

A NEW PERIODICITY LEMMA*

KANGMIN FAN[†], SIMON J. PUGLISI[‡], W. F. SMYTH^{†‡}, AND ANDREW TURPIN[§]

Abstract. Given a string $x = x[1..n]$, a *repetition* of period p in x is a substring $u^r = x[i..i+rp-1]$, $p = |u|$, $r \geq 2$, where neither $u = x[i..i+p-1]$ nor $x[i..i+(r+1)p-1]$ is a repetition. The maximum number of repetitions in any string x is well known to be $\Theta(n \log n)$. A *run* or *maximal periodicity* of period p in x is a substring $u^r t = x[i..i+rp+|t|-1]$ of x , where u^r is a repetition, t is a proper prefix of u , and no repetition of period p begins at position $i-1$ of x or ends at position $i+rp+|t|$. In 2000 Kolpakov and Kucherov [*J. Discrete Algorithms*, 1 (2000), pp. 159–186] showed that the maximum number $\rho(n)$ of runs in any string x is $O(n)$, but their proof was nonconstructive and provided no specific constant of proportionality. At the same time, they presented experimental data strongly suggesting that $\rho(n) < n$. Related work by Fraenkel and Simpson [*J. Combin. Theory Ser. A.*, 82 (1998), pp. 112–120] showed that the maximum number $\sigma(n)$ of *distinct* squares in any string x satisfies $\sigma(n) < 2n$, while experiment again encourages the belief that in fact $\sigma(n) < n$. In this paper, as a first step toward proving these conjectures, we present a periodicity lemma that establishes limitations on the number and range of periodicities that can occur over a specified range of positions in x . We then apply this result to specify corresponding limitations on the occurrence of runs.

Key words. string, word, periodicity, square, repetition, run, maximal periodicity

AMS subject classification. 68R15

DOI. 10.1137/050630180

1. Introduction. The study of strings began with an investigation of periodicity properties [23], and periodicity of various kinds still remains a central theme, important both in theory and practice—for example, in data compression, pattern matching, computational biology, and many other areas. In this paper we present results that specify restrictions on the nature and extent of periodic behavior in strings. Although these results are theoretical, their importance is very much a product of their practical application, as we explain below.

It will be convenient throughout to represent strings in boldface (for example, $x = \mathbf{x}[1..n]$) and their lengths in italics (for example, $x = |x|$).

If $w = u^r$ for some nonempty string u and some integer $r \geq 2$, then w is said to be a *repetition*. Further, a *repetition in x* is a substring $u^r = x[i..i+ru-1]$, $r \geq 2$, in x , where $x[i..i+u-1]$ is not a repetition and $x[i..i+(r+1)u-1] \neq u^{r+1}$. We call u the *generator* of the repetition, u its *period*, and r its *exponent*; and we represent it economically by an integer triple (i, u, r) . In the early 1980s three quite different

*Received by the editors April 28, 2005; accepted for publication (in revised form) March 16, 2006; published electronically September 15, 2006. Preliminary versions of parts of this paper appeared in *Proceedings of the 16th Annual Symposium on Combinatorial Pattern Matching*, Lecture Notes in Comput. Sci. 3537, Springer-Verlag, Berlin, 2005, and in *Proceedings of the 16th Australasian Workshop on Combinatorial Algorithms*, University of Ballarat, Ballarat, Victoria, Australia, 2005. <http://www.siam.org/journals/sidma/20-3/63018.html>

[†]Algorithms Research Group, Department of Computing and Software, McMaster University, Hamilton, ON L8S 4K1, Canada (fank@mcmaster.ca, smyth@mcmaster.ca, www.cas.mcmaster.ca/cas/research/algorithms.htm). The first and third authors were supported in part by grants from the Natural Sciences and Engineering Research Council of Canada.

[‡]Department of Computing, Curtin University, GPO Box U1987, Perth WA 6845, Australia (puglisi@computing.edu.au, smyth@computing.edu.au).

[§]Department of Computer Science and Information Technology, RMIT University, GPO Box 2476V, Melbourne V 3001, Australia (aht@cs.rmit.edu.au). The fourth author was supported by a grant from the Australian Research Council.

$O(x \log x)$ algorithms were published [2, 1, 17] for the computation of all the repetitions in a given string \mathbf{x} . In a sense these algorithms were all asymptotically optimal, since in [2] it was shown that in fact a Fibonacci string \mathbf{f}_n contains $\Theta(f_n \log f_n)$ repetitions.

In [16] Main introduced a more compact encoding of repetitions: a *run* or *maximal periodicity* of period u in \mathbf{x} was defined to be a substring $\mathbf{u}^r \mathbf{t} = \mathbf{x}[i..i+ru+t-1]$ of \mathbf{x} , where \mathbf{u}^r is a repetition, \mathbf{t} is a proper prefix of \mathbf{u} , and no repetition of period u begins at position $i-1$ of \mathbf{x} or ends at position $i+ru+t$. \mathbf{u} is called the *generator* of the run, \mathbf{t} is called its *tail*, and a run is economically represented by a 4-tuple (i, u, r, t) . Computing all the runs in \mathbf{x} permits all the repetitions in \mathbf{x} to be listed in an obvious way. Main [16] showed how to compute all the “leftmost” runs in \mathbf{x} in time $\Theta(x)$, provided that the suffix tree [24, 18] and the Lempel–Ziv (LZ) factorization [14] of \mathbf{x} were both available. In [4] it was shown that a suffix tree could be computed in linear time on an *indexed* (bounded integer) alphabet; since the LZ factorization is computable in linear time from the suffix tree, this meant that the overall worst-case time requirement of Main’s algorithm was $\Theta(x)$ on an indexed alphabet. In [13] Kolpakov and Kucherov took matters a step further by extending Main’s algorithm to also compute nonleftmost runs in \mathbf{x} in time proportional to their number, and then by showing that the maximum number $\rho(x)$ of runs in any string \mathbf{x} was at most

$$(1) \qquad k_1 x - k_2 \log_2 x \sqrt{x},$$

where k_1 and k_2 are positive constants. Thus, at least in principle, all the runs in \mathbf{x} could be determined in linear time.

However, there is a problem with (1): The proof is nonconstructive and gives no information about the magnitude of the constants k_1 and k_2 . Nevertheless Kolpakov and Kucherov provide convincing experimental evidence that

- * $\rho(x) < x$;
- * $\rho(x)$ is achieved by a cube-free string \mathbf{x} on alphabet $\{a, b\}$;
- * $\rho(x + 1) \leq \rho(x) + 2$.

As far as we know, there are only two published works that address these fundamental questions of periodicity. In [7] an infinite family of strings \mathbf{x} is constructed that is conjectured for sufficiently large x to achieve $\rho(x) < x$. This family thus provides a lower bound on $\rho(x)$. More recently, Rytter [21] has used interesting techniques to show that $\rho(n) \leq 5n$, thus establishing an upper bound.

It was mentioned above that Main’s algorithm computes all the leftmost runs in \mathbf{x} , that is, the leftmost occurrence of each distinct run, a collection that certainly includes the leftmost occurrence of each distinct square in \mathbf{x} . This suggests a connection with another well-known problem: the determination of $\sigma(x)$, the maximum number of distinct squares in any string \mathbf{x} , where again experiment strongly suggests that $\sigma(x) < x$. With this problem better progress has been made: Fraenkel and Simpson showed [6] that $\sigma(x) \leq 2x - 2$, a result recently proved somewhat more simply by Ilie [8], then later improved to $\sigma(x) \leq 2x - \Theta(\log x)$ [9].

In order to show that in general $\rho(x) < x$ ($\sigma(x) < x$), it seems to be necessary to establish restrictions on the number of runs (squares) that can occur near a position in \mathbf{x} at which one or two runs (squares) are already known to occur. Perhaps the most famous theoretical result available for such a purpose is the following “periodicity lemma.”

LEMMA 1 (see [5]). *Let p and q be two periods of \mathbf{x} , and let $d = \gcd(p, q)$. If $p+q \leq x+d$, then d is also a period of \mathbf{x} .*

Unfortunately this lemma provides no special information about runs or the squares with which runs must begin, and it places no restrictions on the positions at which periodic substrings may occur. To our knowledge the only result that provides such information is the following “three squares lemma.”

LEMMA 2 (see [3, 15]). *Suppose \mathbf{u} is not a repetition, and suppose $\mathbf{w} \neq \mathbf{u}^j$ for any $j \geq 1$. If \mathbf{u}^2 is a prefix of \mathbf{w}^2 , in turn a proper prefix of \mathbf{v}^2 , then $w \leq v - u$.*

Our main result in this paper is essentially a generalization of this result, which we call a “new periodicity lemma”: We allow \mathbf{w} to be offset by k positions from the start of \mathbf{v}^2 , and we do not always require complete squares \mathbf{v}^2 and \mathbf{w}^2 , only sufficiently long substrings of periods v and w . Moreover, as a corollary of our main result, we are able to specify exactly the periodic behavior in the string.

2. New periodicity lemma. In this section we prove results that establish restrictions on the squares that can occur in the neighborhood of positions in a string at which one or two squares already appear. We begin with three simple definitions.

DEFINITION 3. *A square \mathbf{u}^2 is said to be irreducible if \mathbf{u} is not a repetition.*

DEFINITION 4. *A square \mathbf{u}^2 is said to be regular if no prefix of \mathbf{u} is a square.*

DEFINITION 5. *A square \mathbf{u}^2 is said to be minimal if no proper prefix of \mathbf{u}^2 is a square.*

LEMMA 6. *If \mathbf{u}^2 is minimal, then \mathbf{u}^2 is regular; if \mathbf{u}^2 is regular, then \mathbf{u}^2 is irreducible.*

Proof. The proof of the first statement is immediate. To prove the second, observe that by Definition 4, no prefix of \mathbf{u} is a square. Therefore \mathbf{u} cannot be a repetition, and so by Definition 3 \mathbf{u}^2 is irreducible. \square

The existence of a minimal square already imposes significant limitations on the nature of other squares that can exist, as the following result shows.

LEMMA 7. *If $\mathbf{x} = \mathbf{u}^2$ is minimal, then for all integers $k \geq 0$ and $w \in u/2..u-1$,*

(a) *if*

$$(2) \quad k + w \leq u, \quad k + 3w \geq 2u,$$

$\mathbf{x}[k+1..k+2w]$ *is not a square;*

(b) *if*

$$(3) \quad k + w > u, \quad k + 2w \leq 2u,$$

either $\mathbf{x}[k+1..k+2w]$ is not a square or $\mathbf{x}[w'+1..w'+u]$ has period $u-w$, where

$$w' = (k+w) - u.$$

Proof. Suppose that for some pair of integers k and w satisfying either (2) or (3), $\mathbf{x}[k+1..k+2w] = \mathbf{w}^2$.

First assume that $k = 0$. Then if (2) holds, either $w = u$, a contradiction, or else $w < u$, contradicting the minimality of \mathbf{u}^2 . On the other hand, if (3) holds, then both $w > u$ and $w \leq u$ must hold, again a contradiction. Thus we can assume that $k \geq 1$.

(a) Suppose that (2) holds, let $w' = u - (k+w)$, and consider

$$\widehat{\mathbf{w}} = \mathbf{x}[1..w-w'] = \mathbf{x}[k+w'+1..k+w].$$

Since by (2)

$$(w-w') - (k+w') = k+3w-2u \geq 0,$$

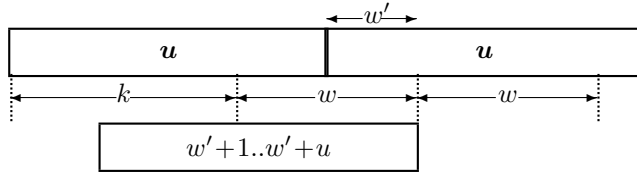


FIG. 1. Lemma 7(b).

the substring $\mathbf{x}[1..k+w]$ has period $k+w'$. Again by (2),

$$(k+w) - 2(k+w') = (k+w) - 2(u-w) \geq 0,$$

so that $\mathbf{x}[1..k+w]$ has prefix $(\mathbf{x}[1..k+w'])^2$, contradicting the minimality of \mathbf{u}^2 . Thus in case (a) no such k and w can exist.

(b) Next we suppose that (3) holds so that $w < u$ and hence that $k-w' = u-w > 0$ (see Figure 1).

Consider

$$\mathbf{w} = \mathbf{x}[w'+1..w'+w] = \mathbf{x}[k+1..u+w'].$$

Since by (3) $w'+w = k+(2w-u) \geq k$, the substring $\mathbf{x}[w'+1..w'+u]$ of length u has period $k-w' = u-w$, as required. \square

To show that in case (a) of Lemma 7 the assumption that $k+3w \geq 2u$ (as well as the weaker condition $w \geq u/2$) is necessary, consider the example $u = 14, k = 6, w = 5$:

$$\mathbf{x} = \mathbf{u}^2 = \text{abbaba}(\text{babab})(\text{bab}||\text{ab})(\text{babab})\text{ababbab}.$$

Here $\mathbf{w} = \text{babab}$, and \mathbf{w}^3 is a substring of \mathbf{x} .

To show that in case (b) of Lemma 7 the substring \mathbf{w}^2 can in fact exist, consider the example $u = 11, k = 4, w = 8$ with $w' = 1$:

$$\mathbf{x} = \mathbf{u}^2 = \text{babc}(\text{abcabca}||\text{b})(\text{abcabca})\text{ca}.$$

The substring $\mathbf{x}[2..12] = (\text{abc})^3\text{ab}$ has period $u-w = 3$.

We turn now to the situation in which a regular square and an irreducible square occur at the same position. We first prove two basic lemmas that describe the relationship between regularity and irreducibility, and then go on to prove our main result.

LEMMA 8. *If \mathbf{v}^2 is irreducible with regular proper prefix \mathbf{u}^2 , then*

$$v > \max\{u+1, 3u/2\}.$$

Proof. Observe that $1 \leq u < v$, and observe further that $u+1 \geq 3u/2$ if and only if $u \leq 2$.

For $u = 1$, $\mathbf{u}^2 = \lambda^2$ for some letter λ and the shortest irreducible square $\mathbf{v}^2 = (\lambda^2\mu)^2$ for some letter $\mu \neq \lambda$. Thus for $u = 1$, $v \geq 3 > u+1$, as required.

For $u = 2$, since \mathbf{u}^2 is regular, $\mathbf{u}^2 = (\lambda\mu)^2$ and the shortest irreducible square $\mathbf{v}^2 = (\lambda\mu\lambda\mu\nu)^2$ for some letter ν . Thus for $u = 2$, $v \geq 5 > u+1$, as required.

Suppose therefore that $u \geq 3$, and suppose further, without loss of generality, that $v < 2u$. Then

$$\mathbf{v} = \mathbf{uu}[1..v-u] = \mathbf{u}[v-u+1..u]\mathbf{v}[2u-v+1..v],$$

where $\mathbf{y} = \mathbf{u}[1..v-u]$ of length $v-u$ is a prefix of \mathbf{u} , and hence of \mathbf{v} , and $\mathbf{z} = \mathbf{u}[v-u+1..u]$ of length $2u-v$ is a prefix of \mathbf{v} , and hence of \mathbf{u} . If we now assume $2v \leq 3u$, it follows that $v-u \leq 2u-v$, so that \mathbf{y} is also a prefix of \mathbf{z} . Thus \mathbf{u} has prefix \mathbf{y}^2 and so \mathbf{u}^2 cannot be regular, a contradiction. We conclude that $2v > 3u$, as required. \square

Observe that if \mathbf{u}^2 is not regular, Lemma 8 may not hold: $\mathbf{u} = ababa$ allows $\mathbf{v} = ababaab$ with $v < 3u/2$.

LEMMA 9. *If $\mathbf{x} = \mathbf{v}^2$ is irreducible with regular proper prefix \mathbf{u}^2 , $v < 2u$, then*

$$\mathbf{x} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_1\mathbf{u}_2,$$

where $u_1 = 2u-v$, $u_2 = 2v-3u$.

Proof. Since $v < 2u$, $u \geq 3$ by Lemma 8. Let \mathbf{u}_1 be the suffix of \mathbf{u} of length $u_1 = 2u-v$ that is a prefix of \mathbf{v} , and hence also a prefix of \mathbf{u} . By the regularity of \mathbf{u} and Lemma 8, $u_1 < u/2$ and so $\mathbf{u} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1$ for some nonempty \mathbf{u}_2 . Then $\mathbf{v} = \mathbf{u}_1\mathbf{u}_2\mathbf{u}_1\mathbf{u}_1\mathbf{u}_2$, so that $u_2 = 2v-3u$, as required. \square

For the proof of our main result, the following definitions will be helpful. If $\mathbf{x} = \mathbf{u}\mathbf{v}$, \mathbf{v} nonempty, then $\mathbf{v}\mathbf{u} = R_u(\mathbf{x})$ is said to be the u th rotation of \mathbf{x} ; also, if \mathbf{u} is both a proper prefix and a suffix of \mathbf{x} , then it is said to be a border of \mathbf{x} .

We frequently make use of the following two well-known results.

LEMMA 10 (see [22, p. 76]). *Let \mathbf{x} be a string of length n and minimum period p , and let $j \in 1..n-1$ be an integer. Then $R_j(\mathbf{x}) = \mathbf{x}$ if and only if \mathbf{x} is a repetition and p divides j .*

LEMMA 11 (see [22, p. 76]). *If a string \mathbf{x} is a repetition, then so is every rotation of \mathbf{x} .*

We first state the new periodicity lemma (NPL) in a rather general and easily understood form: Having gone through the proof, we will then be able to reexpress it to yield stronger conclusions based on weaker premises. A total of 14 cases arise in the proof (see Table 1). For each of these cases, we are able to identify a specific square prefix of \mathbf{u} that is forced by the presence of \mathbf{w}^2 in the string \mathbf{x} , thus contradicting the assumption that \mathbf{u}^2 is regular; therefore, if \mathbf{u}^2 is not regular, the square prefix must exist.

For each of the main cases, we specify the range of values of k (either $k \in 0..u_1$ or $k \in u_1+1..u_1+u_2-1$) and the end position of $\mathbf{w}^{(1)}$ (first occurrence of \mathbf{w}) in \mathbf{x} . To facilitate this latter task, we introduce the notation $\mathbf{u}_1^{(j)}, \mathbf{u}_2^{(j)}$ to denote the j th occurrence of $\mathbf{u}_1, \mathbf{u}_2$, respectively, in \mathbf{x} . Thus “ $\mathbf{w}^{(1)}$ ends in $\mathbf{u}_2^{(2)}$ ” means that the first occurrence of \mathbf{w} in \mathbf{x} ends in the second occurrence of \mathbf{u}_2 in \mathbf{x} . In most of the cases, it is useful to introduce a substring \mathbf{s} that is both a prefix of \mathbf{w} and a suffix of one of the substrings $\mathbf{u}_1^{(j)}$ or $\mathbf{u}_2^{(j)}$ in which $\mathbf{w}^{(1)}$ ends.

LEMMA 12 (NPL). *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.*

Proof. Suppose instead that for some k and w , $\mathbf{w}^2 = \mathbf{x}[k+1..k+2w]$. Recall the definitions of \mathbf{u}_1 and \mathbf{u}_2 given in Lemma 9, with $u_1+u_2 = v-u$.

A. $k \leq u_1$.

I. $\mathbf{w}^{(1)}$ ends in $\mathbf{u}_1^{(2)}$ ($k+w \leq u$, $s = u-(k+w)$).

(a) $\mathbf{w}^{(2)}$ ends in $\mathbf{u}_1^{(3)}$ ($k+2w \leq u+u_1$) (see Figure 2).

Define $\mathbf{q} = \mathbf{u}_1[1..q]$ and $\mathbf{z} = \mathbf{u}_1[1..z]$, which are both prefixes of \mathbf{u}_1 and suffixes of \mathbf{w} :

$$q = u_1 - s = k + w - (u_1 + u_2), \quad z = k + 2w - u.$$

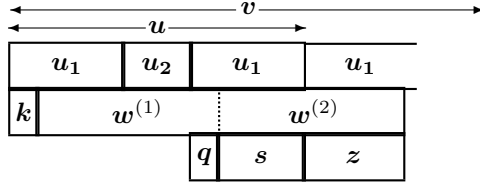


FIG. 2. Case A.I.(a).

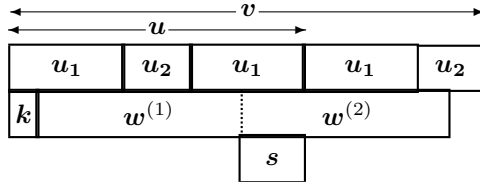


FIG. 3. Case A.I.(b).

Observe that

$$q - k = w - (u_1 + u_2) > 0, \quad z - q = w - u_1 > 0.$$

Since $q < z$, q is a border of z , and thus z has period $z - q$.

(i) $q \geq z/2$ ($k \geq u_2$). Here z , and hence u_1 , has prefix

$$z[1..z - q]^2 = z[1..w - u_1]^2,$$

contradicting the regularity of u^2 .

(ii) $q < z/2$ ($k < u_2$). Here we can set $z = qpq$, where $p > 0$. Since $q > k$, we can also set $q = kt$, where, as noted above, $t = w - (u_1 + u_2) > 0$. Hence $z = ktpkt = ktr$ for $r = pkt$.

Observe now that $tpkt$ is a prefix of $w^{(1)}$, while r is a prefix of $w^{(2)}$. Thus $r = R_t(r)$, so that by Lemmas 10 and 11, r and all of its rotations are repetitions of period t . It follows that z , a prefix of u_1 , is a repetition of period $t = w - (u_1 + u_2)$ and exponent at least 3, contradicting the regularity of u^2 .

(b) $w^{(2)}$ ends in $u_2^{(2)}$ ($k + 2w > u + u_1$) (see Figure 3).

Since $w > u_1 + u_2$, $k + s < u_1$; since $w < u$, $k + s > 0$. Therefore ks is a prefix of u_1 , and since su_1 is a prefix of w , it follows that u has prefix $(ks)^2$, $k + s = u - w$, contradicting the regularity of u^2 .

II. $w^{(1)}$ ends in $u_1^{(3)}$ ($k + w \leq u + u_1$, $s = u + u_1 - (k + w)$) (see Figure 4).

Since $w \neq u$, $k + s \neq u_1$. Observe that $w^{(1)}$ has prefix $R_k(u_1u_2)$, while $w^{(2)}$ has prefix $R_{u_1-s}(u_1u_2)$. Since $u_1 - s \neq k$ (otherwise $w = u$), it follows from Lemma 10 that u_1u_2 is a repetition of period $|k - (u_1 - s)| = |u - w|$, contradicting the regularity of u^2 . Note that if u^2 is not regular, then u must also have period $|u - w|$.

III. $w^{(1)}$ ends in $u_2^{(2)}$ ($k + w \leq v$, $s = v - (k + w)$, $k + s > 0$) (see Figure 5).

$w^{(1)}$ has prefix $R_k(u_1u_2)$, while $w^{(2)}$ has prefix $R_t(u_1u_2)$, where $t = u_1 + u_2 - s$. Since $t = k + w - u > k$, it follows from Lemma 10 that u_1u_2 is a repetition of period $t - k = w - u$, contradicting the regularity of u^2 .

Note that if u^2 is not regular, then u must also have period $w - u$.

IV. $w^{(1)}$ ends in $u_1^{(4)}$ ($k + w \leq 2u$, $s = 2u - (k + w)$) (see Figure 6).

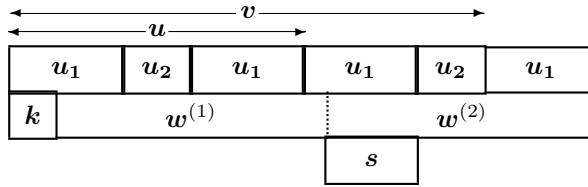


FIG. 4. Case A.II.

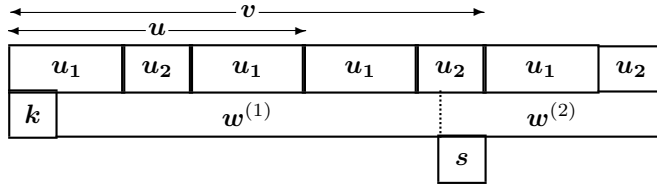


FIG. 5. Case A.III.

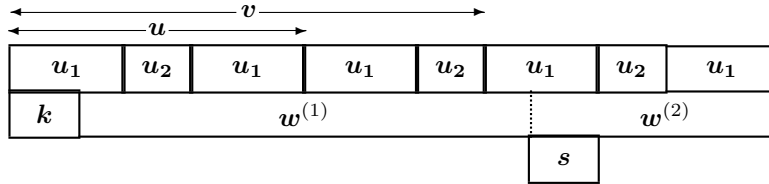


FIG. 6. Case A.IV.

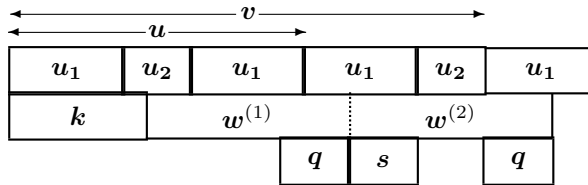


FIG. 7. Case B.I(a).

As in case A.III., $w^{(1)}$ has prefix $R_k(u_1u_2)$, while $w^{(2)}$ has prefix $R_{u_1-s}(u_1u_2)$. It follows from Lemma 10 that u_1u_2 is a repetition of period $k+s-u_1 = v-w$, contradicting the regularity of u^2 . Note that if u^2 is not regular, then u must also have period $v-w$.

B. $k > u_1$.

I. $w^{(1)}$ ends in $u_1^{(3)}$ ($k+w \leq u+u_1$, $s = u+u_1-(k+w)$, $k+s < 2u_1$).

(a) $w^{(2)}$ ends in $u_1^{(4)}$ ($k+2w \leq 2u$) (see Figure 7).

Let q be the prefix of u_1 and suffix of $w^{(2)}$ defined by

$$q = w - u_2 - s = k + 2w - v;$$

then, because it is a prefix of u_1 , q occurs at position $u+1$ of x and, because it is a suffix of $w^{(1)}$, also at position $k+w-q+1$. These two copies of q are offset by period

$$t = u+q-(k+w) = w+u-v.$$

Since

$$\begin{aligned} q-2t &= k+2w-v-2w-2u+2v \\ &= k+v-2u > 0, \end{aligned}$$

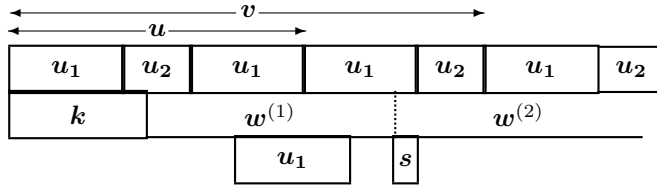


FIG. 8. Case B.I.(b).

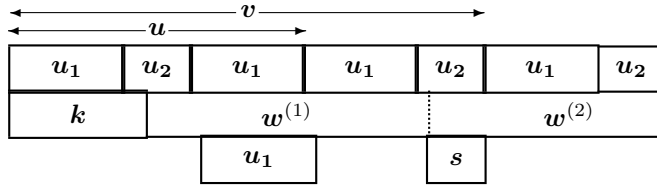


FIG. 9. Case B.II.(a).

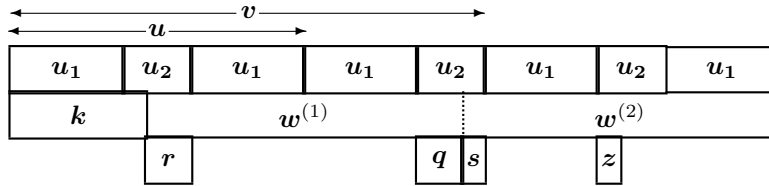


FIG. 10. Case B.II.(b)(i).

therefore q , and hence u_1 , has a square prefix of period $w + u - v$, contradicting the regularity of u^2 .

- (b) $w^{(2)}$ ends in $u_2^{(3)}$ ($k + 2w > 2u$) (see Figure 8).

Observe that since u_1 occurs at position $s + u_2 + 1$ in $w^{(2)}$, and since u_1^2 begins at position $u_1 + u_2 - k + 1$ in $w^{(1)}$, therefore $u_1 = R_t(u_1)$ for $t = k + s - u_1 = u - w$. Hence by Lemma 10, u_1 is a repetition of period $u - w$, contradicting the regularity of u^2 .

- II. $w^{(1)}$ ends in $u_2^{(2)}$ ($k + w \leq v$, $s = v - (k + w)$, $k + s \neq u_1 + u_2$).

- (a) $w < u$ ($k + s > u_1 + u_2$) (see Figure 9).

Observe that u_1^2 occurs at position $u_1 + u_2 - k + 1$ in $w^{(1)}$, while u_1 occurs at position $s + 1$ in $w^{(2)}$. Since $s > u_1 + u_2 - k$, this means that $u_1 = R_t(u_1)$ for $t = k + s - (u_1 + u_2) = u - w > 0$. Hence u_1 is a repetition of period $u - w$, contradicting the regularity of u^2 .

- (b) $w > u$ ($k + s < u_1 + u_2$).

- (i) $w^{(2)}$ ends in $u_1^{(5)}$ ($k + 2w \leq v + u$, $w - s \leq u$) (see Figure 10).

Let $r = u_2[k - u_1 + 1..u_2]$, where $r = u_1 + u_2 - k$ and $r - s = w - u > 0$. Observe that $w^{(1)} = (ru_1)(u_1q)$, where $q = u_2[1..u_2 - s]$. Also $w^{(2)}$ has prefix su_1z , where $z = u_2[1..r - s]$, of length $r + u_1$. Since $w - s \leq u$, the copy of u that begins at position $v + 1$ of x has prefix $(u_1z)(u_1q)$, where $q - z = u_2 - r > 0$. Thus u has prefix $(u_1z)^2$ of period $u_1 + r - s = w - (u_1 + u_2)$, contradicting the regularity of u^2 .

- (ii) $w^{(2)}$ ends in $u_1^{(6)}$ ($k + 2w \leq v + u + u_1$, $u < w - s \leq u + u_1$) (see Figure 11).

Observe that $w^{(1)}$ has suffix $R_{u_2 - s}(u_2u_1^2)$, while $w^{(2)}$ has suffix

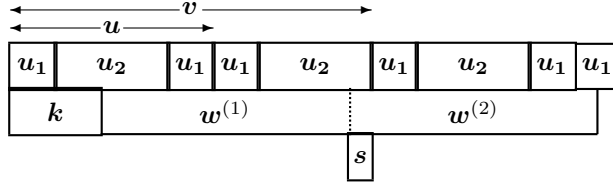


FIG. 11. Case B.II.(b)(ii).

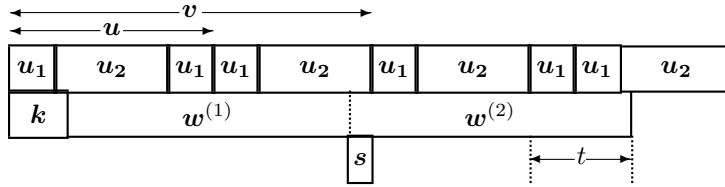


FIG. 12. Case B.II.(b)(iii).

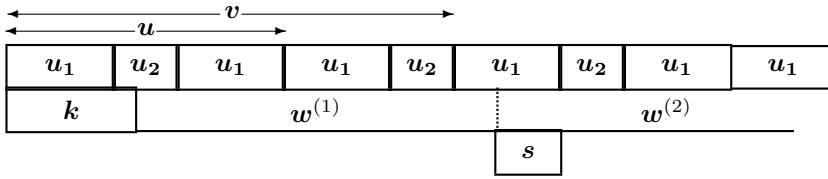


FIG. 13. Case B.III.

$R_t(\mathbf{u}_2\mathbf{u}_1^2)$, where $t = u_1 + u_2 + w - s - u$. Since $t - (u_2 - s) = w - (u_1 + u_2) > 0$, it follows from Lemmas 10 and 11 that $\mathbf{u}_2\mathbf{u}_1^2$, and hence \mathbf{u} , is a repetition of period $w - (u_1 + u_2)$, contradicting the regularity of \mathbf{u}^2 .

(iii) $\mathbf{w}^{(2)}$ ends in $\mathbf{u}_2^{(4)}$ ($k + 2w < 2v$, $u + u_1 < w - s < v$) (see Figure 12).

As in case B.II.(b)(ii), $\mathbf{w}^{(1)}$ has suffix $R_{u_2-s}(\mathbf{u}_2\mathbf{u}_1^2)$, while now $\mathbf{w}^{(2)}$ has suffix $R_t(\mathbf{u}_2\mathbf{u}_1^2)$, where $t = w - s - (u + u_1) > 0$. Since $u_2 - s - t = v - w > 0$, it follows from Lemmas 10 and 11 that $\mathbf{u}_2\mathbf{u}_1^2$, and hence \mathbf{u} , is a repetition of period $v - w$, contradicting the regularity of \mathbf{u}^2 .

III. $\mathbf{w}^{(1)}$ ends in $\mathbf{u}_1^{(4)}$ ($k + w \leq 2u$, $s = 2u - (k + w)$, $k + s < u$) (see Figure 13).

Observe that $\mathbf{w}^{(1)}$ has prefix $R_k(\mathbf{u})$, while $\mathbf{w}^{(2)}$ has prefix $R_{u_1-s}(\mathbf{u})$. Since $u_1 < k + s$, it follows by Lemma 10 that \mathbf{u} is a repetition of period $k + s - u_1 = v - w$, contradicting the regularity of \mathbf{u}^2 .

If $\mathbf{w}^{(2)}$ extends only to the end of $\mathbf{u}_1^{(5)}$, the argument of case B.II.(a) can instead be used to show that \mathbf{u}_1 is a repetition of period $v - w$, again contradicting the regularity of \mathbf{u}^2 .

IV. $\mathbf{w}^{(1)}$ ends in $\mathbf{u}_2^{(3)}$ ($k + w \leq 2u + u_2$, $s = 2u + u_2 - (k + w)$) (see Figure 14).

The arguments of case B.III. apply: \mathbf{u} (or \mathbf{u}_1) is a repetition of period $v - w$, contradicting the regularity of \mathbf{u}^2 .

This completes the proof. \square

In view of this result, and especially its proof, we realize that if \mathbf{u}^2 is not constrained to be regular, the existence of the three squares imposes severe conditions on

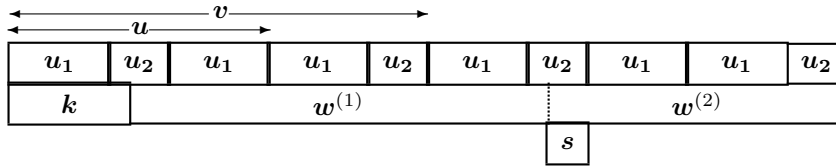


FIG. 14. Case B.IV.

the periodicity of u , as shown in Table 1. In this table we specify, for each of the 14 cases identified in the proof, the prefix of u (u_1 , u_1u_2 , or u itself) that begins with a square, as well as the period of the square. We also indicate cases in which the entire prefix is in fact a repetition. Of course all copies of the prefix in x will have the same periodicity properties. Furthermore, in all cases (3–6 and 10–14) in which a period of u is identified, the periodicity lemma applies, since u also has period $u_1 + u_2 = v - u$. For example, in cases 6 and 12–14 u , hence u^2 , hence all of v^2 , will have period $\gcd(v - w, v - u)$; a similar result holds for cases 4–5 and 11. An alternative form of Lemma 12 may then be given as follows.

LEMMA 13. *Let $u = u_1u_2u_1$, $v = uu_1u_2$, u_1 and u_2 nonempty. If $x = v^2 = ky$, where $k \in 0..v - u - 1$ and y has period $w \in v - u + 1..v - 1$, $w \neq u$, then every occurrence of u in x is determined by cases 1–6 of Table 1.*

Observe that this result holds for every nonempty border u_1 of u such that $u_1 < u/2$.

The rightmost column of Table 1 specifies the length of x that may be required in order to establish the periodicity of the prefix of u . For example, in cases 1–3 not even all of u^2 is required, and even in case 12 not all of v^2 is required. This observation leads to the following weaker, but perhaps still interesting, corollary of the NPL that relates only to u^2 .

LEMMA 14. *Let $u = u_1u_2u_1$, u_1 and u_2 nonempty. If $x = u^2u_1u_2 = ky$, where $k \in 0..u_1$ and y has period $w \in u_1 + u_2 + 1..u_1 + u_2 + u$, $w \neq u$, then every occurrence of u in x is determined by Table 1.*

Again this result holds for every nonempty border u_1 of x .

We can state an equivalent of Lemma 12 for runs. Observe first that by definition every run is irreducible. Observe also that if a run of period u and tail t occurs at position i in x , no run of the same period can occur at any position $j \in i..i + u + t$. Thus, if we define a *regular run* to be a run of generator u where u^2 is a regular square, we can state the following lemma.

LEMMA 15. *Suppose x has a regular run of period u as prefix and another run of period $v < 2u$ as prefix. Then for every integer $k \in 0..v - u - 1$ and for every $w \in u..v$, no run of period w (other than, for $k = 0$, the two given runs) occurs at position $k + 1$ of x .*

Finally, we remark that Lemmas 12 and 15 apply only trivially to the cases $u = 1$ and $u = 2$. As noted earlier for $u = 1$, $v \geq 3 > 2u$, while for $u = 2$, $v \geq 5 > 2u$, contrary to the requirement of the lemmas that $v < 2u$. However, for all $u \geq 3$, the hypothesis of the lemmas can be satisfied—for example, if $u = aba$ of length 3, v may be $abaab$ of length $5 < 2 \times 3$. More generally, we may think of such squares v^2 as being “small,” in contrast to those of period greater than $2u$ that are “large”; thus Lemmas 12 and 15 restrict the occurrences of squares/runs when the second square at some position is small. Note also that if u^2 is in fact minimal (hence by Lemma 6 regular), then the irreducible square v^2 must be regular.

TABLE 1
Periodicity table for k, u, v, w (unconstrained u).

| Case | $k \leq$ | $k+w \leq$ | Subcase | Subsubcase | Square prefix | Repetition? | Period | Required x |
|------|-------------|------------|-------------------|--------------------|----------------|-------------|---------|--------------|
| 1 | u_1 | u | $k+2w \leq u+u_1$ | $k \geq u_2$ | u_1 | | $w-u_1$ | $u+u_1$ |
| 2 | | | | $k < u_2$ | u_1 | | $w-v+u$ | $u+u_1$ |
| 3 | | | | | u | | $u-w$ | v |
| 4 | | $u+u_1$ | $k+2w > u+u_1$ | | u_1u_2 & u | yes | $ u-w $ | $2u$ |
| 5 | | v | | | u_1u_2 & u | yes | $w-u$ | $2v-u$ |
| 6 | | $2u$ | | | u_1u_2 & u | yes | $v-w$ | $v+u$ |
| 7 | u_1+u_2-1 | $u+u_1$ | $k+2w \leq 2u$ | | u_1 | yes | $w-v+u$ | $2u$ |
| 8 | | | $k+2w > 2u$ | | u_1 | yes | $u-w$ | $2v-u$ |
| 9 | | v | $w < u$ | | u_1 | yes | $u-w$ | $2v-u$ |
| 10 | | | $w > u$ | | u | | $w-v+u$ | $v+u$ |
| 11 | | | | $k+2w \leq v+u$ | u | yes | $w-v+u$ | $2v-u_2$ |
| 12 | | | | $k+2w \leq 2v-u_2$ | u | yes | $w-v+u$ | $2v-1$ |
| 13 | | $2u$ | | $k+2w < 2v$ | u | yes | $v-w$ | $v+u$ |
| 14 | | $2v+u_2$ | | | u_1 & u | yes | $v-w$ | $v+u$ |
| | | | | | u_1 & u | yes | $v-w$ | $v+u$ |

3. Discussion. We have proved two main lemmas (Lemmas 7 and 12) that restrict the periods w of squares that can occur at positions $i+k$ in \mathbf{x} when at position i either one (Lemma 7) or two (Lemma 12) squares are known to occur. It seems that, with the exception of [15, Lemma 8.1.14], such properties have not been studied previously. In particular, we hope that with the help of Lemma 12, it will be possible to establish, or at least make progress with, the three conjectures arising out of [13].

The Main/Kolpakov–Kucherov algorithm [16, 13] is the only known linear-time algorithm for computing all the runs in a given string \mathbf{x} . It is complex and, until recently, depended for its worst-case linear behavior on the use of Farach’s algorithm [4], also complex and not space-efficient, for linear-time computation of suffix trees. Since 2003 three worst-case linear-time suffix array construction algorithms [10, 11, 12] have been available for use in the computation of the LZ factorization, but even after the substitution of suffix arrays for suffix trees in the all-runs algorithm, significant complications remain. For instance, it seems clear [19, 20] that due to their recursive nature the linear-time algorithms are not in practice the fastest suffix array construction algorithms available. We hope that, with a more precise understanding of the periodicity of runs, it will become possible to design simpler algorithms that will compute all the runs in a string in a more direct and efficient manner.

Acknowledgment. The authors thank a referee for suggestions that have materially improved the presentation.

REFERENCES

- [1] A. APOSTOLICO AND F. P. PREPARATA, *Optimal off-line detection of repetitions in a string*, Theoret. Comput. Sci., 22 (1983), pp. 297–315.
- [2] M. CROCHEMORE, *An optimal algorithm for computing the repetitions in a word*, Inform. Process. Lett., 12 (1981), pp. 244–250.
- [3] M. CROCHEMORE AND W. RYTTER, *Squares, cubes, and time-space efficient strings searching*, Algorithmica, 13 (1995), pp. 405–425.
- [4] M. FARACH, *Optimal suffix tree construction with large alphabets*, in Proceedings of the 38th Annual IEEE Symposium on Foundation of Computer Science, 1997, pp. 137–143.
- [5] N. J. FINE AND H. S. WILF, *Uniqueness theorems for periodic functions*, Proc. Amer. Math. Soc., 16 (1965), pp. 109–114.
- [6] A. S. FRAENKEL AND R. J. SIMPSON, *How many squares can a string contain?*, J. Combin. Theory Ser. A, 82 (1998), pp. 112–120.
- [7] F. FRANEK, R. J. SIMPSON, AND W. F. SMYTH, *The maximum number of runs in a string*, M. Miller and K. Park, eds., in Proceedings of the 14th Annual Australasian Workshop on Combinatorial Algorithms, 2003, pp. 26–35.
- [8] L. ILIE, *A simple proof that a word of length n has at most $2n$ distinct squares*, J. Combin. Theory Ser. A, 112 (2005), pp. 163–164.
- [9] L. ILIE, *A note on the number of distinct squares in a word*, in Proceedings of the 5th Annual International Conference on Combinatorics on Words, S. Brlek and C. Reutenauer, eds., 2005, pp. 289–294.
- [10] J. KÄRKKÄINEN AND P. SANDERS, *Simple linear work suffix array construction*, in Proceedings of the 30th Annual International Conference on Automata, Languages, and Programming, 2003, pp. 943–955.
- [11] D. K. KIM, J. S. SIM, H. PARK, AND K. PARK, *Linear-time construction of suffix arrays*, in Proceedings of the 14th Annual Symposium on Combinatorial Pattern Matching, Lecture Notes in Comput. Sci. 2676, R. Baeza-Yates, E. Chávez, and M. Crochemore, eds., Springer-Verlag, Berlin, 2003, pp. 186–199.
- [12] P. KO AND S. ALURU, *Space efficient linear time construction of suffix arrays*, in Proceedings of the 14th Annual Symposium on Combinatorial Pattern Matching, Lecture Notes in Comput. Sci. 2676, R. Baeza-Yates, E. Chávez, and M. Crochemore, eds., Springer-Verlag, Berlin, 2003, pp. 200–210.

- [13] R. KOLPAKOV AND G. KUCHEROV, *On maximal repetitions in words*, J. Discrete Algorithms, 1 (2000), pp. 159–186.
- [14] A. LEMPEL AND J. ZIV, *On the complexity of finite sequences*, IEEE Trans. Inform. Theory, 22 (1976), pp. 75–81.
- [15] M. LOTHAIRE, *Algebraic Combinatorics on Words*, Cambridge University Press, Cambridge, UK, 2002.
- [16] M. G. MAIN, *Detecting leftmost maximal periodicities*, Discrete Appl. Math., 25 (1989), pp. 145–153.
- [17] M. G. MAIN AND R. J. LORENTZ, *An $O(n \log n)$ algorithm for finding all repetitions in a string*, J. Algorithms, 5 (1984), pp. 422–432.
- [18] E. M. MCCREIGHT, *A space-economical suffix tree construction algorithm*, J. ACM, 23 (1976), pp. 262–272.
- [19] S. J. PUGLISI, W. F. SMYTH, AND A. TURPIN, *The performance of linear time suffix sorting algorithms*, in Proceedings of the IEEE Data Compression Conference, J. Storer and M. Cohn, eds., 2005, pp. 358–367.
- [20] S. J. PUGLISI, W. F. SMYTH, AND A. TURPIN, *A taxonomy of suffix array construction algorithms*, ACM Comput. Surv., to appear.
- [21] W. RYTTER, *The number of runs in a string: Improved analysis of the linear upper bound*, in Proceedings of the 23rd Symposium on Theoretical Aspects of Computer Science, B. Durand and W. Thomas, eds., Lecture Notes in Comput. Sci. 2884, Springer-Verlag, Berlin, 2006, pp. 184–195.
- [22] B. SMYTH, *Computing Patterns in Strings*, Addison-Wesley, Reading, MA, 2003.
- [23] A. THUE, *Über unendliche zeichenreihen*, Norske Vid. Selsk. Skr. I. Mat. Nat. Kl. Christiana, 7 (1906), pp. 1–22.
- [24] P. WEINER, *Linear pattern matching algorithms*, in Proceedings of the 14th Annual IEEE Symposium on Switching and Automata Theory, 1973, pp. 1–11.