-time, and which together offer various tradeoffs between time and space. All LZ factorization algorithms use an ST, SA, LCP, or other global data structures.

## 2.3 Periodicity

If  $\mathbf{x}[i] = \mathbf{x}[i + p]$  for all  $i \in [1 \dots x - p]$ , then  $\mathbf{x}$  has *period*  $p$ . Every period of a string corresponds to a border:

**Lemma 1** (Lothaire 2002, Section 8.1.1). *If  $\mathbf{v}$  is a border of  $\mathbf{w}$ , then  $\mathbf{w}$  has period  $w - v$ . Conversely, if  $w$  has period  $p$ , then it has border  $\mathbf{w}[1 \dots w - p]$ .*

For example, the string

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \mathbf{x} = & a & b & a & a & b & a & b & a & a & b \end{array} \quad (2.1)$$

has borders *abaab* and *ab*, hence corresponding periods 5 and 8, respectively.

The analysis of periodicity often involves strings of more than one period, or periodic strings that overlap. The next few lemmas express some possible consequences of coincident periodicities. The first follows readily from Lemma 1.

**Lemma 2** (Lothaire 2002, Lemma 8.1.1). *If  $\mathbf{x}$  has periods  $p$  and  $q$  such that  $q < p \leq x$ , then the border of  $\mathbf{x}$  of length  $x - q$  has period  $p - q$ .*

Another basic lemma applies to strings of one period with a substring of another period:

**Lemma 3** (Lothaire 2002, Lemma 8.1.3). *If  $\mathbf{x}$  has period  $p$  and there exists a substring  $\mathbf{u}$  of  $\mathbf{x}$  with  $p \leq u$  that has period  $q$ , where  $q$  divides  $p$ , then  $\mathbf{x}$  has period  $q$ .*

The next lemma, known as the Periodicity Lemma, is one of the most important in combinatorics on words, featuring in many correctness proofs of string algorithms. For strings of two periods  $p$  and  $q$ , the Periodicity Lemma provides the minimum length for which all strings of at least that length also have period  $\gcd(p, q)$ .

**Lemma 4** (Periodicity Lemma (Fine and Wilf 1965; Lothaire 2005)). *If  $\mathbf{x}$  has periods  $p$  and  $q$ , and  $p + q \leq x + \gcd(p, q)$ , then  $\mathbf{x}$  also has period  $\gcd(p, q)$ .*

For example, the string

$$\begin{array}{ccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ \mathbf{x} = & a & b & a & a & b & a & a & b & a & a & b & a & a \end{array}$$

has length  $n = 13$ , and periods  $p = 6$  and  $q = 9$ . Since  $d = \gcd(p, q) = 3$  and  $p + q = 15 < n + d = 16$ , the Periodicity Lemma allows us to infer that the string also has period  $d = 3$ .

In practice, the Periodicity Lemma is often applied via one of the following two corollaries. The first applies to overlapping strings of the same period, the second to overlapping strings of different periods.

**Lemma 5** (Lothaire 2002, Lemma 8.1.2). *If  $\mathbf{x} = \mathbf{uvw}$ , and  $\mathbf{uv}$  and  $\mathbf{vw}$  have period  $p \leq v$ , then  $\mathbf{x}$  has period  $p$ .*

**Lemma 6** (Simpson 2007, Section 1). *If  $\mathbf{x} = \mathbf{uvw}$ , where  $\mathbf{uv}$  has period  $p$ ,  $\mathbf{vw}$  has period  $q$ , and  $p + q \leq v + \gcd(p, q)$ , then  $\mathbf{x}$  has period  $\gcd(p, q)$ .*

The Periodicity Lemma has been generalized to three periods (Castelli et al. 1999), an arbitrary number of periods (Constantinescu and Ilie 2005; Holub 2006; Justin 2000; Tijdeman and Zamboni 2009), multiple dimensions (Simpson and Tijdeman 2003), and the case in which the length of a string with two periods does not satisfy the bound given by the Periodicity Lemma (Fraenkel and Simpson 2005). In particular, the generalization of the Periodicity Lemma to three periods uses the function

$$f(p_1, p_2, p_3) = \frac{1}{2} [p_1 + p_2 + p_3 - 2 \gcd(p_1, p_2, p_3) + h(p_1, p_2, p_3)] \quad (2.2)$$

where  $h$  is a function derived from Euclid’s algorithm for computing the greatest common divisor of three integers.

**Lemma 7** (Periodicity Lemma for Three Periods (Castelli et al. 1999)). *If a string  $\mathbf{x}$  has periods  $p_1$ ,  $p_2$  and  $p_3$ , with  $p_1 \leq p_2 \leq p_3$  and  $f(p_1, p_2, p_3) \leq x$  ( $f$  as defined by 2.2), then  $\mathbf{x}$  also has period  $\gcd(p_1, p_2, p_3)$ .*

As we will see, the New Periodicity Lemma also deals with three periods imposed by three “overlapping squares”. In a sense, Lemma 7 is more general than the New Periodicity Lemma because it only requires three periods rather than squares. On the other hand, the New Periodicity Lemma provides more information and covers a wider range of cases.

## 2.4 Repetitions and runs

If  $\mathbf{x} = \mathbf{v}\mathbf{u}^e\mathbf{w}$ , where  $e \geq 2$  is an integer, and  $\mathbf{u}$  is neither a suffix of  $\mathbf{v}$  nor a prefix of  $\mathbf{w}$  ( $e$  is maximum), then  $\mathbf{u}^e$  is said to be a *repetition* in  $\mathbf{x}$ . The integers  $u$  and  $e$  are the *period* and *exponent*<sup>2</sup>, respectively, of the repetition. The string (2.1) has repetitions  $(aba)^2$ ,  $(abaab)^2$ ,  $a^2$ ,  $(ab)^2$ ,  $(ba)^2$ , each of which is a *square*. In general, every repetition has a square prefix. We say that a square  $\mathbf{u}^2$  is *irreducible*<sup>3</sup> if  $\mathbf{u}$  is not itself a repetition, *regular* if  $\mathbf{u}$  has no square prefix, and *minimal* if no proper prefix of  $\mathbf{u}^2$  is a square. Note that minimality implies regularity, which in turn implies irreducibility.

If  $\mathbf{v} = \mathbf{x}[i..j]$  has period  $u$ , where  $v/u \geq 2$ , and if neither  $\mathbf{x}[i-1..j]$  nor  $\mathbf{x}[i..j+1]$  (whenever these are defined) has period  $u$ , then  $\mathbf{v}$  is said to be a *maximal periodicity* or *run* in  $\mathbf{x}$  (Main 1989) with a (now fractional) *exponent*  $e = v/u$ . All of the repetitions in (2.1) are runs except for  $(ab)^2$  and  $(ba)^2$ : these are substrings of the run  $\mathbf{v} = ababa = (ab)^{5/2}$ . In general, every repetition is a substring of some run; thus computing all the runs implicitly computes all the repetitions.

From Lemmas 1 and 4, it follows that if a string  $\mathbf{x}$  equals its  $u^{\text{th}}$  rotation, then it is a repetition of period  $\gcd(u, x)$ .

**Lemma 8** (Smyth 2003, Theorem 1.4.2). *For any string  $\mathbf{x}$ ,  $\mathbf{x} = \mathbf{u}\mathbf{v} = \mathbf{v}\mathbf{u}$  if and only if  $\mathbf{x}$  is a repetition of period  $\gcd(u, v) = \gcd(u, x - u) = \gcd(u, x)$ .*

The following lemma, used in Chapter 3, also connects repetitions and rotations.

**Lemma 9** (Kopylova and Smyth 2012, Lemma 8). *Suppose both  $\mathbf{x}$  and  $R_v(\mathbf{x})$ ,  $0 < v < x$ , have period  $u$ . Let  $\ell = x \bmod u > 0$ ,  $r = \lfloor \frac{x}{u} \rfloor$ , and  $d = \gcd(u, \ell)$ . Then:*

- (a) *if  $r = 1$  and  $v \geq \ell$ ,  $R_{v-\ell}(\mathbf{x})[1..2\ell]$  is a square of period  $\ell$ ;*
- (b) *if  $r = 1$  and  $v \leq \ell$ ,  $\mathbf{x}[1..v+\ell]$  has period  $\ell$ ;*
- (c) *if  $r > 1$  and  $v < u$ ,  $\mathbf{x}[1..v+\ell]$  has period  $\ell$ ; if moreover  $v+d \geq u$ , then  $\mathbf{x}$  is a repetition of period  $d$ ;*
- (d) *if  $r > 1$  and  $u \leq v \leq x - u$ ,  $\mathbf{x}[1..u+\ell]$ , hence  $\mathbf{x}$ , is a repetition of period  $d$ ;*

<sup>2</sup>We use  $\mathbf{u}^{(h)}$  (with the exponent in parentheses) to denote the  $h^{\text{th}}$  occurrence of  $\mathbf{u}$ .

<sup>3</sup>Others use the term *primitively rooted*.

- (e) if  $r > 1$  and  $x - u < v$ , where  $v' = v - (x - u)$ ,  $\mathbf{x}[v' + 1 .. u + \ell]$  has period  $\ell$ ; if moreover  $v' \leq d$ , then  $\mathbf{x}$  is a repetition of period  $d$ .

### 2.4.1 Algorithms for finding repetitions and runs

Three classical algorithms for computing all the repetitions in a string were proposed in the early 1980s:

1. Crochemore (1981)
2. Apostolico and Preparata (1983)
3. Main and Lorentz (1984)

Smyth (2003, page 340) points out that Crochemore’s algorithm basically constructs a suffix tree. The Apostolico algorithm uses a suffix tree explicitly. Main’s algorithm uses LZ factorization. All of these algorithms execute in  $\mathcal{O}(n \log n)$  time, asymptotically optimal since the *Fibonacci string*  $\mathbf{f}_k$ , defined by

$$\mathbf{f}_k = \begin{cases} b, & k = 0 \\ a, & k = 1 \\ \mathbf{f}_{k-1}\mathbf{f}_{k-2}, & k \geq 2 \end{cases}$$

contains  $\mathcal{O}(f_k \log f_k)$  squares (Crochemore 1981, Lemma 10; Fraenkel and Simpson 1999; Iliopoulos et al. 1997). To obtain a lower time complexity, Main (1989) introduced a new encoding of repetitions, the “maximal periodicity” or run, and described an LZ-factorization-based algorithm to compute all the “leftmost” runs in  $\mathcal{O}(n)$  time. This was extended in Kolpakov and Kucherov (2000) to compute all runs  $\mathcal{O}(n)$  time.

All of these algorithms for computing repetitions and runs use a suffix tree or LZ factorization, both global data structures with significant memory requirements. More efficient algorithms for computing runs have been proposed — for example, Chen et al. (2007) — but still with extensive pre-processing and the same general approach. (For a survey of algorithms for computing repetitions and runs, see Kopylov (2010, Chapter 3).)

While asymptotically optimal, the heavy-handed global approach seems inappropriate because runs are local and often sparsely occurring. Puglisi and Simpson (2008) derives a formula for the expected number of runs in a string as a function of its length  $n$  and alphabet size  $\sigma$ . This value is largest

for binary alphabets, for which the expected number of runs per unit length is 0.41. This value decreases with increasing  $\sigma$ : 0.24 for DNA strings ( $\sigma = 4$ ), 0.04 for protein ( $\sigma = 24$ ), and 0.01 for English-language text. Experiments confirm that these values are good predictions of the average number of runs in real-world data.

A future runs algorithm might dispense with global data structures, instead taking advantage of some combinatorial properties of runs to detect them during a single left-to-right scan.

## 2.4.2 The combinatorics of runs

The linearity of the Kolpakov and Kucherov (2000) algorithm was established by a complex proof that the maximum number  $\rho(n)$  of runs (irreducible and not extendible to the left or right) in any string of length  $n$  satisfies

$$\rho(n) \leq K_1 n - K_2 \sqrt{n} \log_2 n \quad (2.3)$$

for some universal positive constants  $K_1$  and  $K_2$ . The method of proof allowed no bounds to be placed on  $K_1$  and  $K_2$ , but based on computational evidence up to  $n = 60$ , it was conjectured that  $\rho(n) < n$  (for large  $n$ ). This is known as the “Runs Conjecture”. Over the last decade, the bounding of  $\rho(n)/n$  has become a growth industry. Rytter (2006) provided the first upper bound, showing that  $\rho(n)/n < 5$ . The bound improved successively, first to 3.9 (Rytter 2007), then to 3.48 (Puglisi, Simpson, and Smyth 2008), 1.6 (Crochemore and Ilie 2008), 1.52 (Giraud 2009, 2008), 1.048 (Crochemore et al. 2008), and most recently to 1.029 (Crochemore et al. 2011), the last result achieved using three years of CPU time on a supercomputer. The lower bound has progressed from 0.927 (Franek, Simpson, et al. 2003), to 0.944542 (Matsubara et al. 2008), to 0.944575 (Kusano et al. 2013; Simpson 2010). So, for sufficiently large  $n$ , we know that

$$0.944575 < \rho(n)/n \leq 1.029.$$

This work has confined the possible values of  $\rho(n)/n$  to a narrow range, but it seems not to have yielded the combinatorial knowledge needed for a new runs algorithm. Another approach, the subject of this thesis, has sought to find a combinatorial basis for estimating the maximum number of runs in a string by considering the consequences of three overlapping squares.



## 2.5 Three overlapping squares

The Runs Conjecture provides a simple motivation for studying overlapping squares (Fan et al. 2006). If  $\rho(n) \leq n$ , then the average number of runs starting at each of the  $n$  positions is at most one. When two runs occur at the same position, there must be another position at which no run occurs. Perhaps a proof of the Runs Conjecture is to be found in a detailed understanding of how two overlapping runs combine to make a third run impossible. Runs always have a square prefix, so analyzing overlapping squares seems a promising way to proceed<sup>4</sup>. It seems plausible that a position at which a third run is precluded lies within the range of the overlapping square prefixes. If this is true, a thorough understanding of the restrictions imposed by overlapping squares could lead directly to a proof of the Runs Conjecture, as well as a new paradigm for the computation of runs.

Recent research on overlapping squares extends the earlier “Three Squares Lemma” (Crochemore and Rytter 1995) that shows if three squares occur at the same position in a string, one of them must be long:

**Lemma 10** (Three Squares Lemma (Crochemore and Rytter 1995)). *Suppose  $\mathbf{u}$  is irreducible, and suppose  $\mathbf{v} \neq \mathbf{u}^j$  for any  $j \geq 1$ . If  $\mathbf{u}^2$  is a prefix of  $\mathbf{v}^2$ , in turn a proper prefix of  $\mathbf{w}^2$ , then  $w \geq u + v$ .*

The Fibonacci string has three square prefixes of lengths 6, 10, and  $10 + 6 = 16$ , respectively, showing the Three Squares Lemma is best possible:

$$\begin{array}{cccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ \mathbf{f} & = & a & b & a & a & b & \underline{a} & b & a & a & \underline{b} & a & a & b & a & b & \underline{a} \end{array}$$

The Three Squares Lemma has been generalized to the case with two of the squares occurring at the same position in a string and the third nearby, somewhat to the right, producing a result called the “New Periodicity Lemma”.

### 2.5.1 The New Periodicity Lemma

Here is the original NPL, as stated and proved in Fan et al. (2006):

**Lemma 11.** *If  $\mathbf{x}$  has regular prefix  $\mathbf{u}^2$  and irreducible prefix  $\mathbf{v}^2$ ,  $u < v < 2u$ , then for every  $k \in 0..v - u - 1$  and every  $w \in v - u + 1..v - 1$ ,  $w \neq u$ ,  $\mathbf{x}[k + 1..k + 2w]$  is not a square.*

<sup>4</sup>Though not because the maximum number of irreducible squares in a string of length  $n$  might be  $n$ ; as mentioned in Section 2.4.1, the optimal bound is known to be  $\mathcal{O}(n \log n)$ .

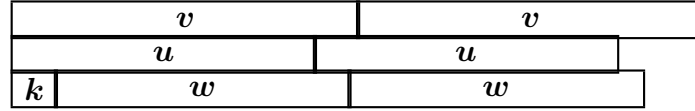


Figure 2.3: Three overlapping squares, as postulated in Lemma 13.

The proof required consideration of 14 subcases based on the magnitudes of  $k$  and  $w$  (see Figure 2.4), each of which led to a proof by contradiction of the regularity of  $\mathbf{u}$ . Subsequent work has split the range  $u < v < 2u$  into two sections  $(u, 3u/2]$  and  $(3u/2, 2u)$ , while eliminating the regularity condition altogether, as we now describe.

In Kopylova and Smyth (2012) it was shown that for  $u < v \leq 3u/2$ , the requirement that  $\mathbf{x} = \mathbf{v}^2$  with prefix  $\mathbf{u}^2$  necessitates

$$\mathbf{x} = \mathbf{t}_1^m \mathbf{t}_2 \mathbf{t}_1^{m+1} \mathbf{t}_2 \mathbf{t}_1, \quad (2.4)$$

where  $t_1 = v - u$ ,  $t_2 = u \bmod t_1$ ,  $m = \lfloor u/t_1 \rfloor \geq 2$  and  $\mathbf{t}_2$  is a proper prefix of  $\mathbf{t}_1$ . It was shown further that, except for  $m + 5$  precisely identified runs that always occur in  $\mathbf{x}$ , there could be no other runs of period greater than  $t_1$ . Thus for  $u < v \leq 3u/2$ , the structure of  $\mathbf{x}$  is well defined, even without reference to  $\mathbf{w}$ .

On the other hand, for  $3u/2 < v < 2u$ , there is a different breakdown:

**Lemma 12** (Fan et al. 2006, Lemma 9). *If  $\mathbf{x} = \mathbf{v}^2$  has prefix  $\mathbf{u}^2$ , then*

$$\mathbf{x} = (\mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 \mathbf{u}_1 \mathbf{u}_2)^2, \quad (2.5)$$

where  $u_1 = 2u - v > 0$  and  $u_2 = 2v - 3u > 0$ , if and only if  $\frac{3u}{2} < v < 2u$ . Observe that  $\mathbf{u} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1$  and  $\mathbf{v} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 \mathbf{u}_1 \mathbf{u}_2$ .

Note that setting  $\mathbf{t}_1 = \mathbf{u}_1 \mathbf{u}_2$ ,  $\mathbf{t}_2 = \mathbf{u}_1$  converts the form (2.5) into (2.4), but with  $m < 2$ . For this case, Kopylova and Smyth (2012) provided, with the assistance of a computer program, conjectures for the breakdown of  $\mathbf{x}$  in each of the 14 subcases. In half of the subcases (1, 2, 5, 6, and 8–10),  $\mathbf{x}$  was conjectured (and, in the same paper, proved) to always be a repetition of period  $d$ . In the other cases,  $\mathbf{x}$  was conjectured to have a different but still highly repetitive structure. An earlier paper (Simpson 2007) had already provided proofs for subcases 5, 6, and 10, as well as results for subcases 11–14 that were later refined in Franek, Fuller, et al. (2012). The latter also proved subcase 4, leaving only two of the 14 subcases unconfirmed. In this thesis, we prove the remaining two, subcases 3 and 7. Thus after much experimental and theoretical work, the revised NPL can be stated as follows:

| Subcase<br>$S$ | $k$                   | $k + w$                     | $k + 2w$                       | Special<br>Conditions |
|----------------|-----------------------|-----------------------------|--------------------------------|-----------------------|
| 1              | $0 \leq k \leq u_1$   | $k + w \leq u$              | $k + 2w \leq u + u_1$          | $k \geq u_2$          |
| 2              | $0 \leq k \leq u_1$   | $k + w \leq u$              | $k + 2w \leq u + u_1$          | $k < u_2$             |
| 3              | $0 \leq k \leq u_1$   | $k + w \leq u$              | $k + 2w > u + u_1$             | —                     |
| 4              | $0 \leq k \leq u_1$   | $u < k + w \leq u + u_1$    | —                              | —                     |
| 5              | $0 \leq k \leq u_1$   | $u + u_1 < k + w \leq v$    | —                              | —                     |
| 6              | $0 \leq k \leq u_1$   | $v < k + w < 2u$            | —                              | —                     |
| 7              | $u_1 < k < u_1 + u_2$ | $k + w \leq u + u_1$        | $k + 2w \leq 2u$               | —                     |
| 8              | $u_1 < k < u_1 + u_2$ | $k + w \leq u + u_1$        | $k + 2w > 2u$                  | —                     |
| 9              | $u_1 < k < u_1 + u_2$ | $u + u_1 < k + w \leq v$    | —                              | $w < u$               |
| 10             | $u_1 < k < u_1 + u_2$ | $k + w \leq v$              | $k + 2w \leq u + v$            | $w > u$               |
| 11             | $u_1 < k < u_1 + u_2$ | $k + w \leq v$              | $u + v < k + 2w \leq 2v - u_2$ | —                     |
| 12             | $u_1 < k < u_1 + u_2$ | $k + w \leq v$              | $2v - u_2 < k + 2w$            | —                     |
| 13             | $u_1 < k < u_1 + u_2$ | $v < k + w \leq 2u$         | —                              | —                     |
| 14             | $u_1 < k < u_1 + u_2$ | $2u < k + w < 2u + u_2 - 1$ | —                              | —                     |

Figure 2.4: The 14 subcases identified in Fan et al. (2006), slightly modified, for three neighbouring squares  $\mathbf{u}$ ,  $\mathbf{v}$ ,  $\mathbf{w}$  (with  $v - u < w < v$ ,  $w \neq u$ ,  $0 \leq k < v - u$ ).

| Subcases $S$     | Conditions                                     | Breakdown of $\mathbf{x}$  |
|------------------|--|--|
| 1, 2, 5, 6, 8–10 | $(\forall \mathbf{x}, \sigma = d)$             | $\mathbf{x} = \mathbf{d}^{x/d}$  |
| 3, 4, 7          | $(\forall \mathbf{x})$<br>specified cases      | $\mathbf{x} = \mathbf{d}_1^{u/d_1} \mathbf{d}_1^{v/d_1} \mathbf{d}_1^{(v-u)/d_1}$<br>$\mathbf{x} = \mathbf{d}^{x/d}$ |
| 11–14            | $\sigma = d$ or $d_2 \leq 2u - v$<br>otherwise | $\mathbf{x} = \mathbf{d}^{x/d}$<br>$\mathbf{x} = ((\mathbf{d}_3^{d_2/d_3})^{v/d_2})^2$                               |

Figure 2.5: Structure of  $\mathbf{x}$  for subcases  $S \in 1..14$ :  $\sigma$  is the largest alphabet size consistent with  $u, v, k, w$  (Franek, Fuller, et al. 2012);  $\mathbf{d}$ ,  $\mathbf{d}_1$  and  $\mathbf{d}_3$  are prefixes of  $\mathbf{x}$  with  $d = \gcd(u, v, w)$ ,  $d_1 = \gcd(u - w, v - u)$ ,  $d_2 = \gcd(u, v - w)$ ,  $d_3 = v \bmod d_2$ .

**Lemma 13.** *Suppose that a string  $\mathbf{x}$  has prefixes  $\mathbf{u}^2$  and  $\mathbf{v}^2$ ,  $3u/2 < v < 2u$ , and suppose further that a third square  $\mathbf{w}^2$  occurs at position  $k + 1$  of  $\mathbf{x}$ , where  $v - u < w < v$ ,  $w \neq u$ , and  $0 \leq k < v - u$ . Then for each of the 14 subcases  $S$  identified in Figure 2.4, the corresponding structure of  $\mathbf{x}$  is given in Figure 2.5.*

In other words,  $\mathbf{x}$  breaks down into repetitions of small period — essentially, the postulate of three such squares cannot be satisfied.

## 2.5.2 The general case characterized

The proof of the New Periodicity Lemma now complete, we believe that further generalization is of interest: what happens when the three squares  $\mathbf{u}^2$ ,  $\mathbf{v}^2$ ,  $\mathbf{w}^2$  are merely constrained to be “neighbouring”, without the requirement that  $\mathbf{u}^2$  and  $\mathbf{v}^2$  occur at the same position? What is an appropriate formulation of such a problem? What relative values of  $k, u, v, w$  are of combinatorial interest?

The following lemma, to appear in a forthcoming paper (Bland and Smyth 2014), begins to address these questions. It states the consequences of a square  $\mathbf{u}^2$  beginning at some position  $i$  in a string and overlapping with a second square  $\mathbf{v}^2$  at position  $i + k$ ,  $k \geq 0$ , to its right.

**Lemma 14.** *Suppose  $\mathbf{x}$  has prefixes  $\mathbf{u}^2$  and  $\mathbf{kv}^2$ ,  $k \geq 0$ , where  $x = \max(2u, k + 2v)$ ,  $k \leq u < 2v$ .*

(a)  $k + v < u < 2v$  ( $k < \min(v - 1, u - v)$ ) :

$$\mathbf{x} = (\mathbf{p}^e \mathbf{z})^2 = \mathbf{p}^e \mathbf{q}^f \mathbf{q}^{f-e} = \mathbf{p}^e \mathbf{q}^f \mathbf{p}[k + 1 \dots u - v],$$

where  $\mathbf{p} = \mathbf{u}[1 \dots u - v]$ ,  $e = \frac{k+v}{u-v} > 1$ ,  $\mathbf{z} = \mathbf{v}[1 \dots u - (k + v)]$ ,  $\mathbf{q} = R_k(\mathbf{p})$ ,  $f = \frac{u}{u-v} > 2$ ,  $f - e \leq 1$ .

(b)  $\frac{k}{2} + v \leq u \leq k + v$  ( $1 \leq u - v \leq k \leq 2(u - v)$ ) :

$$\mathbf{x} = (\mathbf{z}\mathbf{p}^e)^2 = (\mathbf{q}[1 \dots k + v - u]\mathbf{p}^e)^2 = (\mathbf{k}\mathbf{p}^{e-1})^2,$$

where  $\mathbf{z} = \mathbf{u}[1 \dots k + v - u]$ ,  $\mathbf{p} = \mathbf{v}[1 \dots u - v]$ ,  $e = 1 + \frac{u-k}{u-v} \geq 1$ ,  $\mathbf{q} = R_c(\mathbf{p})$ ,  $c = (u - k) \bmod (u - v)$ .

(c)  $v < u < \frac{k}{2} + v$  ( $k > 2(u - v)$ ) :

$$\mathbf{x} = (\mathbf{q}\mathbf{y}\mathbf{p}^e)^2 \mathbf{y},$$

where  $\mathbf{p} = \mathbf{v}[1 \dots u - v]$ ,  $e = 1 + \frac{u-k}{u-v} > 1$ ,  $\mathbf{q} = R_c(\mathbf{p})$ ,  $c = (u - k) \bmod (u - v)$ ,  $\mathbf{y} = \mathbf{v}[2u - (k + v) + 1 \dots v]$ . Moreover, both  $\mathbf{x}$  and  $\mathbf{k}\mathbf{v}$  have border  $\mathbf{q}\mathbf{y}$ .

(d)  $\frac{2(k+v)}{3} \leq u < v$  ( $k \leq \frac{3u}{2} - v < \frac{v}{2}$ ) :

$$\mathbf{x} = (\mathbf{k}\mathbf{p}^e)^2 \mathbf{q}\mathbf{k}\mathbf{p},$$

where  $\mathbf{p} = \mathbf{v}[1 \dots v - u]$ ,  $e = \frac{u-k}{v-u} > 1$ ,  $\mathbf{q} = R_c(\mathbf{p})$ ,  $c = (u - k) \bmod (v - u)$ . Both  $\mathbf{x}$  and  $\mathbf{k}\mathbf{v}$  have border  $\mathbf{k}\mathbf{p}$ .

(e)  $\frac{k+v}{2} < u < \frac{2(k+v)}{3} < v$  ( $\frac{3u-2v}{2} < k < 2u - v < u$ ) :

$$\mathbf{x} = \mathbf{k}(\mathbf{p}^e \mathbf{k}\mathbf{p})^2,$$

where  $\mathbf{p} = \mathbf{v}[1 \dots v - u]$ ,  $e = \frac{u-k}{v-u} > 1$ .

(f)  $k \leq u \leq \frac{k+v}{2}$  ( $\mathbf{u}^2$  a prefix of  $\mathbf{k}\mathbf{v}$ ) :

$$\mathbf{x} = \mathbf{k}(\mathbf{p}^e \mathbf{z})^2,$$

where  $\mathbf{p} = \mathbf{u}[k + 1 \dots u]\mathbf{u}[1 \dots k]$ ,  $e = \frac{2u-k}{u} \geq 1$ ,  $\mathbf{z} = \mathbf{v}[2u - k + 1 \dots v]$ .

Because every instance of three overlapping squares can be seen as two pairs of two overlapping squares, this lemma can characterize three overlapping squares in any configuration. We first use this lemma in Chapter 3 to prove the two remaining subcases of the NPL, then in Chapter 4 we discuss its application to the general case of three overlapping squares.

# Chapter 3

## Completing the New Periodicity Lemma

In this chapter<sup>1</sup>, we prove the two remaining subcases of Lemma 13.

### 3.1 Subcase 3

We first deal with the general case valid for all occurrences of Subcase 3, then go on to identify circumstances in which  $\mathbf{x}$  is constrained to be a repetition of small period  $d = \gcd(u, v, w)$ .

**Lemma 15** (Subcase 3). *Suppose that a string  $\mathbf{x}$  has prefixes  $\mathbf{u}^2$  and  $\mathbf{v}^2$ ,  $3u/2 < v < 2u$ , and suppose further that a third square  $\mathbf{w}^2$ ,  $w \neq u$ , occurs at position  $k + 1$  of  $\mathbf{x}$ , where*

$$0 \leq k \leq u_1 < u_1 + u_2 < w < v \quad (3.1)$$

$$k + w \leq u \quad (3.2)$$

$$k + 2w > u + u_1 \quad (3.3)$$

and  $u_1 = 2u - v$  and  $u_2 = 2v - 3u$ . Then  $\mathbf{x} = \mathbf{d}_1^{u/d_1} \mathbf{d}_1^{v/d_1} \mathbf{d}_1^{(v-u)/d_1}$ , where  $d_1 = \gcd(u - w, v - u)$ .

---

<sup>1</sup>This chapter is to appear in a forthcoming paper (Bland and Smyth 2014).

|                      |              |                |                      |              |
|----------------------|--------------|----------------|----------------------|--------------|
| $\mathbf{u}_1^{(1)}$ |              | $\mathbf{u}_2$ | $\mathbf{u}_1^{(2)}$ |              |
| $\mathbf{k}$         | $\mathbf{w}$ |                |                      | $\mathbf{w}$ |
| $\mathbf{p}^e$       |              |                |                      | $\mathbf{z}$ |

Figure 3.1: String  $\mathbf{u}$  in Subcase 3

*Proof.* By Lemma 12, the overlap of  $\mathbf{u}^2$  and  $\mathbf{v}^2$  forces  $\mathbf{x} = (\mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_1\mathbf{u}_2)^2$ , with  $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$ . By Lemma 14(a),  $\mathbf{u} = \mathbf{p}^e\mathbf{z}$ , where  $\mathbf{z} = \mathbf{w}[1..u-(k+w)]$ ,  $\mathbf{p} = \mathbf{u}[1..u-w]$  and  $e = \frac{k+w}{u-w} > 1$ .

We first show that if  $\mathbf{u}$  has period  $p = u - w$ , the lemma holds. Note that  $\mathbf{u}$  has period  $u_1 + u_2$  and

$$u_1 + u_2 + p = u + u_1 + u_2 - w < u$$

since  $u_1 + u_2 < w$  from (3.1). Therefore, assuming  $\mathbf{u}$  has period  $p$ ,  $\mathbf{u} = \mathbf{x}[1..u]$  has period  $d_1 = \gcd(p, u_1 + u_2)$  by Lemma 4. It follows that  $\mathbf{u}_1\mathbf{u}_2 = \mathbf{x}[u + v + 1..x]$ , a prefix of  $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$ , has period  $d_1$  as well. Finally,  $\mathbf{x}[u + 1..u + v] = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$  has period  $u_1 + u_2$  and prefix  $\mathbf{u}$  of length  $u > u_1 + u_2$  with period  $d_1$ . Since  $d_1 = \gcd(u - w, u_1 + u_2)$  divides  $u_1 + u_2$ ,  $\mathbf{x}[u + 1..u + v]$  has period  $d_1$  by Lemma 3. Thus the lemma holds assuming  $\mathbf{u}$  has period  $p$ .

Note first from (3.1) and (3.2) that  $u_1u_2 < w < u$ , hence that  $u_1 - p = u_1 - (u - w) > 0$ . Then to see that  $\mathbf{u}$  in fact has period  $p$ , consider two cases:

$$u_1 \geq k + w - (u_1 + u_2) \geq p \tag{3.4}$$

and

$$k + w - (u_1 + u_2) < p < u_1. \tag{3.5}$$

In the first case, the prefix  $\mathbf{k}\mathbf{w} = \mathbf{p}^e$  of  $\mathbf{u}$  extends at least  $p$  positions into the suffix  $\mathbf{u}_1^{(2)}$ . Since  $\mathbf{u}_1$  is a prefix of  $\mathbf{k}\mathbf{w}$ ,  $\mathbf{u}_1$  has period  $p$ , and therefore  $\mathbf{u}$  has a prefix and suffix of period  $p$  which overlap by at least  $p$ . Consequently, by Lemma 6,  $\mathbf{u}$  has period  $p$ .

The second case (3.5) is more complicated (see Figure 3.2). Both  $\mathbf{p}$  and  $\mathbf{u}_1$  are prefixes of  $\mathbf{u}$ , so  $\mathbf{p}$  is a proper prefix of  $\mathbf{u}_1$ . Both  $\mathbf{u}_1$  and  $\mathbf{z}$  are suffixes of  $\mathbf{u}$ , and

$$z = u - k - w \leq p < u_1,$$

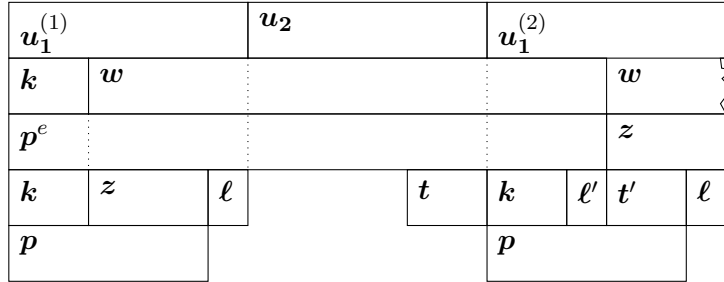


Figure 3.2: String  $u$  in Subcase 3 when (3.5) holds

so  $z$  is a proper suffix of  $u_1$ . The prefix  $p$  and the suffix  $z$  of  $u_1$  must overlap because (3.5) is equivalent to  $u_1 < p + z$ . Noting that  $p = u - w = k + z$  so that  $p = kz$ , we have (Figure 3.2)

$$z = \ell' t' = t' \ell \tag{3.6}$$

and

$$u_1 = p\ell = kz\ell = k\ell'z = k\ell't'\ell, \tag{3.7}$$

where  $\ell$  and  $\ell'$  are respectively the proper suffix and proper prefix of  $z$  of length  $\ell = \ell' = u_1 - p$ , and  $t'$  is the border of  $z$  of length  $t' = z - \ell$ . Since by (3.6)  $z$  has border  $t'$ , it therefore has period  $\ell$ , as does  $\ell'z\ell = \ell'\ell't'\ell = \ell't'\ell\ell$ . Since  $u_1$  has period  $p$ ,  $p$  has prefix  $\ell$ ;  $p = kz$  also has suffix  $\ell$ , so that  $p$  has border  $\ell$ . Observe also that because  $k + w - (u_1 + u_2) = k + \ell$ ,  $k\ell'$  is the prefix of  $u_1$  that overlaps  $w$ .

Let  $t$  be the suffix of  $u_1u_2$  of length  $t = t'$ . Then  $w$  has suffix  $tk\ell'$  in which  $k\ell'$  is a prefix of  $u_1$ . Recall  $u_1 = k\ell't'\ell$  has period  $p = k + t + \ell$ . If  $t = t'$ , then  $tu_1$  has period  $p$ ; moreover,  $kw = p^e$  and  $tu_1$  share substring  $tk\ell'$  of length  $p$ , so  $u$  has period  $p$ , as desired. Hence, it will suffice to show  $t = t'$ .

From  $kw = p^e$ , where  $e > 1$ , a complete copy of  $p$  occurs  $h = \lfloor e \rfloor$  times in  $kw$ . Three cases arise based on where in  $u$  the  $h^{\text{th}}$  occurrence of  $p$  ends:

- (C1)  $p^{(h)}$  ends inside the suffix  $t$  of  $u_2$ .
- (C2)  $p^{(h)}$  ends inside the prefix  $k$  of  $u_1^{(2)}$ .
- (C3)  $p^{(h)}$  ends inside the suffix  $\ell'$  of  $w$ .

We will see that  $t = t'$  in each of these cases.



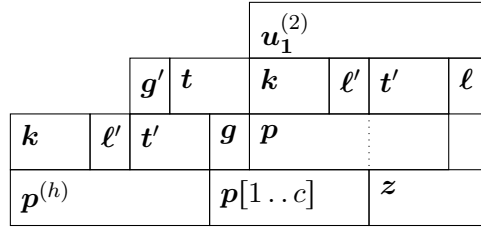


Figure 3.3: Subcase 3 when (C1) holds

**(C1)**

Suppose (C1) holds; that is,  $\mathbf{p}^{(h)}$  ends inside the suffix  $\mathbf{t}$  of  $\mathbf{u}_2$ . We introduce the “gap”  $g = u_1 + u_2 - ph$ , a measure of the overlap between  $\mathbf{t}$  and the suffix  $\mathbf{t}'$  of  $\mathbf{p}^{(h)}$ . Note that if  $g = 0$ , then  $\mathbf{t} = \mathbf{t}'$  immediately. Let  $\mathbf{g} = \mathbf{t}[t - g + 1..t] = \mathbf{p}[1..g]$  be the suffix of  $\mathbf{t}$  that follows  $\mathbf{p}^{(h)}$ , and let  $\mathbf{g}' = \mathbf{t}'[1..g]$  be the prefix of  $\mathbf{t}'$  that precedes  $\mathbf{t}$ . Also, let

$$c = (k + w) \bmod p = g + k + \ell$$

and observe that  $\mathbf{k}\mathbf{w}$  has suffix  $\mathbf{p}[1..c] = \mathbf{g}\mathbf{k}\ell' = \mathbf{k}\ell'\mathbf{g}'$ .

Thus  $\mathbf{p}[1..c]$  has border  $\mathbf{k}\ell'$  and therefore period  $g$ . String  $\ell'\mathbf{g}'$  has period  $\ell$  as a prefix of  $\mathbf{z} = \ell'\mathbf{t}'$ , and period  $g$  as a suffix of  $\mathbf{p}[1..c]$ , so by Lemma 4 it has period  $\gcd(g, \ell)$ . Then  $\mathbf{p}[1..c]$  has period  $g$  and suffix  $\ell'\mathbf{g}'$  of period  $\gcd(g, \ell) \mid g$  and length  $\ell + g \geq g$ , so that by Lemma 3  $\mathbf{p}[1..c]$  itself has period  $\gcd(g, \ell)$ . Both  $\mathbf{p}[1..c]$  and  $\ell'\mathbf{z}$  have period  $\ell$  and share substring  $\ell'$ , so  $\mathbf{p}[1..c]\mathbf{z}$  has period  $\ell$  by Lemma 6. It also has substring  $\mathbf{p}$ , so  $\mathbf{p}$  has period  $\ell$ . Because  $\mathbf{p}$  has border  $\ell$  as well as period  $\ell$ , any power of  $\mathbf{p}$  has period  $\ell$ . It follows that  $\mathbf{k}\mathbf{w} = \mathbf{p}^e$  has period  $\ell$  and, since  $\mathbf{p}[1..c]\mathbf{z}$  has period  $\ell$  and shares with  $\mathbf{k}\mathbf{w}$  a substring of length  $c > \ell$ ,  $\mathbf{u}$  has period  $\ell$  by Lemma 6. Recall that  $\mathbf{u}$  has substring  $\mathbf{p}[1..c]$  of period  $\gcd(g, \ell) \mid \ell$ , so  $\mathbf{u}$  itself has period  $\gcd(g, \ell)$  by Lemma 3. Recalling that  $\mathbf{t}$  is a suffix of  $\mathbf{t}'\mathbf{g}$  and that both are substrings of  $\mathbf{u}$ , we find that  $\mathbf{t}$  and  $\mathbf{t}'$  have period  $\gcd(g, \ell)$  and suffix  $\mathbf{g}$ , so  $\mathbf{t} = \mathbf{t}'$ .

**(C2)**

Suppose (C2) holds; that is,  $\mathbf{p}^{(h)}$  ends inside the prefix  $\mathbf{k}$  of  $\mathbf{u}_1^{(2)}$ . Let

$$g = ph - u_1 - u_2$$

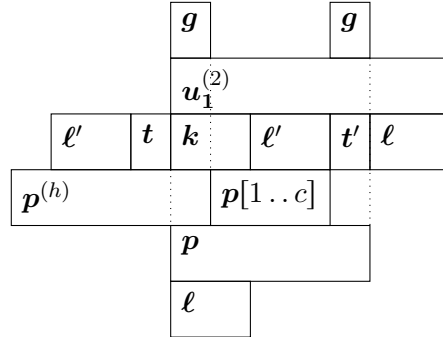


Figure 3.4: Subcase 3 when (C2) holds

be the overlap of  $\mathbf{p}^{(h)}$  with  $\mathbf{k}$ , and let  $\mathbf{g} = \mathbf{k}[1..g]$ . Note that if  $g = 0$ , then  $\mathbf{t} = \mathbf{t}'$  immediately. Otherwise,  $\mathbf{g}$  is a border of  $\mathbf{p}$ . Let

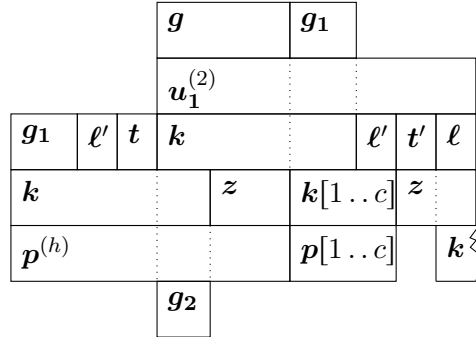
$$c = (k + w) \bmod p = k + \ell - g$$

and observe that, from the overlap of the suffix  $\mathbf{p}[1..c]$  of  $\mathbf{k}\mathbf{w}$  and the prefix  $\mathbf{p}$  of  $\mathbf{u}_1$ ,  $\mathbf{p}[1..k + \ell] = \mathbf{k}\ell'$  has period  $g$ . Also note that since by (3.7)  $\mathbf{p} = \mathbf{k}\ell'\mathbf{t}'$  and since  $\mathbf{k}\mathbf{w} = \mathbf{p}^e$  has period  $p$ , therefore  $\mathbf{p}^{(h)}$  has suffix  $\ell'\mathbf{t}\mathbf{g}$ . Consider four cases:  $0 < g < \ell$ ,  $g = \ell$ ,  $\ell < g \leq z$ , and  $z < g$ .

If  $0 < g < \ell$ , then since  $\ell$  is a prefix of  $\mathbf{u}_1$ ,  $\ell$  has period  $g$  as a prefix of  $\mathbf{k}\ell'$ . Recall that  $\ell$  is also a suffix of  $\mathbf{p}$ , so  $\ell$  has border  $\mathbf{g}$  and period  $\ell - g$ . Hence, by Lemma 4,  $\ell$  has period  $\gcd(g, \ell - g) = \gcd(g, \ell)$ . Recall also that  $\ell'z$  has period  $\ell$  and substring  $\ell$  of period  $\gcd(g, \ell) \mid \ell$ , so by Lemma 3,  $\ell'z$  has period  $\gcd(g, \ell)$ . Prefix  $\ell'$  of  $\ell'z$  then has period  $\gcd(g, \ell) \mid g$ , and since  $\ell'$  is also a suffix of the string  $\mathbf{k}\ell'$  of period  $g$ ,  $\mathbf{k}\ell'$  has period  $\gcd(g, \ell)$  by Lemma 3. Since  $\ell'z$  and  $\mathbf{k}\ell'$  have period  $\gcd(g, \ell)$ ,  $\mathbf{u}_1 = \mathbf{k}\ell'z$  has period  $\gcd(g, \ell)$  by Lemma 6. Both  $\mathbf{t}\mathbf{g}$  and  $\mathbf{t}'\mathbf{g}$  are substrings of  $\mathbf{u}_1$ , so  $\mathbf{t} = \mathbf{t}'$ .

If  $g = \ell$ , then  $\mathbf{p}$  has suffixes  $\mathbf{t}\ell$  and  $\mathbf{z} = \ell'\mathbf{t} = \mathbf{t}'\ell$ , so immediately  $\mathbf{t} = \mathbf{t}'$ . Note that  $\mathbf{k}\ell'$  and  $\ell'z$  have period  $g = \ell$ , so by Lemma 6,  $\mathbf{u}_1 = \mathbf{k}\ell'\mathbf{t}'$  has period  $\ell$ .

If  $\ell < g \leq z$ , then since  $\ell$  is a border of  $\mathbf{g}$ ,  $\mathbf{g}$  has period  $g - \ell$ ; it also has period  $\ell$  as a substring of the suffix  $\mathbf{z}$  of  $\mathbf{p}$ , and thus by Lemma 4 period  $\gcd(g, g - \ell) = \gcd(g, \ell)$ . String  $\mathbf{g}$  is a substring of  $\mathbf{k}\ell'$ , which as we have seen has period  $g$ , so that by Lemma 3,  $\mathbf{k}\ell'$  has period  $\gcd(g, \ell)$ . Since  $\ell'z$  and  $\mathbf{k}\ell'$  have period  $\gcd(g, \ell)$ ,  $\mathbf{u}_1 = \mathbf{k}\ell'z$  has period  $\gcd(g, \ell)$  by Lemma 6. Both  $\mathbf{t}\mathbf{g}$  and  $\mathbf{t}'\mathbf{g}$  are substrings of  $\mathbf{u}_1$ , so  $\mathbf{t} = \mathbf{t}'$ .

Figure 3.5: Subcase 3 when (C2) holds and  $z < g$ 

If  $z < g$ , then, as shown in Figure 3.5, the suffix  $z$  of  $p^{(h)}$  is a substring of the prefix  $k$  of  $u_1$ , and  $\ell't$  is a substring of the prefix  $k$  of  $p^{(h)}$ .  $k$  also has two borders  $g_1$  and  $g_2$ :  $g_1$  is the border of  $k$  of length  $g_1 = k - g$  resulting from the overlap of the prefix  $k$  of  $u_1$  with  $p[1..c] = k[1..c]$ , while  $g_2$  is the border of  $k$  of length  $g_2 = g - z$  resulting from the overlap of the prefix  $k$  of  $p^{(h)}$  with the prefix  $k$  of  $u_1$ . We then have  $k = g_1\ell'tg_2 = g_2\ell't'g_1$ . Also recall that  $\ell$  is a prefix of  $p = kz$ , so that either  $\ell$  is a prefix of  $g_1$  or  $g_1$  is a prefix of  $\ell$ .

If  $g_1 = g_2$ , then  $t = t'$  immediately. If  $g_1 \neq g_2$ , then several cases arise:

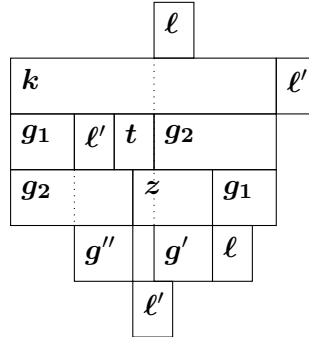
1.  $g_1 < g_2$

Let  $g' = g'' = g_2 - g_1$ , let  $g'$  be the prefix of  $g_2$  such that  $g_2 = g'g_1$ , and let  $g''$  be the suffix of  $g_2$  such that  $g_2 = g_1g''$ . Observe that  $g_1$  is a border of  $g_2$ , so  $g_2$  has period  $g'$ .

- (a)  $g' \leq z$  (Figure 3.6)

The demonstration requires several steps:

- $g'\ell$  has period  $\ell$  as a suffix of  $z\ell$ ; it also has border  $\ell$  and period  $g'$ , so by Lemma 4,  $g'\ell$  has period  $\gcd(g', \ell)$ .
- $z\ell$  has period  $\ell$  and suffix  $g'\ell$  of period  $\gcd(g', \ell) \mid \ell$ , so by Lemma 3,  $z\ell$  has period  $\gcd(g', \ell)$ .

Figure 3.6: Subcase 3 when (C2) holds,  $z < g$ ,  $g_1 < g_2$ , and  $g' \leq z$ 

- $g_2$  has period  $g'$  and prefix  $g'$  of period  $\gcd(g', \ell) \mid g'$ , so by Lemma 3,  $g_2$  has period  $\gcd(g', \ell)$ .
- $z\ell$  and  $g_2$  have period  $\gcd(g', \ell)$  and share substring  $g'$ , so by Lemma 6,  $zg_1$  has period  $\gcd(g', \ell)$ .
- $g''\ell'$  has border  $\ell'$  and period  $g' = g''$ ; it also has prefix  $g''$  which, as a suffix of  $g_2$ , has period  $\gcd(g', \ell) \mid g'$ , so by Lemma 3,  $g''\ell'$  has period  $\gcd(g', \ell)$ .
- $g_2$  and  $g''\ell'$  have period  $\gcd(g', \ell)$  and share substring  $g''$ , so by Lemma 6,  $g_2\ell'$  has period  $\gcd(g', \ell)$ .
- $g_2\ell'$  and  $zg_1$  have period  $\gcd(g', \ell)$  and share substring  $\ell'$ , so by Lemma 6,  $k = g_2zg_1$  has period  $\gcd(g', \ell)$ .

Since  $\ell't$  and  $\ell't$  are substrings of  $k$ , therefore  $t = t'$ .

Note that  $g''z\ell$  has period  $\gcd(g', \ell)$ , so that  $k = g_1\ell'tg_1g''$  and  $g''z\ell$  have period  $\gcd(g', \ell)$  and share substring  $g''$ , so by Lemma 6,  $u_1 = kz\ell$  has period  $\gcd(g', \ell)$ .

(b)  $g' > z$

$k$  has period  $k - g_2 = g_1 + z$ , so  $k$  has a prefix  $g_1\ell'tg'_2 = g'_2zg_1$ , where  $g'_2$  is a prefix of  $g_2$  and  $|g'_2 - g_1| \leq z$ , so one of cases 1(a) and 2(a) applies.

2.  $g_1 > g_2$

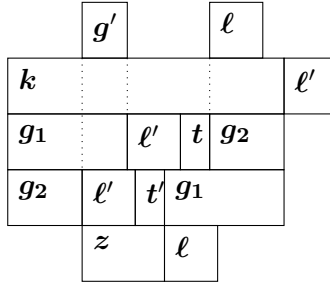


Figure 3.7: Subcase 3 when (C2) holds,  $z < g$ ,  $g_1 > g_2$ , and  $g' \leq \ell$

Let  $g' = g_1 - g_2$ , and let  $g'$  be the suffix of  $g_1$  such that  $g_1 = g_2g'$ . Observe that  $g_2$  is a border of  $g_1$ , so  $g_1$  and  $g_2$  have period  $g'$ .

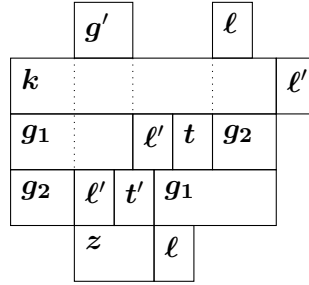
(a)  $g' \leq z$  (Figures 3.7 & 3.8)

Again several steps are required:

- $g'\ell'$  has period  $\ell$  as a prefix of  $z\ell$  and also shares prefix  $\ell'$  with  $z$ , so it has border  $\ell'$  and period  $g'$ .
- $g'\ell'$  has periods  $\ell$  and  $g'$ , so by Lemma 4,  $g'\ell'$  has period  $\gcd(g', \ell)$ .
- $g_1$  has period  $g'$  and suffix  $g'$  of period  $\gcd(g', \ell) \mid g'$ , so by Lemma 3,  $g_1$  has period  $\gcd(g', \ell)$ .
- $z\ell$  has period  $\ell$  and prefix  $g'$  of period  $\gcd(g', \ell) \mid \ell$ , so again by Lemma 3,  $z\ell$  has period  $\gcd(g', \ell)$ .
- $z\ell$  and  $g_1$  have period  $\gcd(g', \ell)$  and share a substring of length at least  $\min(g', \ell)$ , so by Lemma 6,  $zg_1$  has period  $\gcd(g', \ell)$ .

Since  $\ell't$  and  $\ell't'$  are substrings of  $zg_1$ , therefore  $t = t'$ .

Note that  $g_1$  and  $zg_1$  have period  $\gcd(g', \ell)$  and share substring  $g'$ , so by Lemma 6,  $k = g_2zg_1$  has period  $\gcd(g', \ell)$ .  $g'z\ell$  then has period  $\gcd(g', \ell)$ , so that  $k = g_2zg_2g'$  and  $g'z\ell$  have period  $\gcd(g', \ell)$  and share substring  $g'$ , so by Lemma 6,  $u_1 = kz\ell$  has period  $\gcd(g', \ell)$ .

Figure 3.8: Subcase 3 when (C2) holds,  $z < g$ ,  $g_1 > g_2$ , and  $\ell < g' \leq z$ (b)  $g' > z$ 

$k$  has period  $k - g_1 = g_2 + z$ , so  $k$  has a prefix  $g_2 z g'_1 = g'_1 \ell' t g_2$ , where  $g'_1$  is a prefix of  $g_1$  and  $|g'_1 - g_2| \leq z$ , so one of cases 1(a) and 2(a) applies.

(C3)

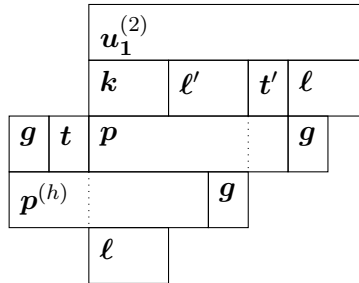


Figure 3.9: Subcase 3 when (C3) holds

Suppose (C3) holds; that is,  $p^{(h)}$  ends inside the suffix  $\ell'$  of  $w$ . Let  $g = k + w - ph$  and let  $g = \ell'[\ell - g + 1 .. \ell]$  be the suffix of  $\ell'$  that follows  $p^{(h)}$ . Because  $kw$  (of which  $\ell'$  is a suffix) has period  $p$ ,  $g$  is a prefix of  $p$ . Recall that  $\ell$  is a prefix of  $p$ , so  $g$  is also a prefix of  $\ell$ . From  $p^{(h)}$  and the occurrence of  $p$  that prefixes  $u_1$ , we have

$$p = g t k (\ell' [1 .. \ell - g]) = k (\ell' [1 .. \ell - g]) g t'$$

and  $\mathbf{p}$  has period  $t + g$ .

Consider the string  $\ell't'g$ , which occurs near the end of  $\mathbf{u}_1$  as a prefix of  $\ell't'\ell = \ell'z$ . As a substring of  $\ell'z = \ell'\ell't'$ , it has period  $\ell$ . Since  $\mathbf{p}$  has period  $t + g$  and suffixes  $\ell't'$  and  $gt'$ ,  $\ell't'g$  also has period  $t + g$ .  $\ell't'g$  has periods  $t + g$  and  $\ell$ , so by Lemma 4, it has period  $\gcd(t + g, \ell)$ . Now  $\mathbf{p}$  has period  $t + g$  and suffix  $gt'$  of period  $\gcd(t + g, \ell) \mid t + g$ , so  $\mathbf{p}$  itself has period  $\gcd(t + g, \ell)$ . Because  $\mathbf{p}$  has border  $\ell$  as well as period  $\gcd(t + g, \ell)$ , any power of  $\mathbf{p}$  has period  $\gcd(t + g, \ell)$ . Thus  $\mathbf{kw} = \mathbf{p}^e$  has period  $\gcd(t + g, \ell)$  and, since  $\ell'z$  has period  $\gcd(t + g, \ell)$  and shares with  $\mathbf{kw}$  a substring of length  $\ell$ ,  $\mathbf{u}$  has period  $\gcd(t + g, \ell)$  by Lemma 6. Since  $t\ell$  and  $t'\ell$  are substrings of  $\mathbf{u}$ , therefore  $t = t'$ .

This completes the proof of Subcase 3.  $\square$

**Lemma 16.** *Suppose the conditions of Subcase 3 hold (Lemma 15). Then  $\mathbf{x} = \mathbf{d}^{x/d}$ , except possibly for*

1.  $k + 2w \geq v$  **or**
2.  $k + 2w < v$  **and**
  - (a)  $v - u = h(u - w)$  **or**
  - (b)  $v - u = (h - \frac{1}{2})(u - w)$ ,

where  $d = \gcd(u, v, w)$  and  $h = \lfloor \frac{k+w}{u-w} \rfloor$ .

*Proof.* Note first that conditions 1 and 2 are simply reformulations of inequalities (3.4) and (3.5), respectively, found at the beginning of the proof of Lemma 15, where  $v - u$  and  $u - w$  replace  $u_1 + u_2$  and  $p$ , respectively. These inequalities constitute the two main cases considered in the proof, and so the result holds for condition 1.

Condition 2 however breaks down into cases (C1)-(C3). For both (C1) and (C3), it is shown that  $\mathbf{u}_1$  has period  $\ell$  (all symbols used as defined in the proof of Lemma 15). The same is true also for the various subcases of (C2), except when

- (a') the gap  $g = 0$  **or**
- (b')  $z < g$  and  $g_1 = g_2$ .

In all other cases, it is shown that  $\mathbf{u}_1 = \mathbf{kz}\ell = \mathbf{p}\ell$  has period  $\ell$ . Then, by Lemma 15,  $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$  has period  $d_1 = \gcd(p, u_1 + u_2)$ . Hence  $\mathbf{p}$  is a suffix of  $\mathbf{u}_1$ , thus a border of  $\mathbf{u}$ . Therefore  $\mathbf{u}^2$  has period  $d_1$ , as well as period  $u$ , so that by Lemma 4,  $\mathbf{u}^2$  has period  $\gcd(u, d_1) = \gcd(u, \gcd(u - w, v - u)) = \gcd(u, v, w) = d$ . This periodicity clearly extends to all of  $\mathbf{x}$ .

Next observe that in the proof of Lemma 15,  $g = u_1 + u_2 - ph$ , so that the condition  $g = 0$  given in (a') converts to the condition of (a) using the indicated substitutions for  $u_1 + u_2$  and  $p$ . Again from the proof of Lemma 15, we find that when  $z < g$ ,  $g_1 = k - g$ ,  $g_2 = g - z$ , from which we conclude that  $p = k + z = 2g$  in (b'). This in turn implies that  $h$  copies of  $\mathbf{p}$  ( $2h$  copies of  $\mathbf{g}$ ) cover  $\mathbf{u}_1\mathbf{u}_2\mathbf{g}$  of length  $v - u + g$ , from which (b) follows.  $\square$

## 3.2 Subcase 7

Here we give results for Subcase 7 corresponding to those for Subcase 3:

**Lemma 17** (Subcase 7). *Suppose that a string  $\mathbf{x}$  has prefixes  $\mathbf{u}^2$  and  $\mathbf{v}^2$ ,  $3u/2 < v < 2u$ , and suppose further that a third square  $\mathbf{w}^2$ ,  $w \neq u$ , occurs at position  $k + 1$  of  $\mathbf{x}$ , where*

$$u_1 < k < u_1 + u_2 < w < v \quad (3.8)$$

$$k + w \leq u + u_1 \quad (3.9)$$

$$k + 2w \leq 2u \quad (3.10)$$

and  $u_1 = 2u - v$  and  $u_2 = 2v - 3u$ . Then  $\mathbf{x} = \mathbf{d}_1^{u/d_1} \mathbf{d}_1^{v/d_1} \mathbf{d}_1^{(v-u)/d_1}$ , where  $d_1 = \gcd(u - w, v - u)$ .

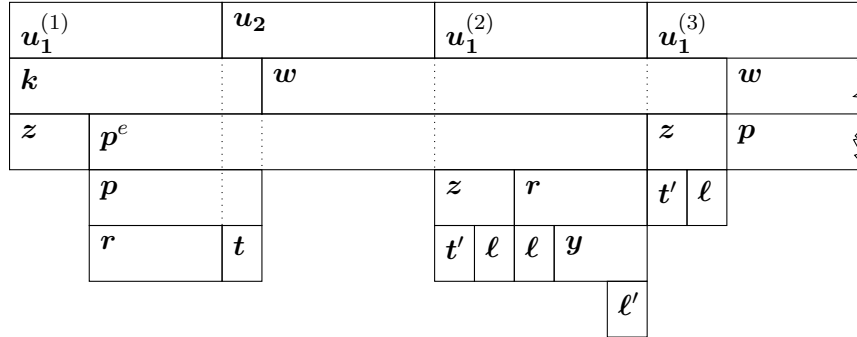
*Proof.* By Lemma 12, the overlap of  $\mathbf{u}^2$  and  $\mathbf{v}^2$  forces  $\mathbf{x} = (\mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_1\mathbf{u}_2)^2$ , with  $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$ . By Lemma 14(b),  $\mathbf{u} = \mathbf{z}\mathbf{p}^e$ , where  $\mathbf{z} = \mathbf{u}[1..k + w - u]$ ,  $\mathbf{p} = \mathbf{w}[1..u - w]$ , and  $e = 1 + \frac{u-k}{u-w}$ . See Figure 3.10.

We first show that if  $\mathbf{u}$  has period  $p = u - w$ , the lemma holds. Note that  $\mathbf{u}$  has period  $u_1 + u_2$  and

$$u_1 + u_2 + p = u + u_1 + u_2 - w < u$$

since  $u_1 + u_2 < w$  from (3.8). Assuming  $\mathbf{u}$  has period  $p$ ,  $\mathbf{u} = \mathbf{x}[1..u]$  has period  $d_1 = \gcd(p, u_1 + u_2)$  by Lemma 4. It follows that  $\mathbf{u}_1\mathbf{u}_2 = \mathbf{x}[u + v +$



Figure 3.10: String  $\mathbf{uu}_1$  in Subcase 7

$1..x]$ , a prefix of  $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$ , has period  $d_1$  as well. Finally,  $\mathbf{x}[u+1..u+v] = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$  has period  $u_1+u_2$  and prefix  $\mathbf{u}$  of length  $u > u_1+u_2$  with period  $d_1$ . Since  $d_1 = \gcd(u-w, u_1+u_2)$  divides  $u_1+u_2$ ,  $\mathbf{x}[u+1..u+v]$  has period  $d_1$  by Lemma 3. Thus the lemma holds assuming  $\mathbf{u}$  has period  $p$ .

We now embark on a demonstration that  $\mathbf{u}$  has period  $p$ . Notice (Figure 3.10) that  $\mathbf{u}_1\mathbf{u}_2$ ,  $\mathbf{k}$ , and  $\mathbf{zp}$  are prefixes of  $\mathbf{u}$ . Given that  $z = k - p$ , that  $k \in (u_1, u_1 + u_2)$  by (3.8), and that  $z \leq u_1$  by (3.9), we have

$$\mathbf{k} = \mathbf{zp} = \mathbf{zrt} = \mathbf{u}_1\mathbf{t}, \quad (3.11)$$

where  $\mathbf{r} = \mathbf{u}_1[z+1, u_1]$ , and  $\mathbf{t} = \mathbf{u}_2[1..k-u_1]$ .

Observe that

$$z - p = (k + w - u) - (u - w) = k - (2u - 2w) \leq 0$$

by (3.10), so that  $p \geq z$ . Also, by (3.8)

$$z = k + w - u > k + u_1 + u_2 - u = k - u_1 > 0,$$

so that in fact  $z \geq 2$ . Note further from (3.8) that

$$p = u - w < u - (u_1 + u_2) = u_1,$$

while from (3.10) and (3.8),

$$p = u - w \geq k/2 > u_1/2. \quad (3.12)$$

Thus  $u_1/2 < k/2 \leq p < u_1$ . Putting these inequalities together, we find

$$2 \leq z \leq p < u_1 = z + r, \quad (3.13)$$

from which we conclude that  $r > 0$ .

Also, since

$$z \leq p = r + t < u_1 = r + z,$$

and recalling from (3.8) that  $t = k - u_1 > 0$ , we see that  $0 < t < z$ , where since  $k = z + p \geq 2z$ ,  $z \leq k/2$ . Hence

$$0 < t < z \leq k/2 \leq p < u_1.$$

Let  $\mathbf{t}'$  be the prefix of  $\mathbf{z}$  of length  $t$ . Since  $t' = t < z$  and  $\mathbf{u}_1\mathbf{z}$  is a suffix of  $\mathbf{w}^{(1)}$ , there exists within  $\mathbf{w}$  a complete occurrence of  $\mathbf{u}_1\mathbf{t}'$ .

Since  $\mathbf{w}^{(1)}$  has prefix  $\mathbf{p}$ , so also does  $\mathbf{w}^{(2)}$ , with  $\mathbf{p} = \mathbf{rt}$ . Furthermore  $\mathbf{k} = \mathbf{zrt} = \mathbf{zrp}$ , so that  $\mathbf{p}$  is a suffix of  $\mathbf{k}$  with nonempty prefix  $\mathbf{r}$  that is a suffix of  $\mathbf{u}_1$ . Since  $\mathbf{u}_1$  is a proper substring of  $\mathbf{w}^{(1)}$  and  $p < u_1$ , it follows that  $\mathbf{u}_1$  has period  $p$ . In fact, the string

$$\mathbf{u}' = R_z(\mathbf{u}) = \mathbf{p}^e\mathbf{z} = \mathbf{pw} = \mathbf{ru}_2\mathbf{u}_1\mathbf{z}$$

has period  $p$ . Then  $\mathbf{u}_1 = \mathbf{zr}$ ,  $\mathbf{rz}$  and  $\mathbf{u}_1\mathbf{z} = \mathbf{zrz}$  all have period  $p$ .

Again by the periodicity of  $\mathbf{u}'$ , there exists a possibly empty  $\mathbf{y}'$  such that  $\mathbf{p}_1 = \mathbf{zy}'$  is a prefix of  $\mathbf{u}_1$ , and a  $\mathbf{y}$ , with  $y = y' = p - z$ , such that  $\mathbf{p}_2 = \mathbf{zy}$  is a suffix of  $\mathbf{u}_1$ , where  $\mathbf{p}_1$  and  $\mathbf{p}_2$  are both rotations of  $\mathbf{p}$ .

Now consider  $\mathbf{u}_1\mathbf{z}$ , a suffix of  $\mathbf{u}'$  with period  $p$ : this string has prefix  $\mathbf{zy}'\mathbf{z}$  and suffix  $\mathbf{zyz}$  which overlap each other by

$$\hat{p} = u_1 + z - 2u_1 + 2p = p + (p + z) - u_1 > p$$

positions. We may therefore conclude that all substrings of length  $p$  in  $\mathbf{zy}'\mathbf{z}$  and  $\mathbf{zyz}$  are rotations of each other. Then  $\mathbf{u}_1$  and  $R_z(\mathbf{u}_1)$  both have period  $p$ , and so, since  $\ell = u_1 - p = z - t < z$ , we can apply Lemma 9(a) (with  $(x, v, u) \sim (u_1, z, p)$ ) to conclude that  $R_t(\mathbf{u}_1)[1..2(z-t)]$  is a square of period  $\ell$ . Thus we may write  $\mathbf{u}_1 = \mathbf{t}'\ell^2 \dots$ , where  $\mathbf{t}'\ell = \mathbf{z}$ . In fact, since  $p = z + y = u_1 - \ell$ , so that  $u_1 = z + \ell + y$ , we find that  $\mathbf{u}_1 = \mathbf{t}'\ell^2\mathbf{y}$  with  $\mathbf{z} = \mathbf{t}'\ell$ ,  $\mathbf{r} = \ell\mathbf{y}$  and  $\mathbf{p} = \ell\mathbf{y}\mathbf{t}'$ .

Since  $\mathbf{u}' = \mathbf{pw}$  has period  $p$ ,  $\mathbf{w}$  is a prefix of  $\mathbf{u}'$ . As we see from Figure 3.10, this prefix  $\mathbf{w}$  ends distance  $r + t = \ell + y + t' = y + t' + \ell$  before the end of  $\mathbf{w}^{(1)}$ , from which we conclude that  $\mathbf{w}$  has suffix  $\mathbf{t}'\ell\ell$  as well as suffix  $\mathbf{z} = \mathbf{t}'\ell$ . Thus  $\mathbf{t}'$  is a border of  $\mathbf{z}$ . Now let  $\ell'$  be the prefix of  $\mathbf{z}$  of length  $\ell$ , so that  $\mathbf{z} = \mathbf{t}'\ell = \ell'\mathbf{t}'$  has period  $\ell$ . Note further that since  $\mathbf{w}$  has suffix  $\mathbf{y}\mathbf{z}$ , which in turn has suffix  $\mathbf{t}'\ell\ell = \ell'\mathbf{t}'\ell$ , therefore  $\ell'$  is a suffix of  $\mathbf{y}$ .

Assume  $t = t'$ . Then  $k = zrt = zrt'$  occurs in  $w$ , and as  $w$  has period  $p$ , so does  $k$ .  $u$  then has prefix  $k = zp$  and suffix  $p^e$  both of period  $p$  and both including  $p$ , so by Lemma 6,  $u$  has period  $p$ , as desired. Hence it will suffice to show  $t = t'$ .

From (3.13) it follows that a complete copy of  $p$  occurs  $h \geq 2$  times in  $u'$ . Several cases arise, based on the position of the suffix  $t$  of the  $h^{\text{th}}$  occurrence of  $p$ :

- (C1)  $t$  ends inside the prefix  $z$  of  $u_1^{(3)}$
- (C2)  $t$  is a substring of the suffix  $y$  of  $u_1^{(2)}$ , but  $t\ell$  is not.
- (C3)  $t\ell$  is a substring of the suffix  $y$  of  $u_1^{(2)}$ .
- (C4)  $t$  begins to the left of the suffix  $y$  of  $u_1^{(2)}$  and ends inside  $y$ .

We will see  $t = t'$  in all of these cases.

(C1)

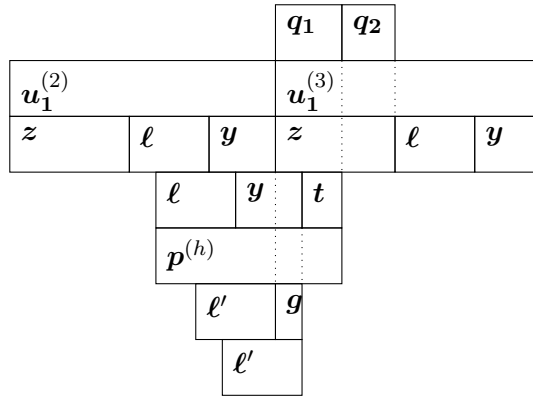


Figure 3.11: Subcase 7 when (C1) holds and  $g > 0$

Suppose first that (C1) holds, and write  $z = q_1q_2$ , where  $q_1$  is a nonempty suffix of  $p$  and, by the periodicity of  $u'$ ,  $q_2$  is a prefix of  $p$ .

We have shown that  $u' = p^h q_2$ , where  $p = \ell y t$ , and so  $u' = p^{h-1}(\ell y)(t)q_2$ . As in the proof of Lemma 15, we introduce the “gap”  $g = q_1 - t$ , a measure of the overlap between the prefix  $q_1$  of  $u_1^{(3)}$  and the suffix  $t$  of  $p^{(h)}$ . If  $g \geq 0$ ,

then  $\mathbf{t}$  is a substring of  $\mathbf{z}$ ; otherwise,  $\mathbf{t}$  ends inside  $\mathbf{z}$  but begins before it. Note that if  $g = 0$ , then  $\mathbf{q}_1 = \mathbf{t} = \mathbf{t}'$  and the remainder of the proof follows.

Suppose then that  $g > 0$  (Figure 3.11), so that  $\mathbf{q}_1 = \mathbf{g}\mathbf{t}$  for some string  $\mathbf{g}$  of length  $g$ . In this case, note that  $\ell'\mathbf{t}$  and  $\ell'\mathbf{t}'$  are substrings of  $\ell'\mathbf{z} = \ell'\ell'\mathbf{t}'$ , as we have seen of period  $\ell$ , and so both these strings also have period  $\ell$ , implying that  $\mathbf{t} = \mathbf{t}'$ , as required.

We now show further that for  $g > 0$ ,  $\mathbf{u}_1$  has period  $\gcd(g, \ell)$ . Since  $\mathbf{t}$  is a substring of  $\mathbf{z} = \mathbf{t}'\ell$ ,  $g \leq \ell$ . Therefore, since  $\ell\mathbf{y}$  is a suffix of  $\mathbf{u}_1$  and  $\mathbf{p} = \ell\mathbf{y}\mathbf{t}$ ,  $\ell\mathbf{y}$  and  $\ell$  both have period  $g$ , as does  $\ell'$ , since it is a suffix of  $\ell\mathbf{y}$ . Observe that  $\ell'\mathbf{g}$  has period  $\ell$  as a prefix of  $\ell'\mathbf{z}$ , as well as period  $g$  as a suffix of  $\ell\mathbf{y}$ , so that by Lemma 4,  $\ell'\mathbf{g}$  has period  $\gcd(g, \ell)$ .  $\mathbf{z}\ell$  then has period  $\ell$  and a substring  $\ell'$  of period  $\gcd(g, \ell) \mid \ell$ , so by Lemma 3,  $\mathbf{z}\ell$  has period  $\gcd(g, \ell)$ .  $\ell\mathbf{y}$  has period  $g$  and suffix  $\mathbf{g}$  of period  $\gcd(g, \ell)$ , so by Lemma 3, it has period  $\gcd(g, \ell)$ .  $\mathbf{z}\ell$  and  $\ell\mathbf{y}$  have period  $\gcd(g, \ell)$  and share substring  $\ell$ , so by Lemma 6,  $\mathbf{u}_1 = \mathbf{z}\ell\mathbf{y}$  has period  $\gcd(g, \ell)$ .

Suppose next that  $g < 0$ , so that  $\mathbf{t} = \mathbf{g}\mathbf{q}_1$  for some string  $\mathbf{g}$  of length  $|g|$ , as shown in Figure 3.12. Again  $\ell\mathbf{y}$  and  $\ell$  both have period  $|g|$ . If  $|g| \leq \ell$ , then  $\mathbf{t}\ell$  is a substring of  $\ell'\mathbf{z}$ , so  $\mathbf{t} = \mathbf{t}'$ . However, when  $g < 0$ , it is possible that  $|g| > \ell$ . In general, let  $\mathbf{g}'$  be the suffix of  $\mathbf{z}$  of length  $|g|$ . The suffix  $\mathbf{q}_2 = \ell\mathbf{g}'$  of  $\mathbf{z}$  has border  $\ell$  and thus period  $q_2 - \ell = |g|$ . It also has period  $\ell$  as a suffix of  $\mathbf{z}$ , so by Lemma 4, it has period  $\gcd(g, \ell)$ .  $\ell'\mathbf{z}$  then has period  $\ell$  and suffix  $\ell\mathbf{g}'$  of period  $\gcd(g, \ell) \mid \ell$ , so that by Lemma 3,  $\ell'\mathbf{z}$  has period  $\gcd(g, \ell)$ . Also by Lemma 3,  $\ell\mathbf{y}$  has period  $\gcd(g, \ell)$  since it has period  $|g|$  and, by the periodicity of  $\mathbf{u}'$ , substring  $\mathbf{g}'$  of period  $\gcd(g, \ell) \mid |g|$ . Both  $\ell\mathbf{y}$  and  $\ell'\mathbf{z}$  have period  $\gcd(g, \ell)$  and share substring  $\ell'$ , so that by Lemma 6,  $\ell\mathbf{y}\mathbf{z}$  has period  $\gcd(g, \ell)$ . Since  $\mathbf{t}\ell$  and  $\mathbf{t}'\ell$  are both substrings of  $\ell\mathbf{y}\mathbf{z}$ , therefore again  $\mathbf{t} = \mathbf{t}'$ .

Note finally that since  $\mathbf{z}\ell$  and  $\ell\mathbf{y}$  both have period  $\gcd(g, \ell)$  and share substring  $\ell$ , therefore by Lemma 6,  $\mathbf{u}_1 = \mathbf{z}\ell\mathbf{y}$  again has period  $\gcd(g, \ell)$ , as it did also for  $g > 0$ .

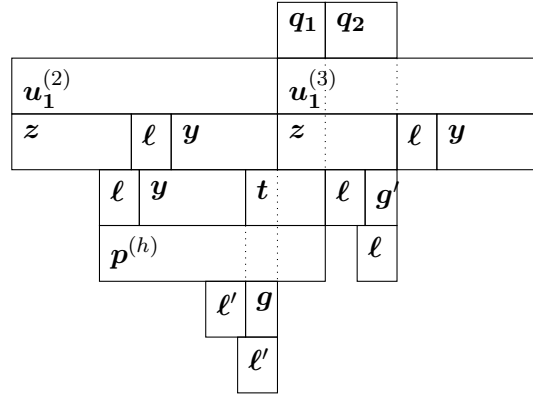


Figure 3.12: Subcase 7 when (C1) holds and  $g < 0$

(C2)

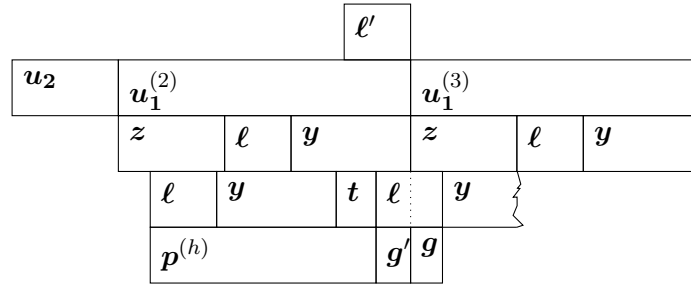


Figure 3.13: Subcase 7 when (C2) holds

Suppose (C2) holds; that is,  $t$  is a substring of  $y$  and ends within distance  $\ell$  of the end of  $y$ . By the periodicity of  $u'$ , a prefix of  $p = \ell y t$  follows  $p^{(h)}$ . Let  $g$  be the suffix of  $\ell$  that overlaps  $u_1^{(3)}$ , and let  $g'$  be the possibly empty prefix of  $\ell$  such that  $\ell = g'g$ .  $\ell y$  has period  $t + g'$  because the suffix  $\ell y$  of  $u_1^{(2)}$  and the prefix  $\ell y$  of  $p^{(h)}$  are offset by length  $t + g'$ .  $g'z$  is a prefix of  $\ell y$  since, by assumption,  $y \geq t + g' = z - g$ .

$g'z$  has period  $t + g'$  as a prefix of  $\ell y$  and period  $\ell$  as a suffix of  $\ell'z$ , so by lemma 4 it has period  $\gcd(t + g', \ell)$ . Thus  $\ell y$  has period  $t + g'$ , and a prefix  $g'z$  of period  $\gcd(t + g', \ell) \mid t + g'$ , so  $\ell y$  has period  $\gcd(t + g', \ell)$  by Lemma 3. Since  $z\ell$  has period  $\ell$  and prefix  $z$  of period  $\gcd(t + g', \ell) \mid \ell$ , therefore by Lemma 3,  $z\ell$  has period  $\gcd(t + g', \ell)$ .  $z\ell$  and  $\ell y$  have period  $\gcd(t + g', \ell)$  and share substring  $\ell$ , so by Lemma 6,  $u_1 = z\ell y$  has period  $\gcd(t + g', \ell)$ .

Remarking that  $\ell't$  and  $\ell't'$  are both substrings of  $\mathbf{u}_1$ , we again conclude that  $\mathbf{t} = \mathbf{t}'$ .

(C3)

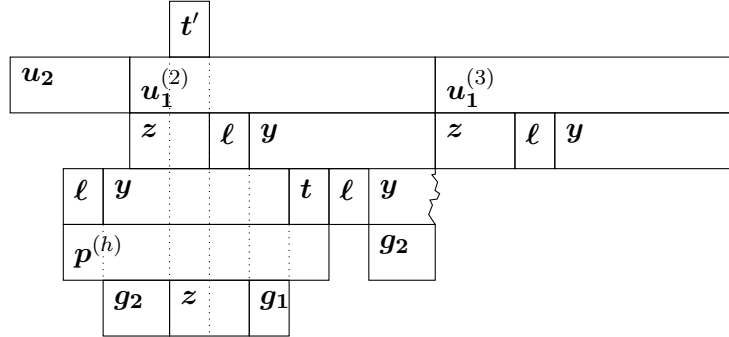


Figure 3.14: Subcase 7 when (C3) holds

Suppose (C3) holds; that is,  $\mathbf{t}\ell$  is a substring of  $\mathbf{y}$ . In this case,  $\mathbf{z} = \mathbf{t}'\ell$  is also a substring of  $\mathbf{y}$  because  $y \geq t + \ell$  and the prefix  $\ell\mathbf{y}$  of  $\mathbf{p}^{(h)}$  ends at least  $z$  and at most  $y$  positions from the end of  $\mathbf{u}_1^{(2)} = \ell't'\ell\mathbf{y}$ .

Let  $\mathbf{g}_1$  and  $\mathbf{g}_2$  be the (possibly empty) substrings of  $\mathbf{y}$  immediately before and after  $\mathbf{t}\ell$  such that  $\mathbf{y} = \mathbf{g}_1\mathbf{t}\ell\mathbf{g}_2$ . Since  $\mathbf{u}_1\mathbf{z}$  has period  $p$ ,  $\mathbf{g}_1$  and  $\mathbf{g}_2$  are borders of  $\mathbf{y}$  such that  $\mathbf{y} = \mathbf{g}_1\mathbf{t}\ell\mathbf{g}_2 = \mathbf{g}_2\mathbf{t}'\ell\mathbf{g}_1$ .

Recall that  $\ell'$  is a suffix of  $\mathbf{y}$ , so  $\ell'$  is a suffix of  $\ell\mathbf{g}_2$ . Since the prefix  $\ell\mathbf{g}_2$  of  $\mathbf{p}$  occurs before the substring  $\mathbf{z}$  of  $\mathbf{y}$ ,  $\ell'$  also occurs before  $\mathbf{z}$ .

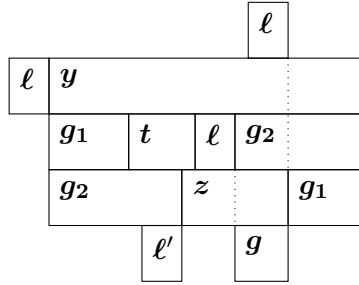
If  $g_1 = g_2$ , then  $\mathbf{t}$  and  $\mathbf{t}'$  occur at the same positions in two copies of  $\mathbf{y}$ , so that  $\mathbf{t} = \mathbf{t}'$ . If  $g_1 \neq g_2$ , several cases arise:

1.  $g_1 < g_2$  ( $g = g_2 - g_1$ )

Let  $g = g_2 - g_1$ , and let  $\mathbf{g}$  be the nonempty prefix of  $\mathbf{g}_2$  such that  $\mathbf{g}_2 = \mathbf{g}\mathbf{g}_1$ .  $\mathbf{g}_1$  is a border of  $\mathbf{g}_2$ , so that  $\mathbf{g}_2$  has period  $g$ .

- (a)  $g \leq z$  (Figure 3.15)

The proof requires several steps:

Figure 3.15: Subcase 7 when (C3) holds,  $g_1 < g_2$ , and  $g \leq z$ 

- As a suffix of  $l'z$ ,  $lg$  has period  $l$  and suffix  $l$ , hence border  $l$  and period  $g$ , therefore by Lemma 4 period  $\gcd(g, l)$ .
- Then  $l'z$  has period  $l$  and a suffix  $lg$  of period  $\gcd(g, l) \mid l$ , so by Lemma 3,  $l'z$  has period  $\gcd(g, l)$ .
- Since  $g_2$  has period  $g$  and prefix  $g$  of period  $\gcd(g, l) \mid g$ , therefore by Lemma 3,  $g_2$  has period  $\gcd(g, l)$ .
- Since prefix  $g_2$  of  $y$  and  $l'z$  have period  $\gcd(g, l)$  and share substring  $l'$ , therefore  $g_2z$  has period  $\gcd(g, l)$  by Lemma 6.
- Thus  $t\ell$  and  $t'\ell$  both have period  $\gcd(g, l)$  as substrings of  $g_2z$ , implying that  $t = t'$ .

Note that  $g_2z$  and  $g_2$  have period  $\gcd(g, l)$  and share substring  $g$ , so that by Lemma 6,  $y = g_2zg_1$  has period  $\gcd(g, l)$ . Since  $lg_2$  has period  $\gcd(g, l)$  as a substring of  $y$ , and since  $g_2$  is a prefix of  $y$ , therefore by Lemma 6,  $ly$  has period  $\gcd(g, l)$ .  $z\ell = l'z$  and  $ly$  have period  $\gcd(g, l)$  and share substring  $l$ , so by Lemma 6,  $u_1 = zly$  has period  $\gcd(g, l)$ .

(b)  $g > z$

$y$  has period  $y - g_2 = z + g_1$ , so  $y$  has a prefix  $g_1t\ell g'_2 = g'_2zg_1$ , where  $g'_2$  is a prefix of  $g_2$  and  $|g'_2 - g_1| \leq z$ , so one of cases 1(a) and 2(a) applies.

2.  $g_1 > g_2$  ( $g = g_1 - g_2$ )

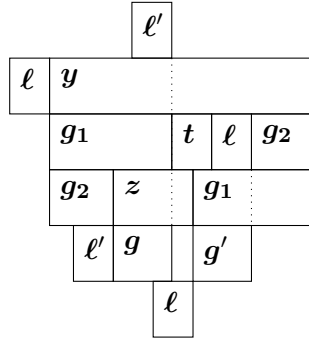


Figure 3.16: Subcase 7 when (C3) holds,  $g_1 > g_2$  and  $g \leq z$

Let  $g = g' = g_1 - g_2$ , let  $\mathbf{g}$  be the nonempty suffix of  $\mathbf{g}_1$  such that  $\mathbf{g}_1 = \mathbf{g}_2\mathbf{g}$ , and let  $\mathbf{g}'$  be the nonempty prefix of  $\mathbf{g}_1$  such that  $\mathbf{g}_1 = \mathbf{g}'\mathbf{g}_2$ . Since  $\mathbf{g}_2$  is a border of  $\mathbf{g}_1$ , therefore  $\mathbf{g}_1$  and  $\mathbf{g}_2$  have period  $g = g'$ .

(a)  $g \leq z$  (Figure 3.16)

Again there are several steps:

- $\ell'\mathbf{g}$  has period  $\ell$  as a prefix of  $\ell'z$  and shares suffix  $\ell'$  with  $\ell\mathbf{g}_1$ ; accordingly,  $\ell'\mathbf{g}$  has border  $\ell'$  and period  $g$ , hence periods  $g$  and  $\ell$ , thus by Lemma 4 period  $\gcd(g, \ell)$ .
- $\ell'z$  has period  $\ell$  and prefix  $\ell'\mathbf{g}$  of period  $\gcd(g, \ell) \mid \ell$ , hence by Lemma 3, it also has period  $\gcd(g, \ell)$ .
- $\mathbf{g}_1$  has period  $g$  and suffix  $\mathbf{g}$  of period  $\gcd(g, \ell) \mid g$ , so again by Lemma 3, it also has period  $\gcd(g, \ell)$ .
- $\mathbf{g}_1$  and  $\ell'z$  have period  $\gcd(g, \ell)$  and share substring  $\mathbf{g}$ , so that by Lemma 6,  $\mathbf{g}_2z$  has period  $\gcd(g, \ell)$ .
- Prefix  $\ell\mathbf{g}'$  of  $\ell\mathbf{g}_1$  shares suffix  $\ell$  with  $t\ell$ , so  $\ell\mathbf{g}'$  has border  $\ell$  and period  $g$ . Moreover  $\ell\mathbf{g}'$  has suffix  $\mathbf{g}'$  which, as a prefix of  $\mathbf{g}_1$ , has period  $\gcd(g, \ell) \mid g$ , implying that  $\ell\mathbf{g}'$  also has period  $\gcd(g, \ell)$  by Lemma 3.
- $\mathbf{g}_2z$  and  $\ell\mathbf{g}'$  have period  $\gcd(g, \ell)$  and share substring  $\ell$ , so by Lemma 6,  $\mathbf{g}_2z\mathbf{g}'$  has period  $\gcd(g, \ell)$ .



- Then  $g_2zg'$  and  $g_1$  have period  $\gcd(g, \ell)$  and share substring  $g'$ , so that by Lemma 6 the entire string  $y = g_2zg_1$  has period  $\gcd(g, \ell)$ .
- Therefore  $t\ell$  and  $t'\ell$  both have period  $\gcd(g, \ell)$  as substrings of  $y$ , so that  $t = t'$ , as required.

Since, as a substring of  $y$ ,  $\ell g_1$  has period  $\gcd(g, \ell)$ , and since  $g_1$  is also a prefix of  $y$ , it follows from Lemma 6 that  $\ell y$  has period  $\gcd(g, \ell)$ . Note then that  $z\ell = \ell'z$  and  $\ell y$  have period  $\gcd(g, \ell)$  and share substring  $\ell$ , so that by Lemma 6  $u_1 = z\ell y$  has period  $\gcd(g, \ell)$ .

(b)  $g > z$

$y$  has period  $y - g_1 = g_2 + z$ , so  $y$  has a prefix  $g_2zg'_1 = g'_1t\ell g_2$ , where  $g'_1$  is a prefix of  $g_1$  and  $|g'_1 - g_2| \leq z$ , so one cases 1(a) and 2(a) applies.

(C4)

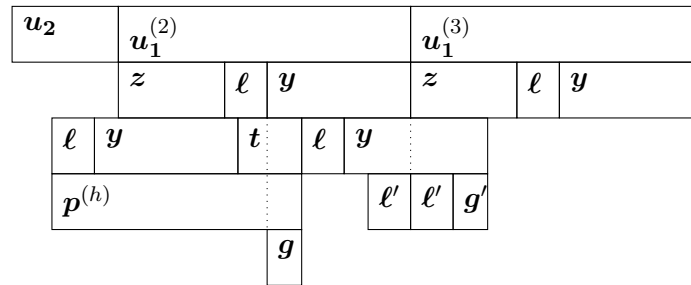


Figure 3.17: Subcase 7 when (C4) holds

Suppose (C4) holds; that is,  $t$  begins to the left of  $y$  and ends inside it. Let  $g$  be the suffix of  $t$  that is also a prefix of  $y$ . Let  $g'$  be the suffix of  $y$  of length  $g' = g$ . By the periodicity of  $u'$ , a copy of  $\ell y$  follows  $t$ , extending  $\ell + g$  positions into  $u_1^{(3)}$ .  $\ell'$  is a suffix of  $\ell y$ , so  $\ell' \ell' g'$  is a suffix of  $\ell y$ .

The suffix  $\ell\mathbf{y}$  of  $\mathbf{u}_1^{(2)}$  and the occurrence of  $\ell\mathbf{y}$  that follows  $\mathbf{t}$  are offset by length  $\ell+g$ , so  $\ell\mathbf{y}$  has period  $\ell+g$ . Since  $\ell'\ell'g'$  has period  $\ell+g$  as a suffix of  $\ell\mathbf{y}$  and period  $\ell$  as a prefix of  $\ell'\mathbf{z}$ , it therefore has period  $\gcd(\ell+g, \ell) = \gcd(g, \ell)$  by Lemma 4.  $\ell\mathbf{y}$  has period  $\ell+g$  and suffix  $\ell'\ell'g'$  of period  $\gcd(g, \ell) \mid \ell+g$ , implying that it has period  $\gcd(g, \ell)$  by Lemma 3.  $\mathbf{z}\ell$  has period  $\ell$  and substring  $\ell$  of period  $\gcd(g, \ell) \mid \ell$ , so by Lemma 3, it has period  $\gcd(g, \ell)$ .  $\mathbf{z}\ell$  and  $\ell\mathbf{y}$  have period  $\gcd(g, \ell)$  and share substring  $\ell$ , so by Lemma 6,  $\mathbf{u}_1 = \mathbf{z}\ell\mathbf{y}$  has period  $\gcd(g, \ell)$ . Since  $\mathbf{t}\ell$  and  $\mathbf{t}'\ell$  are substrings of  $\mathbf{u}_1$ , we conclude finally that  $\mathbf{t} = \mathbf{t}'$ .

This completes the proof of Subcase 7.  $\square$

**Lemma 18.** *Suppose the conditions of Subcase 7 hold (Lemma 17). Then  $\mathbf{x} = \mathbf{d}^{x/d}$ , except possibly for*

$$(a) \ v - u = (h - 1)(u - w) \ \mathbf{or}$$

$$(b) \ v - u = (h - \frac{1}{2})(u - w),$$

where  $d = \gcd(u, v, w)$  and  $h = \lfloor \frac{u}{p} \rfloor$ .

*Proof.* (All symbols used as defined in the proof of Lemma 17; refer to Figure 3.10.) Notice that in the proof of the lemma,  $\mathbf{u}_1$  has period  $\ell$  in all cases except when

(a') (C1) holds and  $g = 0$  **and**

(b') (C3) holds and  $g_1 = g_2$ .

Suppose then that  $\mathbf{u}_1 = \mathbf{z}\mathbf{r} = \mathbf{z}\ell\mathbf{y}$  does indeed have period  $\ell$ . Since  $\mathbf{z} = \ell'\mathbf{t}'$  is a prefix of  $\mathbf{u}_1$ , it follows that  $\mathbf{t}'\mathbf{r} = \mathbf{t}'\ell\mathbf{y}$  is a suffix of  $\mathbf{u}_1$ . Since  $\mathbf{u}_1$  has period  $\ell$ ,  $\mathbf{r} = \ell\mathbf{y}$  has prefix  $\mathbf{y}$ , and so  $\mathbf{u}_1$  has border  $\mathbf{t}'\ell\mathbf{y}$  of length  $p = t + \ell + y$ , as therefore  $\mathbf{u}$  does also. By Lemma 17,  $\mathbf{u}$  has period  $p$ , so that by Lemma 15,  $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$  has period  $d_1 = \gcd(p, u_1 + u_2) = \gcd(u - w, v - u)$ . Since  $\mathbf{u}$  has a border of length  $p$ , it follows that  $\mathbf{u}^2$  also has period  $d_1$ , as well as period  $u$ , so that by Lemma 4,  $\mathbf{u}^2$  has period  $\gcd(u, d_1) = \gcd(u, \gcd(u - w, v - u)) = \gcd(u, v, w) = d$ . This periodicity clearly extends to all of  $\mathbf{x}$ .

Now consider the exceptional cases. For (a'), recall that in (C1) the gap  $g$  is the difference between the two prefixes of  $\mathbf{x}$ ,  $\mathbf{z}\mathbf{p}^h$  and  $\mathbf{u}$ , where  $p = u - w$ , so that  $g = 0$  implies  $hp + z = u + t$ . Substituting  $z = k + w - u$ ,  $t = k - u_1$  yields

$$h(u - w) = u - k + (u - w) + k - u_1,$$

from which, with a little manipulation, (a) follows. For (b'), from Figure 3.14  $g_1 = g_2$  in (C3) implies

$$z + ph - t - (u_1 + u_2 + z + \ell) = u - (z + ph + \ell),$$

which since  $z = \ell + t$  and  $\ell = u_1 - p$  becomes

$$2ph = u + u_1 + u_2 - u_1 + p.$$

A bit more manipulation yields (b), completing the proof. □

# Chapter 4

## On to the General Case

In this thesis, we have proven the last two remaining subcases of the New Periodicity Lemma, which describes the regularity that must result from three overlapping squares of which two begin at the same position and the third begins to the right. Future work should generalize this result to cases in which three squares occur close to each other, but with no two of them necessarily at the same position. As explained in Section 2.5.2, Lemma 14 provides six cases (*a–f*) covering all possible configurations of two overlapping squares. This allows the characterization of any instance of the general case of three overlapping squares  $\mathbf{u}^2, \mathbf{v}^2, \mathbf{w}^2$  as a pair  $[ij]$ ,  $i, j \in [a..f]$ , in which case  $i$  corresponds to the overlap of  $\mathbf{u}^2$  with  $\mathbf{v}^2$ ,  $j$  the overlap of  $\mathbf{v}^2$  with  $\mathbf{w}^2$ .

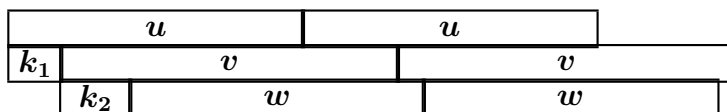


Figure 4.1:  $\mathbf{u}^2$  overlapping  $\mathbf{v}^2$  (case (d)) that in turn overlaps  $\mathbf{w}^2$  (case (b)): what is the combined effect?

As an example, we present a lemma that combines cases *d* and *b* (illustrated in Figure 4.1). According to our computational experiments, this lemma applies to about three-quarters of cases in which the maximum alphabet size  $\sigma = \gcd(u, v, w)$ .

**Lemma 19.** *In case  $[db]$ , if  $k_2 \leq 2u - v - k_1 + d_1$ , then  $\mathbf{x}$  has period  $d$ , where  $d = \gcd(u, v, w)$  and  $d_1 = \gcd(v - u, v - w) = \gcd(u - w, v - w)$ .*

*Proof.* We use subscripts  $_1$  to identify variables for  $\mathbf{u}$  and  $\mathbf{v}$ , subscripts  $_2$  for those of  $\mathbf{v}$  and  $\mathbf{w}$ . Observe then that for  $e_1 > 1$ ,  $e_2 \geq 1$ ,

$$\mathbf{v} = \mathbf{p}_1^{e_1} \mathbf{k}_1 \mathbf{p}_1 = \mathbf{z}_2 \mathbf{p}_2^{e_2},$$

where the variables subscripted  $_1$  relate to case (d) of Lemma 14, those subscripted  $_2$  to case (b). The substring  $\mathbf{v}' = \mathbf{v}[z_2 + 1, v - k_1 - p_1]$  has two periods,  $p_1 = v - u$  and  $p_2 = v - w$ . To apply Lemma 4, we must have

$$\begin{aligned} p_1 + p_2 - \gcd(p_1, p_2) &\leq v - k_1 - p_1 - z_2 \\ k_2 &\leq 2u - v - k_1 + d_1 \end{aligned}$$

Thus if  $k_2 \leq 2u - v - k_1 + d_1$ , then  $\mathbf{v}'$  has period  $d_1$ . Moreover,  $\mathbf{v}$  has a prefix of period  $p_1$  that includes  $\mathbf{v}'$ , as well as a suffix of period  $p_2$  that includes  $\mathbf{v}'$ , so  $\mathbf{v}$  itself has period  $d_1$ . Since  $\mathbf{p}_1$  is a border of  $\mathbf{v}$ ,  $\mathbf{v}$  is a repetition of period  $d_1$ . Because  $\mathbf{u}$  is a substring of  $\mathbf{v}^2$ ,  $\mathbf{u}$  has period  $d_1$ . Therefore  $\mathbf{x}$  has prefix  $\mathbf{u}$  and suffix  $\mathbf{v}^2$  both of period  $d_1$  that include  $\mathbf{v}[1..u - k_1]$ . In case (d),  $\frac{2(k_1 + v)}{3} \leq u$  which implies  $d_1 \leq u - k_1$ , so  $\mathbf{x}$  has period  $d_1$ . The substrings  $\mathbf{u}^2$  and  $\mathbf{w}^2$  then have periods  $\gcd(d_1, u)$  and  $\gcd(d_1, w)$ , respectively, and so  $\mathbf{x}$  itself has those periods. Finally,  $\mathbf{x}$  has period  $\gcd(d_1, u, w) = d$ .  $\square$

The preceding lemma is a sample of the combinatorial information that may be obtained from considering all cases  $[ij]$  as specified above. Much work remains to be done to state and prove similar results for all 36  $[ij]$  pairs. To date, all the results given for the New Periodicity Lemma (Crochemore and Rytter 1995; Fan et al. 2006; Franek, Fuller, et al. 2012; Kopylova and Smyth 2012; Simpson 2007) deal only with the special cases  $[ij]$ ,  $i \in \{a, d\}$ , that arise for  $k_1 = 0$ . Kopylova and Smyth (2012) used computer simulations for small values of  $k, u, v, w$  to help generate conjectures, and it seems that similar techniques can profitably be used for the general case.

Once the combinatorics of overlapping squares is well understood, we may find new and more combinatorially-sophisticated approaches to the bounding of the number of runs in any string of given length. Moreover, it may be possible to design an algorithmic approach to the computation of runs in a manner consistent with their sparseness of occurrence.

# Bibliography

- [1] Mohamed Ibrahim Abouelhoda, Stefan Kurtz, and Enno Ohlebusch. Replacing suffix trees with enhanced suffix arrays. *Journal of Discrete Algorithms*, 2(1):53–86, 2004 (cited on page 5).
- [2] Alberto Apostolico. The myriad virtues of subword trees. In, *Combinatorial Algorithms on Words*, pages 85–96. Springer, 1985 (cited on page 4).
- [3] Alberto Apostolico and Franco P. Preparata. Optimal off-line detection of repetitions in a string. *Theoretical Computer Science*, 22(3):297–315, 1983 (cited on page 9).
- [4] Widmer Bland and W.F. Smyth. Three overlapping squares: The general case characterized & applications, 2014. Submitted (cited on pages 14, 16).
- [5] M. Gabriella Castelli, Filippo Mignosi, and Antonio Restivo. Fine and Wilf’s theorem for three periods and a generalization of Sturmian words. *Theoretical Computer Science*, 218(1):83–94, 1999 (cited on page 7).
- [6] Gang Chen, Simon J. Puglisi, and William F. Smyth. Fast and practical algorithms for computing all the runs in a string. In *Combinatorial Pattern Matching*, 2007, pages 307–315 (cited on page 9).
- [7] Sorin Constantinescu and Lucian Ilie. Generalised Fine and Wilf’s theorem for arbitrary number of periods. *Theoretical Computer Science*, 339(1):49–60, 2005 (cited on page 7).
- [8] Max Crochemore. An optimal algorithm for computing the repetitions in a word. *Information Processing Letters*, 12(5):244–250, 1981 (cited on page 9).

- [9] Maxime Crochemore and Lucian Ilie. Maximal repetitions in strings. *Journal of Computer and System Sciences*, 74(5):796–807, 2008 (cited on page 10).
- [10] Maxime Crochemore, Lucian Ilie, and Liviu Tinta. The “runs” conjecture. *Theoretical Computer Science*, 412(27):2931–2941, 2011 (cited on pages 1, 10).
- [11] Maxime Crochemore, Lucian Ilie, and Liviu Tinta. Towards a solution to the “runs” conjecture. In *Combinatorial Pattern Matching*, 2008, pages 290–302 (cited on page 10).
- [12] Maxime Crochemore and Wojciech Rytter. Squares, cubes, and time-space efficient string searching. *Algorithmica*, 13(5):405–425, 1995 (cited on pages i, 1, 11, 39).
- [13] Kangmin Fan, Simon J. Puglisi, William F. Smyth, and Andrew Turpin. A new periodicity lemma. *SIAM Journal on Discrete Mathematics*, 20(3):656–668, 2006 (cited on pages i, 1, 11–13, 39).
- [14] Martin Farach. Optimal suffix tree construction with large alphabets. In *Proceedings of the 38<sup>th</sup> Annual Symposium on Foundations of Computer Science*. IEEE, 1997, pages 137–143 (cited on page 4).
- [15] Nathan J. Fine and Herbert S. Wilf. Uniqueness theorems for periodic functions. *Proceedings of the American Mathematical Society*, 16(1):109–114, 1965 (cited on page 6).
- [16] Aviezri S. Fraenkel and Jamie Simpson. An extension of the periodicity lemma to longer periods. *Discrete Applied Mathematics*, 146(2):146–155, 2005 (cited on page 7).
- [17] Aviezri S. Fraenkel and Jamie Simpson. The exact number of squares in Fibonacci words. *Theoretical Computer Science*, 218(1):95–106, 1999 (cited on page 9).
- [18] Frantisek Franek, Robert C.G. Fuller, Jamie Simpson, and William F. Smyth. More results on overlapping squares. *Journal of Discrete Algorithms*, 17:2–8, 2012 (cited on pages i, 1, 12, 14, 39).
- [19] Frantisek Franek, R.J. Simpson, and William F. Smyth. The maximum number of runs in a string. In *Proceedings 14<sup>th</sup> Australasian Workshop on Combinatorial Algorithms*, 2003, pages 26–35 (cited on page 10).

- [20] Edward Fredkin. Trie memory. *Communications of the ACM*, 3(9):490–499, 1960 (cited on page 4).
- [21] Robert Giegerich and Stefan Kurtz. From ukkonen to mcreight and weiner: a unifying view of linear-time suffix tree construction. *Algorithmica*, 19(3):331–353, 1997 (cited on page 4).
- [22] Mathieu Giraud. Asymptotic behavior of the numbers of runs and microruns. *Information and Computation*, 207(11):1221–1228, 2009 (cited on page 10).
- [23] Mathieu Giraud. Not so many runs in strings. In, *Language and Automata Theory and Applications*, pages 232–239, 2008 (cited on page 10).
- [24] Dan Gusfield. *Algorithms on Strings, Trees and Sequences: Computer science and computational biology*. Cambridge University Press, 1997 (cited on page 4).
- [25] Anisa Al-Hafeedh, Maxime Crochemore, Lucian Ilie, Evguenia Kopylova, William F. Smyth, German Tischler, and Munina Yusufu. A comparison of index-based Lempel-Ziv LZ77 factorization algorithms. *ACM Computing Surveys*, 45(1):5, 2012 (cited on page 5).
- [26] Štěpán Holub. On multiperiodic words. *RAIRO — Theoretical Informatics and Applications*, 40(04):583–591, 2006 (cited on page 7).
- [27] Costas S. Iliopoulos, Dennis Moore, and William F. Smyth. A characterization of the squares in a Fibonacci string. *Theoretical Computer Science*, 172(1):281–291, 1997 (cited on page 9).
- [28] Jacques Justin. On a paper by Castelli, Mignosi, Restivo. *RAIRO — Theoretical Informatics and Applications*, 34(5):373–377, 2000 (cited on page 7).
- [29] Juha Kärkkäinen and Peter Sanders. Simple linear work suffix array construction. In, *Automata, Languages and Programming*, pages 943–955. Springer, 2003 (cited on page 5).
- [30] Dong Kyue Kim, Jeong Seop Sim, Heejin Park, and Kunsoo Park. Linear-time construction of suffix arrays. In *Combinatorial Pattern Matching*. Springer, 2003, pages 186–199 (cited on page 5).
- [31] Pang Ko and Srinivas Aluru. Space efficient linear time construction of suffix arrays. In *Combinatorial Pattern Matching*. Springer, 2003, pages 200–210 (cited on page 5).



- [32] Roman Kolpakov and Gregory Kucherov. On maximal repetitions in words. *Journal on Discrete Algorithms*, 1(1):159–186, 2000 (cited on pages 1, 9, 10).
- [33] Evguenia Kopylov. Computing repetitions in strings: Current algorithms & the combinatorics of future ones. McMaster University, 2010 (cited on page 9).
- [34] Evguenia Kopylova and William F. Smyth. The three squares lemma revisited. *Journal of Discrete Algorithms*, 11:3–14, 2012 (cited on pages i, 1, 8, 12, 39).
- [35] Stefan Kurtz. Reducing the space requirement of suffix trees. *Software—Practice and Experience*, 29(13):1149–71, 1999 (cited on page 4).
- [36] Kazuhiko Kusano, Kazuyuki Narisawa, and Ayumi Shinohara. On morphisms generating run-rich strings. In *Prague Stringology Conference 2013*, 2013, page 35 (cited on pages 1, 10).
- [37] Abraham Lempel and Jacob Ziv. On the complexity of finite sequences. *IEEE Transactions on Information Theory*, 22(1):75–81, 1976 (cited on page 5).
- [38] M. Lothaire. *Algebraic Combinatorics on Words*. Cambridge University Press, 2002 (cited on pages 6, 7).
- [39] M. Lothaire. *Applied Combinatorics on Words*. Cambridge University Press, 2005 (cited on page 6).
- [40] Michael G. Main. Detecting leftmost maximal periodicities. *Discrete Applied Mathematics*, 25(1):145–153, 1989 (cited on pages 8, 9).
- [41] Michael G. Main and Richard J. Lorentz. An  $\mathcal{O}(n \log n)$  algorithm for finding all repetitions in a string. *Journal of Algorithms*, 5(3):422–432, 1984 (cited on page 9).
- [42] Udi Manber and Gene Myers. Suffix arrays: A new method for on-line string searches. *SIAM Journal on Computing*, 22(5):935–948, 1993 (cited on page 5).
- [43] Michael A. Maniscalco and Simon J. Puglisi. An efficient, versatile approach to suffix sorting. *Journal of Experimental Algorithmics*, 12:1–2, 2008 (cited on page 5).

- [44] Michael A. Maniscalco and Simon J. Puglisi. Faster lightweight suffix array construction. In *Proceedings of the 17<sup>th</sup> International Workshop on Combinatorial Algorithms*, 2006, pages 16–29 (cited on page 5).
- [45] Giovanni Manzini and Paolo Ferragina. Engineering a lightweight suffix array construction algorithm. *Algorithmica*, 40(1):33–50, 2004 (cited on page 5).
- [46] Wataru Matsubara, Kazuhiko Kusano, Akira Ishino, Hideo Bannai, and Ayumi Shinohara. New lower bounds for the maximum number of runs in a string. In *Prague Stringology Conference 2008*, 2008, page 140 (cited on page 10).
- [47] Edward M. McCreight. A space-economical suffix tree construction algorithm. *Journal of the ACM*, 23(2):262–272, 1976 (cited on page 4).
- [48] Ge Nong, Sen Zhang, and Wai Hong Chan. Linear time suffix array construction using d-critical substrings. In *Combinatorial Pattern Matching*. Springer, 2009, pages 54–67 (cited on page 5).
- [49] Simon J. Puglisi and Jamie Simpson. The expected number of runs in a word. *Australasian Journal of Combinatorics*, 42:45–54, 2008 (cited on pages 2, 9).
- [50] Simon J. Puglisi, Jamie Simpson, and William F. Smyth. How many runs can a string contain? *Theoretical Computer Science*, 401(1):165–171, 2008 (cited on page 10).
- [51] Simon J. Puglisi, William F. Smyth, and Andrew H. Turpin. A taxonomy of suffix array construction algorithms. *ACM Computing Surveys*, 39(2):4, 2007 (cited on page 5).
- [52] Wojciech Rytter. The number of runs in a string. *Information and Computation*, 205(9):1459–1469, 2007 (cited on page 10).
- [53] Wojciech Rytter. The number of runs in a string: Improved analysis of the linear upper bound. In *Proceedings of the 23<sup>rd</sup> Symposium on Theoretical Aspects of Computer Science*. Volume 3884, 2006, page 184 (cited on page 10).
- [54] Jamie Simpson. Intersecting periodic words. *Theoretical Computer Science*, 374(1):58–65, 2007 (cited on pages i, 1, 7, 12, 39).

- [55] Jamie Simpson. Modified Padovan words and the maximum number of runs in a word. *Australasian Journal of Combinatorics*, 46:129–145, 2010 (cited on pages 1, 10).
- [56] R.J. Simpson and Robert Tijdeman. Multi-dimensional versions of a theorem of Fine and Wilf and a formula of Sylvester. *Proceedings of the American Mathematical Society*, 131(6):1661–1671, 2003 (cited on page 7).
- [57] Bill Smyth. *Computing Patterns in Strings*. Pearson Education, 2003 (cited on pages 3, 8, 9).
- [58] Robert Tijdeman and Luca Q. Zamboni. Fine and Wilf words for any periods II. *Theoretical Computer Science*, 410(30):3027–3034, 2009 (cited on page 7).
- [59] Esko Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995 (cited on page 4).
- [60] Peter Weiner. Linear pattern matching algorithms. In *Proceedings of the 14<sup>th</sup> Annual Symposium on Switching and Automata Theory*. IEEE, 1973, pages 1–11 (cited on page 4).