



Murdoch
UNIVERSITY

MURDOCH RESEARCH REPOSITORY

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

The definitive version is available at :

<http://dx.doi.org/10.1016/j.tcs.2015.06.037>

Bland, W. and Smyth, W.F. (2015) Three overlapping squares: The general case characterized & applications. Theoretical Computer Science . In Press.

<http://researchrepository.murdoch.edu.au/27491/>

Copyright: © 2015 Elsevier B.V.
It is posted here for your personal use. No further distribution is permitted.

Accepted Manuscript

Three overlapping squares: The general case characterized & applications

Widmer Bland, W.F. Smyth

PII: S0304-3975(15)00547-2
DOI: <http://dx.doi.org/10.1016/j.tcs.2015.06.037>
Reference: TCS 10294

To appear in: *Theoretical Computer Science*

Received date: 13 March 2014
Revised date: 15 June 2015
Accepted date: 18 June 2015

Please cite this article in press as: W. Bland, W.F. Smyth, Three overlapping squares: The general case characterized & applications, *Theoret. Comput. Sci.* (2015), <http://dx.doi.org/10.1016/j.tcs.2015.06.037>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



Three Overlapping Squares: The General Case Characterized & Applications

Widmer Bland¹ and W. F. Smyth^{1,2}

¹ Algorithms Research Group, Department of Computing & Software
McMaster University, Hamilton, Ontario, Canada L8S 4K1
smyth@mcmaster.ca

² School of Engineering & Information Technology
Murdoch University, Murdoch WA 6150, Australia

Abstract. The “Three Squares Lemma” [Crochemore & 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; Simpson, 2007; Kopylova & Smyth, 2012; Franek *et al.*, 2012] 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 — in these cases also, similar restrictions apply. In this paper an alternative strategy is proposed: the consequences of having only *two* squares at neighbouring positions are carefully analyzed, and then the observation is made that the analysis applies in a straightforward way (though perhaps with complicated details) to the three neighbouring squares problem in its full generality. We then apply these new insights, first to proofs of the final two remaining unproved subcases (out of a total of 14) of the NPL [Fan *et al.*, 2012], then to an instance of the more general problem.

Keywords: string, word, overlapping squares, repetition, run, maximal periodicity.

1 Introduction

Beginning with the “Three Squares Lemma” of Crochemore & Rytter [9], there has for several years been considerable interest in the limitations that may exist on periodicity in strings. An early survey of this topic by Mignosi & Restivo, with useful suggestions for future research directions, appears as Chapter 8 in [22]. In [9] it was shown 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. Over the last decade, a sequence of papers [10, 30, 21, 13] greatly generalized

this result and also made it more precise by considering two squares u^2 and v^2 at the same position, with however the third square w^2 offset a distance $k \geq 0$ to the right. First stated and proved as the “New Periodicity Lemma” (NPL) in [10], the main theorem has since been made more specific, with 12 of 14 subcases proved [30, 21, 13] — a main achievement of this paper is to establish the two that remain. Thus the assumption that three neighbouring squares of well-defined size exist within these well-defined bounds has been shown to lead to the conclusion that locally the string breaks down into repetitions of small period. In this paper we begin by proving a lemma that deals in a precise way with just two overlapping squares; we then apply this result to complete the proof of the final two cases of the NPL. We are as a consequence able to characterize the general case of three overlapping squares — no two constrained to begin at the same position — and therefore we can make a start on considering the combinatorial consequences.

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 & Kucherov [19] showed that the maximum number of runs (usually denoted $\rho(n)$) was linear in n , and moreover they described a linear-time algorithm to compute all the runs in any given string; but their proof was non-constructive — the maximum number of runs was shown to be $\Theta(n)$ but no constant of proportionality was specified. As briefly described in Section 2, the resulting research has led to the conclusion that $\rho(n)$ is at least $0.9445757n$ [31, 20] and no more than $n-1$ [2] — 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 is a form of brute force: global data structures (suffix array, longest common prefix array, Lempel-Ziv decomposition) need to be computed in an extended preprocessing phase, when of course runs are generally a *local* phenomenon. Moreover, it has been shown [26] that the *expected* number of runs in a string is much less than string length: runs generally occur sparsely. A global approach is necessitated by the absence of a detailed understanding of the combinatorics of overlapping occurrences of runs in strings.

In Section 2 terminology, notation and the relevant background are reviewed; Section 3 shows how to express the general case of three overlapping squares,

making use of a careful analysis of two overlapping squares; Section 4 makes use of the new result to prove the two remaining subcases (3 & 7) of the NPL; then in Section 5 a further application to the general case of three overlapping squares is proved; finally, in Section 6 we briefly discuss future research directions.

2 Preliminaries

(Usage generally follows [32].) A *string* is a finite sequence of symbols (*letters*) drawn from some finite or infinite set Σ called the *alphabet*. The alphabet *size* is $\sigma = |\Sigma|$. We write a string \mathbf{x} in mathbold, and we represent it as an array $\mathbf{x}[1..n]$ for some $n \geq 0$. We call $n = x$ the *length* of \mathbf{x} . For $x = 0$, $\mathbf{x} = \epsilon$, the *empty string*.

If $\mathbf{x} = \mathbf{uvw}$, then \mathbf{u} is said to be a *prefix*, \mathbf{v} a *substring* (or *factor*) and \mathbf{w} a *suffix* of \mathbf{x} . If $\mathbf{x} = \mathbf{uv}$, $0 \leq u < x$, then \mathbf{vu} is said to be the u^{th} *rotation* of \mathbf{x} , written $R_u(\mathbf{x})$. If $\mathbf{x} = \mathbf{uv} = \mathbf{wu}$ for $u < x$, then \mathbf{u} is a *border* of \mathbf{x} , and \mathbf{x} has *period* $p = x - u$; that is, for every $i \in 1..u$, $\mathbf{x}[i] = \mathbf{x}[i+p]$. The string

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

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

If $\mathbf{v} = \mathbf{x}[i..j]$ has period p , where $v/p \geq 2$, and if neither $\mathbf{x}[i-1..j]$ nor $\mathbf{x}[i..j+1]$ (whenever these are defined) has period p , then the range $i..j$ in \mathbf{x} is said to be a *maximal periodicity* or *run* in \mathbf{x} [23]. A run is identified by a 4-tuple (i, p, e, t) , where we choose p to be the *minimum period* of \mathbf{v} , $e = \lfloor v/p \rfloor \geq 2$ is its *exponent*, and $t = v \bmod p \in 0..p-1$ is its *tail*. Then $j = i + pe + (t-1)$. The string (1) has five runs

$$(1, 3, 2, 0), (1, 5, 2, 0), (3, 1, 2, 0), (4, 2, 2, 1), (8, 1, 2, 0)$$

corresponding to $(aba)^2$, $(abaab)^2$, a^2 , $(ab)^2a$, a^2 , respectively.

Every run in \mathbf{x} determines $t+1$ *repetitions*,

$$(i, p, e), (i+1, p, e), \dots, (i+t, p, e),$$

where (i', p, e) , $i \leq i' \leq i+t$, identifies the substring

$$\mathbf{x}[i'..i'+pe-1] = \mathbf{x}[i'..i'+p-1]^e.$$

Thus every repetition in \mathbf{x} is a subrange of exactly one run in \mathbf{x} . For example, (1) has six repetitions

$$(1, 3, 2), (1, 5, 2), (3, 1, 2), (4, 2, 2), (5, 2, 2), (8, 1, 2)$$

corresponding to $(aba)^2$, $(abaab)^2$, a^2 , $(ab)^2$, $(ba)^2$, a^2 , respectively. Where no ambiguity arises, we will generally refer to runs and repetitions as substrings (for example, $(aba)^2$, $(ab)^2a$) rather than as ranges in \mathbf{x} (1..6, 4..8). If $e = 2$, we say

that the repetition is a *square*. We say that a square u^2 is *irreducible* if u is not itself a repetition, *regular* if u has no square prefix.

There were three classical algorithms proposed [5, 1, 24] for computing all the repetitions in a string of length n , each executing in $O(n \log n)$ time, asymptotically optimal since the *Fibonacci string* f_k , defined by

$$f_0 = b, f_1 = a; k \geq 2 \implies f_k = f_{k-1}f_{k-2},$$

contains $O(f_k \log f_k)$ repetitions [5, 18, 12]. In [23] Main proposed an algorithm to compute all the “leftmost” runs, extended by Kolpakov & Kucherov in [19] to compute all runs. As mentioned in Section 1, this approach makes extensive use of preprocessing, but still executes in linear time, based on a complex proof that the maximum number $\rho(n)$ of runs in any string of length n satisfies

$$\rho(n) \leq K_1 n - K_2 \sqrt{n} \log_2 n \quad (2)$$

for some universal positive constants K_1 and K_2 . Even though [19] provided computational evidence (up to $n = 60$) that $\rho(n) \leq n$, the method of proof allowed no bounds to be placed on K_1 and K_2 . Over the last decade, the bounding of $\rho(n)/n$ has become a growth industry, leading to a lower bound 0.9445757 [15, 14, 25, 31, 20] and an asymptotic upper bound 1.029 [28, 27, 6, 16, 17, 7, 8], the latter result achieved using three years of CPU time on a supercomputer [29]. Very recently, Bannai *et al.* [2] have published a remarkably simple proof, using Lyndon words, that in fact $\rho(n)/n < 1$ for all n . Meanwhile, more efficient and truly linear (independent of alphabet size σ) algorithms for computing runs have been proposed — for example, [3, 4] — but still with heavy preprocessing and the same general approach. Since, as noted in Section 1, runs are expected to be sparse in strings, even for small σ [26], a heavy-handed global approach seems inappropriate.

A parallel approach has sought to find a combinatorial basis for estimating the maximum number of runs in a string, specifically by considering the consequences of assuming that two squares occur at the same position in a string, with a third nearby, somewhat to the right. This generalizes the “Three Squares Lemma” [9] that considers three squares at the same position in the string:

Lemma 1. *Suppose u^2 is irreducible, and suppose $v \neq u^j$ for any $j \geq 1$. If u^2 is a prefix of v^2 , in turn a proper prefix of w^2 , then $w \geq u+v$.*

Here is the original NPL, as stated and proved in [10]:

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

The proof required consideration of 14 subcases based on the magnitudes of k and w (see Table 1), each of which led to a proof by contradiction of the regularity of u . Figure 1 shows two of these subcases.

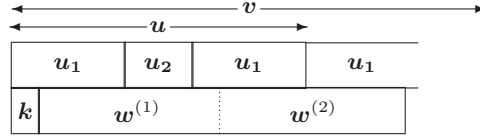


Fig. 1. Subcases 1 & 2 of Lemma 2

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 [21] it was shown that for $u < v \leq 3u/2$, the requirement that $\mathbf{x} = \mathbf{v}^2$ with prefix \mathbf{u}^2 necessitates

$$\mathbf{x} = \mathbf{t}_1^m \mathbf{t}_2 \mathbf{t}_1^{m+1} \mathbf{t}_2 \mathbf{t}_1, \quad (3)$$

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} .

Table 1. The 14 subcases identified in [10], slightly modified, for three neighbouring squares \mathbf{u} , \mathbf{v} , \mathbf{w} (with $v - u < w < v$, $w \neq u$, $0 \leq k < v - u$).

Subcase S	k	$k+w$	$k+2w$	Special 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	$u_1 < k < u_1+u_2$	$2u < k+w < 2u+u_2-1$	—	—

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

$$\mathbf{x} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 \mathbf{u}_1 \mathbf{u}_2, \quad (4)$$

where $\mathbf{u} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1$, $\mathbf{v} = \mathbf{u}_1 \mathbf{u}_2 \mathbf{u}_1 \mathbf{u}_1 \mathbf{u}_2$ and $u_1 = 2u - v$, $u_2 = 2v - 3u$. Note that setting $\mathbf{t}_1 = \mathbf{u}_1 \mathbf{u}_2$, $\mathbf{t}_2 = \mathbf{u}_1$ converts the form (4) into (3), but with $m < 2$. For

Table 2. Structure of \mathbf{x} for subcases $S \in 1..14$: σ is the largest alphabet size consistent with u, v, k, w [13]; \mathbf{d} , \mathbf{d}_1 and \mathbf{d}_3 are prefixes of \mathbf{x} with $d = \gcd(u, v, w)$, $d_1 = \gcd(u - w, v - u)$, $d_2 = \gcd(u, v - w)$, $d_3 = v \bmod d_2$.

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	$(\forall \mathbf{x})$ specified cases	$\mathbf{x} = \mathbf{d}_1^{u/d_1} \mathbf{d}_1^{v/d_1} \mathbf{d}_1^{(v-u)/d_1}$ $\mathbf{x} = \mathbf{d}^{x/d}$
11–14	$\sigma = d$ or $d_2 \leq 2u - v$ otherwise	$\mathbf{x} = \mathbf{d}^{x/d}$ $\mathbf{x} = ((\mathbf{d}_3^{d_2/d_3})^{v/d_2})^2$

this case, after much experimental and theoretical work [30, 21, 13], the revised NPL can be stated as follows:

Lemma 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 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 Table 1, the corresponding structure of \mathbf{x} is given in Table 2.*

In other words, \mathbf{x} breaks down into repetitions of small period — essentially, the postulate of three such squares cannot be satisfied. Twelve of the 14 subcases have been proved [30, 21, 13]; in this paper we prove the remaining two, subcases 3 and 7.

We believe moreover 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?

In this paper we begin to answer these questions by first considering only two overlapping squares in some detail, then making the observation that three overlapping squares can always be thought of as two sets of two overlapping squares. In Section 3 a general lemma for two squares is stated and proved, a result used to establish subcases 3 and 7 in Section 4. As a further application, a “sample” three squares lemma is proved in Section 5. Section 6 briefly discusses future work.

We conclude this section by stating results that will be useful in what follows:

Lemma 4. (“The Periodicity Lemma” [11]) *Let p and q be two periods of $\mathbf{x} = \mathbf{x}[1..n]$, and let $d = \gcd(p, q)$. If $p+q \leq n+d$, then d is also a period of \mathbf{x} .*

Lemma 5. (Corollary to Lemma 4.) *If $\mathbf{x} = \mathbf{uvw}$, where \mathbf{uv} and \mathbf{vw} have period $p \leq v$, then \mathbf{x} has period p .*

Lemma 6. ([22], Lemma 8.1.3) *If a string \mathbf{x} of period p has a substring \mathbf{u} , $u \geq p$, of period q , where $q \mid p$, then \mathbf{x} has period q .*

Lemma 7. ([21], Lemma 8) Suppose both \mathbf{x} and $R_v(\mathbf{x})$, $0 < v < x$, have period u , and let $\ell = x \bmod u > 0$ and $r = \lfloor x/u \rfloor$. Let \mathbf{x}_v denote $R_v(\mathbf{x})$, and let $d = \gcd(u, \ell)$. Then

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

Proof (Lemma 7(a)). Let $\mathbf{v} = \mathbf{x}[1..v] = \mathbf{x}_v[x-v+1..x]$, and let $\mathbf{u} = \mathbf{x}_v[1..u]$, as shown in Figure 2. Since \mathbf{x}_v has period u and $v \geq \ell$, $\mathbf{u}[1..\ell]$ is a suffix of \mathbf{x}_v and \mathbf{v} . Recalling that \mathbf{v} is also a prefix of \mathbf{x} , we see that $\mathbf{x}_{v-\ell} = \mathbf{v}[v-\ell+1..v]\mathbf{u}$ has prefix $(\mathbf{u}[1..\ell])^2$. \square

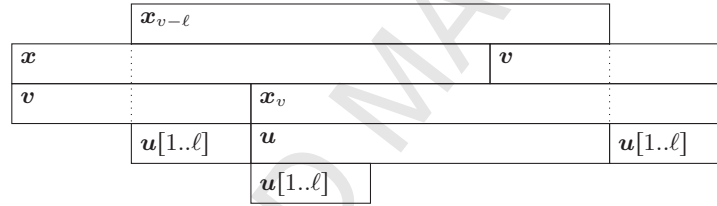


Fig. 2. Lemma 7(a)

3 Characterizing the General Case

We are interested in the cases that arise when a square \mathbf{u}^2 beginning at some position i in a string overlaps with a second square \mathbf{v}^2 at position $i+k$, $k \geq 0$, to its right.

Lemma 8. 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) $k+v < u < 2v$ ($k < \min(v-1, u-v)$):

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

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

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

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

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

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

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

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

(d) $v = u$ ($k \geq 0$) :

$$\mathbf{v} = R_k(\mathbf{u}).$$

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

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

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

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

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

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

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

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

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

Proof.

(a) Let $\mathbf{z} = \mathbf{u}[k+v+1..u] = \mathbf{v}[1..u-(k+v)]$, suffix of \mathbf{u} and prefix of \mathbf{v} .

\mathbf{u}		\mathbf{u}	
\mathbf{k}	\mathbf{v}	\mathbf{v}	

Observe that

$$\mathbf{u}[k+j] = \mathbf{v}[j] = \mathbf{u}[j-z], \quad z+1 \leq j \leq v,$$

so that $\mathbf{u}[1..k+v] = \mathbf{k}\mathbf{v}$ has period $k+z = u-v = p$ (where $\mathbf{p} = \mathbf{u}[1..u-v]$). Consequently, we may write $\mathbf{x} = (\mathbf{p}^e \mathbf{z})^2$, where $e = \frac{k+v}{u-v} > 1$ (since $k+v < u$). Noting that $\mathbf{v} = \mathbf{u}[k+1..u-z]$, with $k < u-v$, we see also that $\mathbf{u} = \mathbf{k}\mathbf{q}^{f-1}\mathbf{z}$, where $\mathbf{q} = R_k(\mathbf{p})$,

$$f = \frac{u}{u-v} = \frac{v}{u-v} + 1 > 2.$$

Hence $\mathbf{x} = \mathbf{p}^e \mathbf{z}\mathbf{k}\mathbf{q}^{f-1}\mathbf{z}$. But $\mathbf{z}\mathbf{k}$ is a prefix of the second copy of \mathbf{v} of length p , and comparing with the first copy of \mathbf{v} , we see that therefore $\mathbf{z}\mathbf{k} = R_k(\mathbf{p}) = \mathbf{q}$. Since moreover $\mathbf{z} = \mathbf{q}^g$, where

$$g = \frac{z}{u-v} = \frac{u}{u-v} - \frac{k+v}{u-v} = f - e \leq 1,$$

we find $\mathbf{x} = \mathbf{p}^e \mathbf{q}^f \mathbf{q}^{f-e}$, as claimed. (Note that $g = 1$ iff $k = 0$.) Finally, writing $\mathbf{q} = \mathbf{p}[k+1..u-v]\mathbf{p}[1..k]$ and $f-e = 1 - \frac{k}{u-v}$, we find that $\mathbf{q}^{f-e} = \mathbf{p}[k+1..u-v]$.

- (b) Let $\mathbf{z} = \mathbf{u}[1..k+v-u]$, a possibly empty prefix of \mathbf{u} and suffix of \mathbf{v} .

\mathbf{u}		\mathbf{u}	
\mathbf{k}	\mathbf{v}	\mathbf{v}	

Observe that

$$\mathbf{u}[k+j] = \mathbf{v}[j] = \mathbf{u}[z+j], \quad 1 \leq j \leq v-z = u-k,$$

where $z = k+v-u < k$. Thus $\mathbf{u}[z+1..u]$ and \mathbf{v} have period $k-z = u-v = p$. Consequently, setting $\mathbf{p} = \mathbf{u}[z+1..k] = \mathbf{v}[1..k-z]$, we may write $\mathbf{x} = (\mathbf{z}\mathbf{p}^e)^2$, where $e = \frac{u-z}{u-v} \geq 1$, since $z \leq v$. Noting that

$$u-z = u-k-v+u = (u-v) + (u-k), \quad (5)$$

we see that $e = 1 + \frac{u-k}{u-v}$.

Since \mathbf{p} is a prefix of \mathbf{v} and $z+p = k$, it follows that $\mathbf{k} = \mathbf{z}\mathbf{p}$. Thus we can also write $\mathbf{x} = (\mathbf{k}\mathbf{p}^{e-1})^2$.

Finally, setting $\mathbf{y} = \mathbf{u}[k+2v-u+1..u]$, since $\mathbf{z}\mathbf{y}$ is a suffix of \mathbf{u} of length

$$k+v-u+2(u-v)-k = u-v = p,$$

it follows that $\mathbf{z}\mathbf{y}$ is a rotation of \mathbf{p} . In fact, $\mathbf{z}\mathbf{y} = \mathbf{q} = R_c(\mathbf{p})$, where by (5)

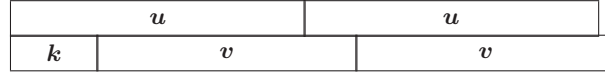
$$c = (u-z) \bmod (u-v) = (u-k) \bmod (u-v).$$

Then $\mathbf{z} = \mathbf{q}^f$, where

$$f = \frac{z}{p} = \frac{(k+v)-u}{u-v} = \frac{k}{u-v} - 1 \leq 1,$$

so that $\mathbf{q}^f = \mathbf{q}[1..z]$ and $\mathbf{x} = (\mathbf{q}[1..k+v-u]\mathbf{p}^e)^2$, as required.

- (c) Let $z = \mathbf{u}[1..k+v-u]$, nonempty prefix of \mathbf{u} and suffix of \mathbf{v} .



As in (b), observe that

$$\mathbf{u}[k+j] = \mathbf{v}[j] = \mathbf{u}[z+j], \quad 1 \leq j \leq v-z = u-k,$$

with $z = k+v-u < k$. Again $\mathbf{u}[z+1..u]$ has period $k-z = u-v = p$, where $\mathbf{p} = \mathbf{u}[z+1..k] = \mathbf{v}[1..u-v]$. However, unlike (b), not \mathbf{v} , but only $\mathbf{v}[1..v-y]$, has period p , where $\mathbf{y} = \mathbf{v}[2u-(k+v)+1..v] = \mathbf{v}[u-z+1..v]$ and

$$y = k+2v-2u = z-(u-v) = z-p < z.$$

Thus, noting that $z < v$ and setting

$$e = \frac{u-z}{u-v} = \frac{(u-v)+(u-k)}{u-v} = 1 + \frac{u-k}{u-v} > 1,$$

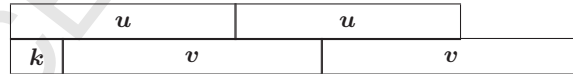
we can write $\mathbf{x} = (z\mathbf{p}^e)^2\mathbf{y}$. Since $\mathbf{u}[z+1..u]$ has prefix \mathbf{p} and $k = z+p$, we see that $\mathbf{k} = z\mathbf{p}$; further, since z has suffix \mathbf{y} and $z = p+y$, it follows that $z = \mathbf{qy}$ for some rotation \mathbf{q} of \mathbf{p} . In fact, $\mathbf{q} = R_c(\mathbf{p})$, where

$$c = (v-y) \bmod (u-v) = ((v-z)+(u-v)) \bmod (u-v) = (u-k) \bmod (u-v).$$

Noting that $\mathbf{k} = \mathbf{qyp}$, we see that both \mathbf{kv} and \mathbf{x} have border \mathbf{qy} , while $\mathbf{x} = (\mathbf{qyp}^e)^2\mathbf{y}$, as required.

- (d) Obvious.

- (e) Again let $z = \mathbf{u}[1..k+v-u]$, nonempty prefix of \mathbf{u} and suffix of \mathbf{v} .



Observe that

$$\mathbf{u}[k+j] = \mathbf{v}[j] = \mathbf{u}[z+j], \quad 1 \leq j \leq u-z,$$

where $z = k+v-u > k$. Thus $\mathbf{u}[k+1..u]$ and $\mathbf{v}[1..u-k]$ have period $z-k = v-u = p$. Therefore, setting $\mathbf{p} = \mathbf{v}[1..z-k]$, $e = \frac{u-k}{v-u}$, we can write $\mathbf{x} = (\mathbf{kp}^e)^2\mathbf{y}$, where $\mathbf{y} = \mathbf{v}[2u-(k+v)+1..v]$ is a suffix of \mathbf{v} , $y = (k+v-u)+(v-u) = z+p > z$.

Since z is a suffix of \mathbf{y} and $y-z = p$, it follows that $\mathbf{y} = \mathbf{qz}$, where $\mathbf{q} = R_c(\mathbf{p})$, $c = (u-k) \bmod (v-u)$. Similarly, since \mathbf{k} is a prefix of z and $z-k = p$, we see

that $z = kp$, hence that $y = qkp$. Thus $x = (kp^e)^2qkp$, as claimed, and x and kv both have border kp .

To see that $e > 1$, note that $k < \frac{v}{2}$, so that $e > \frac{v-u/2}{v-u} > 1$.

- (f) Let $z = v[1..2u - (k+v)]$, nonempty prefix of v and suffix of u .

u		u			
k	v			v	

Observe that

$$u[k+j] = v[j] = u[(u-z)+j], \quad 1 \leq j \leq z.$$

Thus $u[k+1..u]$ has period $p = (u-z) - k = v - u$. Consider

$$y = v[u-k+1..v] = u[1..(k+v)-u] = u[1..k+p],$$

nonempty suffix of v and prefix of u . Since y has prefix k , it follows that $y = kp$, where $p = u[k+1..k+p]$. Thus for $e = \frac{u-k}{v-u}$, $v = p^e kp$, $x = k(p^e kp)^2$, as stated. Since $k+v < 2u$, $u-k > v-u$, and so $e > 1$.

- (g) Let $z = v[2u - k + 1..k+v]$, suffix of v .

u		u			
k	v			v	

Observe that

$$v[j] = u[k+j] = v[u+k+j], \quad 1 \leq j \leq u-k,$$

so that $v[1..2u-k]$ has period $p = u$. Then for $p = u[k+1..u]u[1..k]$ and $e = \frac{2u-k}{u} \geq 1$, $v = p^e z$ and $x = k(p^e z)^2$. \square

Case (g) of Lemma 8 is not a true overlap, but is included for completeness. Note also that cases (b), (c) and (f) require $k > 0$, and so do not exist if it is assumed that u^2 and v^2 (or v^2 and w^2) occur at the same position. Overall, it turns out that analysis of overlapping squares requires breaking down the interval $[k, 2v]$ for u into seven highly nonuniform subintervals:

$$\left[k, \frac{k+v}{2} \right], \left(\frac{k+v}{2}, \frac{2(k+v)}{3} \right), \left[\frac{2(k+v)}{3}, v \right), [v], \left(v, \frac{k}{2} + v \right), \left[\frac{k}{2} + v, k+v \right], (k+v, 2v).$$

We make the observation that if a third square w^2 begins to the right of the starting position of v^2 , sufficiently near to satisfy the postulates of Lemma 8,

then the analysis of the three squares $\mathbf{u}^2, \mathbf{v}^2, \mathbf{w}^2$ reduces to a simultaneous consideration of two of the lemma's cases. Thus, for example, the analysis of the situation shown in Figure 20 of Section 5 would take place in terms of the simultaneous occurrence of cases (e) (for \mathbf{u}^2 and \mathbf{v}^2) and (b) (for \mathbf{v}^2 and \mathbf{w}^2). Indeed, all cases of three overlapping squares can be represented by pairs $[ij]$, $a \leq i, j \leq f$, referring to the cases (a)-(g) arising in Lemma 8.

4 Subcases 3 & 7 of the NPL

In this section we prove the two remaining subcases of Lemma 3.

4.1 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 9 (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*

$$0 \leq k \leq u_1 < u_1 + u_2 < w < v \quad (6)$$

$$k + w \leq u \quad (7)$$

$$k + 2w > u + u_1 \quad (8)$$

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. As we have seen (4), 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 8(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 $e = \frac{k+w}{u-w} > 1$.

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

Fig. 3. String \mathbf{u} in Subcase 3

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 (6). 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 6. Thus the lemma holds assuming \mathbf{u} has period p .

Note first from (6) and (7) that $u_1 + u_2 < w < u$, hence that $u_1 - p = u_1 - (u - w) > 0$. Then to see that \mathbf{u} in fact has period p , consider two cases:

$$u_1 \geq k + w - (u_1 + u_2) \geq p \tag{9}$$

and

$$k + w - (u_1 + u_2) < p < u_1. \tag{10}$$

In the first case, the prefix $\mathbf{k}w = \mathbf{p}^e$ of \mathbf{u} extends at least p positions into the suffix $\mathbf{u}_1^{(2)}$. Since \mathbf{u}_1 is a prefix of $\mathbf{k}w$, \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 5, \mathbf{u} has period p .

The second case (10) is more complicated (see Figure 4). 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

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

so \mathbf{z} is a proper suffix of \mathbf{u}_1 . The prefix \mathbf{p} and the suffix \mathbf{z} of \mathbf{u}_1 must overlap because (10) is equivalent to $u_1 < p + z$. Noting that $p = u - w = k + z$ so that $\mathbf{p} = \mathbf{kz}$, we have (Figure 4)

$$z = \ell't' = t'\ell \tag{11}$$

and

$$\mathbf{u}_1 = \mathbf{p}\ell = \mathbf{kz}\ell = \mathbf{k}\ell'z = \mathbf{k}\ell't'\ell, \tag{12}$$

where ℓ and ℓ' are respectively the proper suffix and proper prefix of \mathbf{z} of length $\ell = \ell' = u_1 - p$, and t' is the border of \mathbf{z} of length $t' = z - \ell$. Since by (11) \mathbf{z} has border t' , it therefore has period ℓ , as does $\ell'z\ell = \ell'\ell't'\ell = \ell't'\ell\ell$. Since \mathbf{u}_1 has period p , \mathbf{p} has prefix ℓ ; $\mathbf{p} = \mathbf{kz}$ also has suffix ℓ , so that \mathbf{p} has border ℓ . Observe also that because $k + w - (u_1 + u_2) = k + \ell$, $\mathbf{k}\ell'$ is the prefix of \mathbf{u}_1 that overlaps \mathbf{w} .

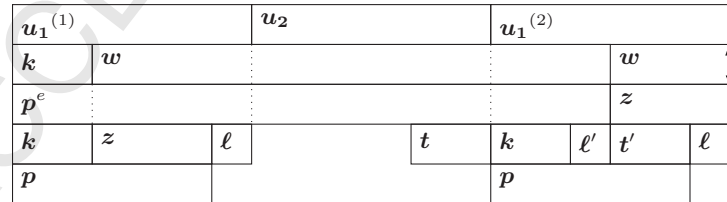


Fig. 4. String \mathbf{u} in Subcase 3 when (10) holds

Let \mathbf{t} be the suffix of $\mathbf{u}_1\mathbf{u}_2$ of length $t = t'$. Then \mathbf{w} has suffix $\mathbf{tk}\ell'$ in which $\mathbf{k}\ell'$ is a prefix of \mathbf{u}_1 . Recall $\mathbf{u}_1 = \mathbf{k}\ell't'\ell$ has period $p = k + t + \ell$. If $\mathbf{t} = \mathbf{t}'$, 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^{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 ℓ' of w .

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

(C1)

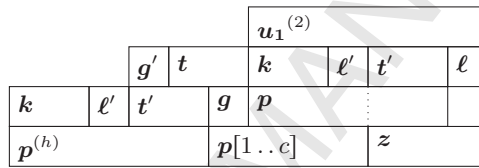


Fig. 5. Subcase 3 when (C1) holds

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

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

and observe that kw has suffix $p[1..c] = gk\ell' = k\ell'g'$.

Thus $p[1..c]$ has border $k\ell'$ and therefore period g . String $\ell'g'$ has period ℓ as a prefix of $z = \ell't'$, and period g as a suffix of $p[1..c]$, so by Lemma 4 it has period $\gcd(g, \ell)$. Then $p[1..c]$ has period g and suffix $\ell'g'$ of period $\gcd(g, \ell) \mid g$ and length $\ell + g \geq g$, so that by Lemma 6 $p[1..c]$ itself has period $\gcd(g, \ell)$. Both $p[1..c]$ and $\ell'z$ have period ℓ and share substring ℓ' , so $p[1..c]z$ has period ℓ by Lemma 5. It also has substring p , so p has period ℓ . Because p has border ℓ as well as period ℓ , any power of p has period ℓ . It follows that $kw = p^e$ has period ℓ and, since $p[1..c]z$ has period ℓ and shares with kw a substring of length $c > \ell$, u has period ℓ by Lemma 5. Recall that u has substring $p[1..c]$ of period $\gcd(g, \ell) \mid \ell$, so u itself has period $\gcd(g, \ell)$ by Lemma 6. Recalling that t is a suffix of $t'g$ and that both are substrings of u , we find that t and t' have period $\gcd(g, \ell)$ and suffix g , so $t = t'$.

(C2)

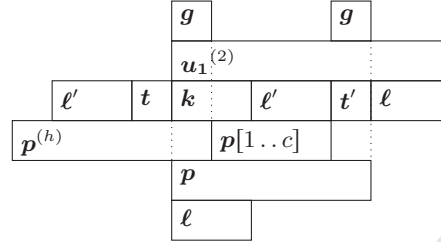


Fig. 6. Subcase 3 when (C2) holds

Suppose (C2) holds; that is, $p^{(h)}$ ends inside the prefix k of $u_1^{(2)}$. Let

$$g = ph - u_1 - u_2$$

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

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

and observe that, from the overlap of the suffix $p[1..c]$ of kw and the prefix p of u_1 , $p[1..k+\ell] = k\ell'$ has period g . Also note that since by (12) $p = k\ell't'$ and since $kw = p^e$ has period p , therefore $p^{(h)}$ has suffix $\ell'tg$. Consider four cases: $0 < g < \ell$, $g = \ell$, $\ell < g \leq z$, and $z < g$.

If $0 < g < \ell$, then since ℓ is a prefix of u_1 , ℓ has period g as a prefix of $k\ell'$. Recall that ℓ is also a suffix of p , so ℓ has border g and period $\ell - g$. Hence, by Lemma 4, ℓ has period $\gcd(g, \ell - g) = \gcd(g, \ell)$. Recall also that $\ell'z$ has period ℓ and substring ℓ of period $\gcd(g, \ell) \mid \ell$, so by Lemma 6, $\ell'z$ has period $\gcd(g, \ell)$. Prefix ℓ' of $\ell'z$ then has period $\gcd(g, \ell) \mid g$, and since ℓ' is also a suffix of the string $k\ell'$ of period g , $k\ell'$ has period $\gcd(g, \ell)$ by Lemma 6. Since $\ell'z$ and $k\ell'$ have period $\gcd(g, \ell)$, $u_1 = k\ell'z$ has period $\gcd(g, \ell)$ by Lemma 5. Both tg and $t'g$ are substrings of u_1 , so $t = t'$.

If $g = \ell$, then p has suffixes $t\ell$ and $z = \ell't = t'\ell$, so immediately $t = t'$. Note that $k\ell'$ and $\ell'z$ have period $g = \ell$, so by Lemma 5, $u_1 = k\ell't'$ has period ℓ .

If $\ell < g \leq z$, then since ℓ is a border of g , g has period $g - \ell$; it also has period ℓ as a substring of the suffix z of p , and thus by Lemma 4 period $\gcd(g, g - \ell) = \gcd(g, \ell)$. String g is a substring of $k\ell'$, which as we have seen has period g , so that by Lemma 6, $k\ell'$ has period $\gcd(g, \ell)$. Since $\ell'z$ and $k\ell'$ have period $\gcd(g, \ell)$, $u_1 = k\ell'z$ has period $\gcd(g, \ell)$ by Lemma 5. Both tg and $t'g$ are substrings of u_1 , so $t = t'$.

If $z < g$, then, as shown in Figure 7, 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 ℓ is a prefix of $p = kz$, so that either ℓ is a prefix of g_1 or g_1 is a prefix of ℓ .

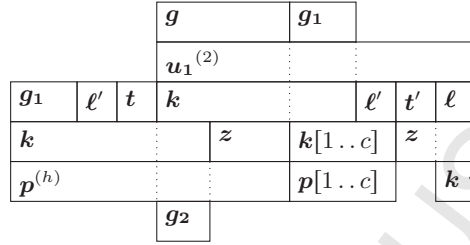


Fig. 7. Subcase 3 when (C2) holds and $z < g$

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) $g' \leq z$

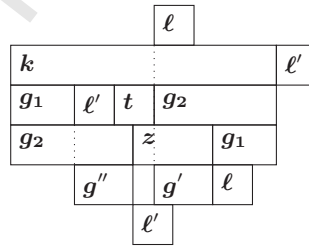


Fig. 8. Subcase 3 when (C2) holds, $z < g$, and $g_1 < g_2$

The demonstration requires several steps:

- $g'\ell$ has period ℓ as a suffix of $z\ell$; it also has border ℓ and period g' , so by Lemma 4, $g'\ell$ has period $\gcd(g', \ell)$.
- $z\ell$ has period ℓ and suffix $g'\ell$ of period $\gcd(g', \ell) \mid \ell$, so by Lemma 6, $z\ell$ has period $\gcd(g', \ell)$.

- g_2 has period g' and prefix g' of period $\gcd(g', \ell) \mid g'$, so by Lemma 6, g_2 has period $\gcd(g', \ell)$.
- $z\ell$ and g_2 have period $\gcd(g', \ell)$ and share substring g' , so by Lemma 5, zg_1 has period $\gcd(g', \ell)$.
- $g''\ell'$ has border ℓ' 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 6, $g''\ell'$ has period $\gcd(g', \ell)$.
- g_2 and $g''\ell'$ have period $\gcd(g', \ell)$ and share substring g'' , so by Lemma 5, $g_2\ell'$ has period $\gcd(g', \ell)$.
- $g_2\ell'$ and zg_1 have period $\gcd(g', \ell)$ and share substring ℓ' , so by Lemma 5, $k = g_2zg_1$ 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 $k = g_1\ell'tg_2$ and $g''z\ell$ have period $\gcd(g', \ell)$ and share substring g'' , so by Lemma 5, $u_1 = kz\ell$ has period $\gcd(g', \ell)$.

(b) $g' > z$

k has period $k - g_2 = g_1 + z$, so k has a prefix $g_1\ell'tg_2' = g_2'zg_1$, 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. $g_1 > g_2$

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 9 & 10)

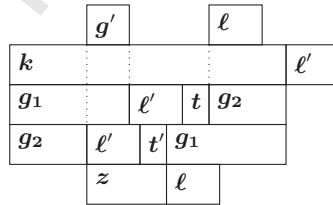


Fig. 9. Subcase 3 when (C2) holds, $z < g$, $g_1 > g_2$, and $g' \leq \ell$

Again several steps are required:

- $g'\ell'$ has period ℓ as a prefix of $z\ell$ and also shares prefix ℓ' with z , so it has border ℓ' and period g' .
- $g'\ell'$ has periods ℓ 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 6, g_1 has period $\gcd(g', \ell)$.

- $z\ell$ has period ℓ and prefix g' of period $\gcd(g', \ell) \mid \ell$, so again by Lemma 6, $z\ell$ 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 5, 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 5, $k = g_2zg_1$ has period $\gcd(g', \ell)$. $g'z\ell$ then has period $\gcd(g', \ell)$, so that $k = g_2zg_1$ and $g'z\ell$ have period $\gcd(g', \ell)$ and share substring g' , so by Lemma 5, $u_1 = kz\ell$ has period $\gcd(g', \ell)$.

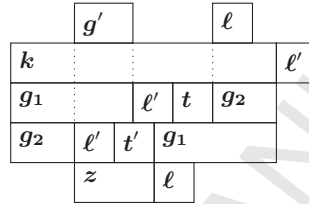


Fig. 10. Subcase 3 when (C2) holds, $z < g$, $g_1 > g_2$, and $\ell < g' \leq z$

(b) $g' > z$

k has period $k - g_1 = g_2 + z$, so k has a prefix $g_2zg'_1 = g'_1\ell'tg_2$, 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)

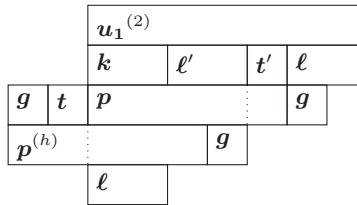


Fig. 11. Subcase 3 when (C3) holds

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

p , so g is also a prefix of ℓ . From $p^{(h)}$ and the occurrence of p that prefixes u_1 , we have

$$p = gk(\ell'[1.. \ell - g]) = k(\ell'[1.. \ell - g])gt'$$

and p has period $t+g$.

Consider the string $\ell't'g$, which occurs near the end of u_1 as a prefix of $\ell't'\ell = \ell'z$. As a substring of $\ell'z = \ell'\ell't'$, it has period ℓ . Since 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 ℓ , so by Lemma 4, it has period $\gcd(t+g, \ell)$. Now p has period $t+g$ and suffix gt' of period $\gcd(t+g, \ell) \mid t+g$, so p itself has period $\gcd(t+g, \ell)$. Because p has border ℓ as well as period $\gcd(t+g, \ell)$, any power of p has period $\gcd(t+g, \ell)$. Thus $kw = p^e$ has period $\gcd(t+g, \ell)$ and, since $\ell'z$ has period $\gcd(t+g, \ell)$ and shares with kw a substring of length ℓ , u has period $\gcd(t+g, \ell)$ by Lemma 5. Since $t\ell$ and $t'\ell$ are substrings of u , therefore $t = t'$.

This completes the proof of Subcase 3. \square

Lemma 10. *Suppose the conditions of Subcase 3 hold (Lemma 9). Then $x = d^{x/d}$, 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 (9) and (10), respectively, found at the beginning of the proof of Lemma 9, 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 u_1 has period ℓ (all symbols used as defined in the proof of Lemma 9). 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 $u_1 = kz\ell = p\ell$ has period ℓ . Then, by Lemma 9, $u = u_1u_2u_1$ has period $d_1 = \gcd(p, u_1+u_2)$. Hence p is a suffix of u_1 , thus a border of u . Therefore u^2 has period d_1 , as well as period u , so that by Lemma 4, 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 x .

Next observe that in the proof of Lemma 9, $g = u_1+u_2-ph$, 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 9, 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 p ($2h$ copies of g) cover u_1u_2g of length $v-u+g$, from which (b) follows. \square

4.2 Subcase 7

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

Lemma 11 (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*

$$u_1 < k < u_1 + u_2 < w < v \tag{13}$$

$$k + w \leq u + u_1 \tag{14}$$

$$k + 2w \leq 2u \tag{15}$$

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. As we have seen (4), 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 8(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 12.

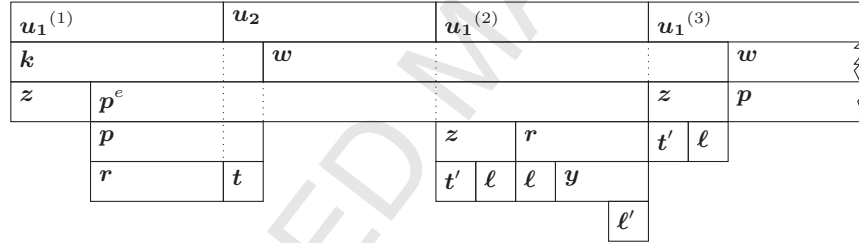


Fig. 12. String \mathbf{uu}_1 in Subcase 7

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 (13). 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 6. Thus the lemma holds assuming \mathbf{u} has period p .

We now embark on a demonstration that \mathbf{u} has period p . Notice (Figure 12) that $\mathbf{u}_1 \mathbf{u}_2$, k , and $\mathbf{z} \mathbf{p}$ are prefixes of \mathbf{u} . Given that $z = k - p$, that $k \in (u_1, u_1 + u_2)$ by (13), and that $z \leq u_1$ by (14), we have

$$\mathbf{k} = \mathbf{z} \mathbf{p} = \mathbf{z} \mathbf{r} \mathbf{t} = \mathbf{u}_1 \mathbf{t}, \tag{16}$$

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 (15), so that $p \geq z$. Also, by (13)

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

so that in fact $z \geq 2$. Note further from (13) that

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

while from (15) and (13),

$$p = u - w \geq k/2 > u_1/2. \quad (17)$$

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

$$2 \leq z \leq p < u_1 = z + r, \quad (18)$$

from which we conclude that $r > 0$.

Also, since

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

and recalling from (13) 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 $\mathbf{k} = \mathbf{zrt} = \mathbf{zp}$, 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 7(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 ℓ . Thus we may write $\mathbf{u}_1 = \mathbf{t}'\ell^2 \dots$, where $\mathbf{t}'\ell = z$. 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 $z = \mathbf{t}'\ell$, $r = \ell \mathbf{y}$ and $\mathbf{p} = \ell \mathbf{y} \mathbf{t}$.

Since $\mathbf{u}' = \mathbf{p} \mathbf{w}$ has period p , \mathbf{w} is a prefix of \mathbf{u}' . As we see from Figure 12, 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 $z = \mathbf{t}'\ell$. Thus \mathbf{t}' is a border of z . Now let ℓ' be the prefix of z of length ℓ , so that $z = \mathbf{t}'\ell = \ell'\mathbf{t}'$ has period ℓ . Note further that since \mathbf{w} has suffix $\mathbf{y}z$, which in turn has suffix $\mathbf{t}'\ell\ell = \ell'\mathbf{t}'\ell$, therefore ℓ' is a suffix of \mathbf{y} .

Assume $t = t'$. Then $\mathbf{k} = \mathbf{z} \mathbf{r} \mathbf{t} = \mathbf{z} \mathbf{r} \mathbf{t}'$ occurs in \mathbf{w} , and as \mathbf{w} has period p , so does \mathbf{k} . \mathbf{u} then has prefix $\mathbf{k} = \mathbf{z} \mathbf{p}$ and suffix \mathbf{p}^e both of period p and both including \mathbf{p} , so by Lemma 5, \mathbf{u} has period p , as desired. Hence it will suffice to show $t = t'$.

From (18) it follows that a complete copy of \mathbf{p} occurs $h \geq 2$ times in \mathbf{u}' . Several cases arise, based on the position of the suffix \mathbf{t} of the h^{th} occurrence of \mathbf{p} :

- (C1) \mathbf{t} ends inside the prefix \mathbf{z} of $\mathbf{u}_1^{(3)}$
- (C2) \mathbf{t} is a substring of the suffix \mathbf{y} of $\mathbf{u}_1^{(2)}$, but $\mathbf{t}\ell$ is not.
- (C3) $\mathbf{t}\ell$ is a substring of the suffix \mathbf{y} of $\mathbf{u}_1^{(2)}$.
- (C4) \mathbf{t} begins to the left of the suffix \mathbf{y} of $\mathbf{u}_1^{(2)}$ and ends inside \mathbf{y} .

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

(C1)

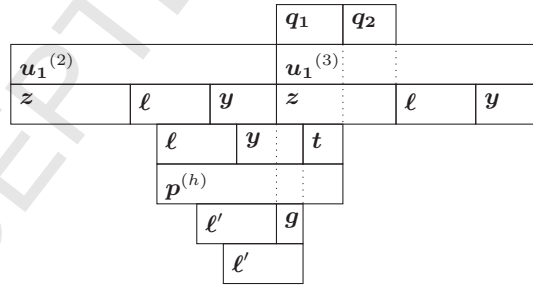


Fig. 13. Subcase 7 when (C1) holds and $g > 0$

Suppose first that (C1) holds, and write $z = \mathbf{q}_1 \mathbf{q}_2$, where \mathbf{q}_1 is a nonempty suffix of \mathbf{p} and, by the periodicity of \mathbf{u}' , \mathbf{q}_2 is a prefix of \mathbf{p} .

We have shown that $\mathbf{u}' = \mathbf{p}^h \mathbf{q}_2$, where $\mathbf{p} = \ell \mathbf{y} \mathbf{t}$, and so $\mathbf{u}' = \mathbf{p}^{h-1}(\ell \mathbf{y})(\mathbf{t}) \mathbf{q}_2$. As in the proof of Lemma 9, we introduce the *gap* $g = q_1 - t$, a measure of the overlap between the prefix \mathbf{q}_1 of $\mathbf{u}_1^{(3)}$ and the suffix \mathbf{t} of $\mathbf{p}^{(h)}$. If $g \geq 0$, then \mathbf{t}

is a substring of z ; otherwise, t ends inside z but begins before it. Note that if $g = 0$, then $q_1 = t = t'$ and the remainder of the proof follows.

Suppose then that $g > 0$ (Figure 13), so that $q_1 = gt$ for some string g of length g . In this case, note that $l't$ and $l't'$ are substrings of $l'z = l'l't'$, as we have seen of period l , and so both these strings also have period l , implying that $t = t'$, as required.

We now show further that for $g > 0$, u_1 has period $\gcd(g, l)$. Since t is a substring of $z = t'l$, $g \leq l$. Therefore, since ly is a suffix of u_1 and $p = lyt$, ly and l both have period g , as does l' , since it is a suffix of ly . Observe that $l'g$ has period l as a prefix of $l'z$, as well as period g as a suffix of ly , so that by Lemma 4, $l'g$ has period $\gcd(g, l)$. zl then has period l and a substring l' of period $\gcd(g, l) \mid l$, so by Lemma 6, zl has period $\gcd(g, l)$. ly has period g and suffix g of period $\gcd(g, l)$, so by Lemma 6, it has period $\gcd(g, l)$. zl and ly have period $\gcd(g, l)$ and share substring l , so by Lemma 5, $u_1 = zly$ has period $\gcd(g, l)$.

Suppose next that $g < 0$, so that $t = gq_1$ for some string g of length $|g|$, as shown in Figure 14. Again ly and l both have period $|g|$. If $|g| \leq l$, then tl is a substring of $l'z$, so $t = t'$. However, when $g < 0$, it is possible that $|g| > l$. In general, let g' be the suffix of z of length $|g|$. The suffix $q_2 = lg'$ of z has border l and thus period $q_2 - l = |g|$. It also has period l as a suffix of z , so by Lemma 4, it has period $\gcd(g, l)$. $l'z$ then has period l and suffix lg' of period $\gcd(g, l) \mid l$, so that by Lemma 6, $l'z$ has period $\gcd(g, l)$. Also by Lemma 6, ly has period $\gcd(g, l)$ since it has period $|g|$ and, by the periodicity of u' , substring g' of period $\gcd(g, l) \mid |g|$. Both ly and $l'z$ have period $\gcd(g, l)$ and share substring l' , so that by Lemma 5, lyz has period $\gcd(g, l)$. Since tl and $t'l$ are both substrings of lyz , therefore again $t = t'$.

Note finally that since zl and ly both have period $\gcd(g, l)$ and share substring l , therefore by Lemma 5, $u_1 = zly$ again has period $\gcd(g, l)$, as it did also for $g > 0$.

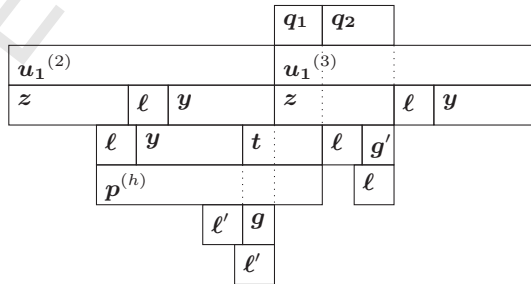


Fig. 14. Subcase 7 when (C1) holds and $g < 0$

(C2)

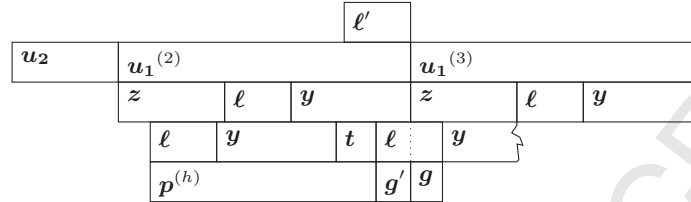


Fig. 15. Subcase 7 when (C2) holds

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

$g'z$ has period $t + g'$ as a prefix of ℓy and period ℓ as a suffix of $\ell'z$, so by lemma 4 it has period $\gcd(t + g', \ell)$. Thus ℓy has period $t + g'$, and a prefix $g'z$ of period $\gcd(t + g', \ell) \mid t + g'$, so ℓy has period $\gcd(t + g', \ell)$ by Lemma 6. Since $z\ell$ has period ℓ and prefix z of period $\gcd(t + g', \ell) \mid \ell$, therefore by Lemma 6, $z\ell$ has period $\gcd(t + g', \ell)$. $z\ell$ and ℓy have period $\gcd(t + g', \ell)$ and share substring ℓ , so by Lemma 5, $u_1 = z\ell y$ has period $\gcd(t + g', \ell)$. Remarking that $\ell't$ and $\ell't'$ are both substrings of u_1 , we again conclude that $t = t'$.

(C3)

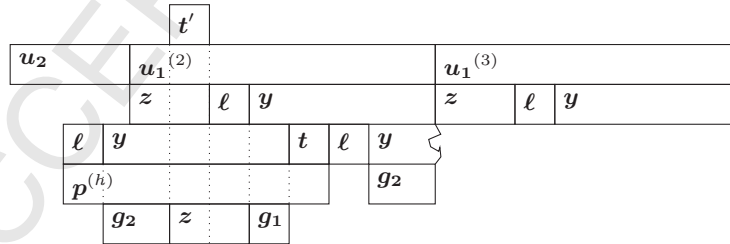


Fig. 16. Subcase 7 when (C3) holds

Suppose (C3) holds; that is, $t\ell$ is a substring of y . In this case, $z = t'\ell$ is also a substring of y because $y \geq t + \ell$ and the prefix ℓy of $p^{(h)}$ ends at least z and at most y positions from the end of $u_1^{(2)} = \ell't'\ell y$.

Let g_1 and g_2 be the (possibly empty) substrings of y immediately before and after $t\ell$ such that $y = g_1 t \ell g_2$. Since $u_1 z$ has period p , g_1 and g_2 are borders of y such that $y = g_1 t \ell g_2 = g_2 t' \ell g_1$.

Recall that ℓ' is a suffix of y , so ℓ' is a suffix of ℓg_2 . Since the prefix ℓg_2 of p occurs before the substring z of y , ℓ' also occurs before z .

If $g_1 = g_2$, then t and t' occur at the same positions in two copies of y , so that $t = t'$. 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 g be the nonempty prefix of g_2 such that $g_2 = g g_1$. g_1 is a border of g_2 , so that g_2 has period g .

- (a) $g \leq z$

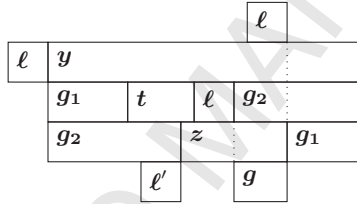


Fig. 17. Subcase 7 when (C3) holds, $g_1 < g_2$, and $g \leq z$

The proof requires several steps:

- As a suffix of $\ell'z$, ℓg has period ℓ and suffix ℓ , hence border ℓ and period g , therefore by Lemma 4 period $\gcd(g, \ell)$.
- Then $\ell'z$ has period ℓ and a suffix ℓg of period $\gcd(g, \ell) \mid \ell$, so by Lemma 6, $\ell'z$ has period $\gcd(g, \ell)$.
- Since g_2 has period g and prefix g of period $\gcd(g, \ell) \mid g$, therefore by Lemma 6, g_2 has period $\gcd(g, \ell)$.
- Since prefix g_2 of y and $\ell'z$ have period $\gcd(g, \ell)$ and share substring ℓ' , therefore $g_2 z$ has period $\gcd(g, \ell)$ by Lemma 5.
- Thus $t\ell$ and $t'\ell$ both have period $\gcd(g, \ell)$ as substrings of $g_2 z$, implying that $t = t'$.

Note that $g_2 z$ and g_2 have period $\gcd(g, \ell)$ and share substring g , so that by Lemma 5, $y = g_2 z g_1$ has period $\gcd(g, \ell)$. Since ℓg_2 has period $\gcd(g, \ell)$ as a substring of y , and since g_2 is a prefix of y , therefore by Lemma 5, ℓy has period $\gcd(g, \ell)$. $z\ell = \ell'z$ and ℓy have period $\gcd(g, \ell)$ and share substring ℓ , so by Lemma 5, $u_1 = z\ell y$ has period $\gcd(g, \ell)$.

(b) $g > z$

y has period $y - g_2 = z + g_1$, so y has a prefix $g_1 t \ell g'_2 = g'_2 z g_1$, 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. $g_1 > g_2$ ($g = g_1 - g_2$)

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

(a) $g \leq z$

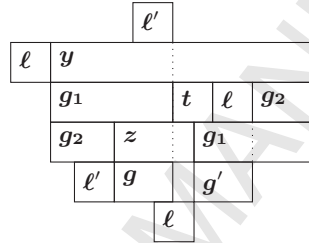


Fig. 18. Subcase 7 when (C3) holds, $g_1 > g_2$ and $g \leq z$

Again there are several steps:

- $l'g$ has period l as a prefix of $l'z$ and shares suffix l' with lg_1 ; accordingly, $l'g$ has border l' and period g , hence periods g and l , thus by Lemma 4 period $\gcd(g, l)$.
- $l'z$ has period l and prefix $l'g$ of period $\gcd(g, l) \mid l$, hence by Lemma 6, it also has period $\gcd(g, l)$.
- g_1 has period g and suffix g of period $\gcd(g, l) \mid g$, so again by Lemma 6, it also has period $\gcd(g, l)$.
- g_1 and $l'z$ have period $\gcd(g, l)$ and share substring g , so that by Lemma 5, g_2z has period $\gcd(g, l)$.
- Prefix lg' of lg_1 shares suffix l with tl , so lg' has border l and period g . Moreover lg' has suffix g' which, as a prefix of g_1 , has period $\gcd(g, l) \mid g$, implying that lg' also has period $\gcd(g, l)$ by Lemma 6.
- g_2z and lg' have period $\gcd(g, l)$ and share substring l , so by Lemma 5, g_2zg' has period $\gcd(g, l)$.
- Then g_2zg' and g_1 have period $\gcd(g, l)$ and share substring g' , so that by Lemma 5 the entire string $y = g_2zg_1$ has period $\gcd(g, l)$.
- Therefore tl and $t'l$ both have period $\gcd(g, l)$ as substrings of y , so that $t = t'$, as required.

Since, as a substring of \mathbf{y} , $\ell\mathbf{g}_1$ has period $\gcd(g, \ell)$, and since \mathbf{g}_1 is also a prefix of \mathbf{y} , it follows from Lemma 5 that $\ell\mathbf{y}$ has period $\gcd(g, \ell)$. Note then that $z\ell = \ell'z$ and $\ell\mathbf{y}$ have period $\gcd(g, \ell)$ and share substring ℓ , so that by Lemma 5 $\mathbf{u}_1 = z\ell\mathbf{y}$ has period $\gcd(g, \ell)$.

(b) $g > z$

\mathbf{y} has period $y - g_1 = g_2 + z$, so \mathbf{y} has a prefix $g_2z\mathbf{g}'_1 = \mathbf{g}'_1t\ell\mathbf{g}_2$, where \mathbf{g}'_1 is a prefix of \mathbf{g}_1 and $|g'_1 - g_2| \leq z$, so one cases 1(a) and 2(a) applies.

(C4)

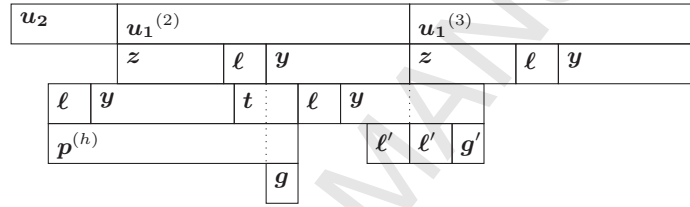


Fig. 19. Subcase 7 when (C4) holds

Suppose (C4) holds; that is, t begins to the left of \mathbf{y} and ends inside it. Let g be the suffix of t that is also a prefix of \mathbf{y} . Let \mathbf{g}' be the suffix of \mathbf{y} of length $g' = g$. By the periodicity of \mathbf{u}' , a copy of $\ell\mathbf{y}$ follows t , extending $\ell + g$ positions into $\mathbf{u}_1^{(3)}$. ℓ' is a suffix of $\ell\mathbf{y}$, so $\ell'\ell'\mathbf{g}'$ is a suffix of $\ell\mathbf{y}$.

The suffix $\ell\mathbf{y}$ of $\mathbf{u}_1^{(2)}$ and the occurrence of $\ell\mathbf{y}$ that follows t are offset by length $\ell + g$, so $\ell\mathbf{y}$ has period $\ell + g$. Since $\ell'\ell'\mathbf{g}'$ has period $\ell + g$ as a suffix of $\ell\mathbf{y}$ and period ℓ as a prefix of $\ell'z$, it therefore has period $\gcd(\ell + g, \ell) = \gcd(g, \ell)$ by Lemma 4. $\ell\mathbf{y}$ has period $\ell + g$ and suffix $\ell'\ell'\mathbf{g}'$ of period $\gcd(g, \ell) \mid \ell + g$, implying that it has period $\gcd(g, \ell)$ by Lemma 6. $z\ell$ has period ℓ and substring ℓ of period $\gcd(g, \ell) \mid \ell$, so by Lemma 6, it has period $\gcd(g, \ell)$. $z\ell$ and $\ell\mathbf{y}$ have period $\gcd(g, \ell)$ and share substring ℓ , so by Lemma 5, $\mathbf{u}_1 = z\ell\mathbf{y}$ has period $\gcd(g, \ell)$. Since $t\ell$ and $t'\ell$ are substrings of \mathbf{u}_1 , we conclude finally that $t = t'$.

This completes the proof of Subcase 7. \square

Lemma 12. *Suppose the conditions of Subcase 7 hold (Lemma 11). Then $\mathbf{x} = \mathbf{d}^{x/d}$, except possibly for*

- (a) $v - u = (h - 1)(u - w)$ 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 11; refer to Figure 12.) Notice that in the proof of the lemma, \mathbf{u}_1 has period ℓ in all cases except when

- (a') (C1) holds and $g = 0$ or
- (b') (C3) holds and $g_1 = g_2$.

Suppose then that $\mathbf{u}_1 = \mathbf{zr} = \mathbf{zly}$ does indeed have period ℓ . Since $\mathbf{z} = \ell't'$ is a prefix of \mathbf{u}_1 , it follows that $t'r = t'ly$ is a suffix of \mathbf{u}_1 . Since \mathbf{u}_1 has period ℓ , $\mathbf{r} = \mathbf{ly}$ has prefix \mathbf{y} , and so \mathbf{u}_1 has border $t'ly$ of length $p = t + \ell + y$, as therefore \mathbf{u} does also. By Lemma 11, \mathbf{u} has period p , so that by Lemma 9, $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$ 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} , $\mathbf{z}p^h$ 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 16 $g_1 = g_2$ 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. □

5 General Case [eb]

In this section we present a sample lemma that considers three squares occurring close to each other, but with no two of them necessarily at the same position.

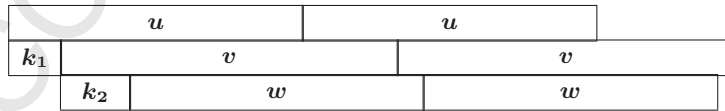


Fig. 20. \mathbf{u}^2 overlapping \mathbf{v}^2 (case (e)) that in turn overlaps \mathbf{w}^2 (case (b)): what is the combined effect?

According to our computational experiments, the following lemma applies to about three-quarters of cases in which the maximum alphabet size $\sigma = \gcd(u, v, w)$. Figure 20 illustrates case [eb].

Lemma 13. *In case [eb], 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 (e) of Lemma 8, 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 the Periodicity Lemma [11], 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 (e), $\frac{2(k_1 + v)}{3} \leq u$ 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. To date, all the results given for the NPL in [9, 10, 30, 21, 13] deal only with the special cases [ij], $i = a, d$, that arise for $k_1 = 0$.

6 Commentary & Future Research

The proofs of Lemmas 3 and 8 may, as indicated in the Introduction, open the way for new and more combinatorially sophisticated approaches to the bounding of the number of runs in any string of given length. Of course it would be highly desirable to find an approach to this lemma that would permit a much simplified and more suggestive proof.

In addition, much work remains to be done to state and prove results such as Lemma 13. In proving the results of [21], it turned out to be very helpful to look at the results of computer simulations for small values of k, u, v, w . It seems that similar techniques can profitably be used to generate conjectures for the cases [ij] of three overlapping squares that arise from Lemma 8.

More generally, once the combinatorics of overlapping squares is well understood, it may well be possible to begin to design an algorithmic approach to the computation of runs that handles the various cases that arise without the need for elaborate preprocessing. This would for the first time permit direct analysis and computation of the local periodicities of a string in a manner consistent with their sparseness of occurrence.

Acknowledgements

This work was supported in part by the Natural Sciences & Engineering Research Council of Canada. The authors would like to thank two anonymous referees for their careful and constructive comments.

References

1. ALBERTO APOSTOLICO & FRANCO P. PREPARATA, **Optimal off-line detection of repetitions in a string**, *Theoret. Comput. Sci.* 22 (1983) 297–315.
2. HIDEO BANNAI, TOMOHIRO I, SHUNSUKE INENAGA, YUTO NAKASHIMA, MASAYUKI TAKEDA & KAZUYA TSURUTA, **The “Runs” Theorem**, [arXiv:1406.0263](https://arxiv.org/abs/1406.0263) (2014).
3. GANG CHEN, SIMON J. PUGLISI & W. F. SMYTH, **Fast & practical algorithms for computing all the runs in a string**, *Proc. 18th Annual Symp. Combinatorial Pattern Matching*, B. Ma & K. Zhang (eds.), LNCS 4580, Springer-Verlag (2007) 307–315.
4. GANG CHEN, SIMON J. PUGLISI & W. F. SMYTH, **Lempel-Ziv factorization using less time & space**, *Math. in Computer Science1–4* (2008) 605–623.
5. MAXIME CROCHEMORE, **An optimal algorithm for computing all the repetitions in a word**, *Inform. Process. Lett.* 12–5 (1981) 244–248.
6. MAXIME CROCHEMORE & LUCIAN ILIE, **Maximal repetitions in strings**, *J. Comput. Sys. Sci.* (2008) 796–807.
7. MAXIME CROCHEMORE, LUCIAN ILIE & LIVIU TINTA, **Towards a solution to the “runs” conjecture**, *Proc. 19th Annual Symp. Combinatorial Pattern Matching*, P. Ferragina & G. Landau (eds.), LNCS 5029, Springer-Verlag (2008) 290–302.
8. MAXIME CROCHEMORE, LUCIAN ILIE & LIVIU TINTA, **The “runs” conjecture**, *TCS 412–27* (2011) 2931–2941.
9. MAXIME CROCHEMORE AND WOJCIECH RYTTER, **Squares, cubes, and time-space efficient strings searching**, *Algorithmica* 13 (1995) 405–425.
10. KANGMIN FAN, SIMON J. PUGLISI, W. F. SMYTH & ANDREW TURPIN, **A new periodicity lemma**, *SIAM J. Discrete Math.* 20–3 (2006) 656–668.
11. N. J. FINE AND H. S. WILF, **Uniqueness theorems for periodic functions**, *Proc. Amer. Math. Soc.* 16 (1965) 109–114.
12. AVIEZRI S. FRAENKEL & JAMIE SIMPSON, **The exact number of squares in Fibonacci words**, *Theoret. Comput. Sci.* 218–1 (1999) 95–106.
13. FRANTISEK FRANEK, ROBERT C. G. FULLER, JAMIE SIMPSON & W. F. SMYTH, **More results on overlapping squares**, *J. Discrete Algorithms* 17 (2012) 2–8.
14. FRANTISEK FRANEK & Q. YANG, **An asymptotic lower bound for the maximal number of runs in a string**, *Internat. J. Foundations of Computer Science19–1* (2008) 195–203.
15. FRANTISEK FRANEK, R. J. SIMPSON & W. F. SMYTH, **The maximum number of runs in a string**, *Proc. 14th Australasian Workshop on Combinatorial Algs.*, Mirka Miller & Kunsoo Park (eds.) (2003) 26–35.
16. MATHIEU GIRAUD, **Not so many runs in strings**, *Proc. 2nd Internat. Conf. on Language & Automata Theory & Applications*, Carlos Martín-Vide, Friedrich Otto & Henning Fernau (eds.), LNCS 5196, Springer-Verlag (2008) 232–239.
17. MATHIEU GIRAUD, **Asymptotic behavior of the numbers of runs and microruns**, *Inform. & Computation* 207–11 (2009) 1221–1228.

18. COSTAS S. ILIOPOULOS & W. F. SMYTH, **A characterization of the squares in a Fibonacci string**, *Theoret. Comput. Sci.* 172 (1997) 281–291.
19. ROMAN KOLPAKOV & GREGORY KUCHEROV, **On maximal repetitions in words**, *J. Discrete Algorithms* 1 (2000) 159–186.
20. KAZUHIKO KUSANO, KAZUYUKI NARISAWA & AYUMI SHINOHARA, **On morphisms generating run-rich strings**, *Proc. Prague Stringology Conf.2013* (2013) 35–47.
21. EVGUENIA KOPYLOVA & W. F. SMYTH, **The three squares lemma revisited**, *J. Discrete Algorithms* 11 (2012) 3–14.
22. M. LOTHAIRE, *Algebraic Combinatorics on Words*, Cambridge University Press (2002) 504 pp.
23. MICHAEL G. MAIN, **Detecting leftmost maximal periodicities**, *Discrete Applied Maths.* 25 (1989) 145–153.
24. MICHAEL G. MAIN & RICHARD J. LORENTZ, **An $O(n \log n)$ algorithm for finding all repetitions in a string**, *J. Algorithms* 5 (1984) 422–432.
25. WATARU MATSUBARA, KAZUHIKO KUSANO, AKIRA ISHINO, HIDEO BANNAI & AYUMI SHINOHARA, **New lower bounds for the maximum number of runs in a string**, *PSC* (2008) 140–145.
26. SIMON J. PUGLISI & R. J. SIMPSON, **The expected number of runs in a word**, *Australasian J. Combinatorics* 42 (2008) 45–54.
27. SIMON J. PUGLISI, R. J. SIMPSON & W. F. SMYTH, **How many runs can a string contain?**, *Theoret. Comput. Sci.* 401 (2008) 165–171.
28. WOJCIECH RYTTER, **The number of runs in a string: improved analysis of the linear upper bound**, *Proc. 23rd Symp. Theoretical Aspects of Computer Science*, B. Durand & W. Thomas (eds.), LNCS 2884, Springer-Verlag (2006) 184–195.
29. SHARCNET, <https://www.sharcnet.ca/my/front/>
30. R. J. SIMPSON, **Intersecting periodic words**, *Theoret. Comput. Sci.* 374 (2007) 58–65.
31. JAMIE SIMPSON, **Modified Padovan words and the maximum number of runs in a word**, *Australasian J. Combinatorics* 46 (2010) 129–145.
32. BILL SMYTH, *Computing Patterns in Strings*, Pearson Addison-Wesley (2003) 423pp.