

Computing Periodicities in Strings — A New Approach^{*}

W. F. Smyth^{1,2}

¹ Algorithms Research Group, Department of Computing & Software
McMaster University, Hamilton, Ontario, Canada L8S 4K1
smyth@mcmaster.ca
www.cas.mcmaster.ca/cas/research/groups.shtml

² Department of Computing, Curtin University, GPO Box U1987
Perth WA 6845, Australia
smyth@computing.edu.au

November 14, 2007

Abstract. The most efficient methods currently available for the computation of repetitions or repeats in a string $x = x[1..n]$ all depend on the prior computation of a suffix tree/array ST_x/SA_x . Although these data structures can be computed in asymptotic $\Theta(n)$ time, nevertheless in practice they involve significant overhead, both in time and space. Since the number of repetitions/repeats in x can be reported in a way that is at most linear in string length, it therefore seems that it should be possible to devise less roundabout means of computing repetitions/repeats that take advantage of their infrequent occurrence. This survey paper provides background for these ideas and explores the possibilities for more efficient computation of periodicities in strings.

1 Introduction

The mathematical study of *strings* (*words*) began with an investigation of their periodicity properties [25], a focus that has been sustained over the intervening 100 years. But in addition, periodicity has turned out to be important from a computational/algorithmic point of view: in bioinformatics (repeating subsequences of DNA or protein), for data compression, in cryptanalysis, and in various other contexts. Thus algorithms to compute or recognize periodicities in strings were among the earliest to be developed when the importance of string processing was first recognized about 35–40 years ago.

^{*} This research was supported in part by grants from the Natural Sciences & Engineering Research Council of Canada.

This survey paper begins in Section 2 with an overview of the most important kinds of periodicity and the algorithms that compute them. Then in Section 3 a start is made on a “theory of periodicities” that may provide clues about their more efficient computation. Finally, Section 4 briefly discusses future directions for research.

2 Periodicities

This section defines the most common and important kinds of periodicity in strings and characterizes the algorithms that compute them. Essentially, we discover that periodicities are in some sense infrequent (at most linear in string length), but that the algorithms that compute them, though asymptotically efficient, are also rather complex, perhaps more complex than they need to be. Terminology and notation generally follow [24].

Suffix Structures

Since existing algorithms almost always use suffix structures to achieve effectiveness and efficiency, we briefly define them here.

A *suffix tree* $ST_{\mathbf{x}}$ of a given string \mathbf{x} is a *compacted trie* on the nonempty suffixes of \mathbf{x} [26, 22], in which every terminal node identifies a suffix and every internal node I identifies the *least common prefix* (LCP) of all the terminal nodes of the subtree rooted at I . Terminal nodes are labelled i , identifying suffix $\mathbf{x}[i..n]$, and internal nodes are labelled $\text{lcp} = |\text{LCP}|$, the length of the LCP of the terminal nodes in that subtree. Thus the terminal nodes of each subtree identify suffixes whose prefixes of length lcp are repeating substrings in \mathbf{x} . Figure 1 shows the suffix tree of $\mathbf{x} = \text{abaababa}$. As is usual with suffix trees, the terminal nodes of ST_{abaababa} taken in left-to-right order give the suffixes in lexicographic order; however, for the calculation of periodicities, this property is not necessary.

The *suffix array* $SA_{\mathbf{x}}$ is just the array of the suffixes i of \mathbf{x} read off from $ST_{\mathbf{x}}$ in a preorder (left-to-right) traversal. Sometimes also used in conjunction with $ST_{\mathbf{x}}$ is an lcp array $\text{lcp}_{\mathbf{x}}$ in which, for every $j \in 2..n$,

$$\text{lcp}_{\mathbf{x}}[j] = \text{lcp}(SA_{\mathbf{x}}[j-1], SA_{\mathbf{x}}[j]).$$

The lcp array can also be read off from $ST_{\mathbf{x}}$ in a preorder traversal. For example:

$$\begin{array}{cccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \mathbf{x} & = & a & b & a & a & b & a & b & a \\ SA_{\mathbf{x}} & = & 8 & 3 & 6 & 1 & 4 & 7 & 2 & 5 \\ lcp_{\mathbf{x}} & = & - & 1 & 1 & 3 & 3 & 0 & 2 & 2 \end{array}$$

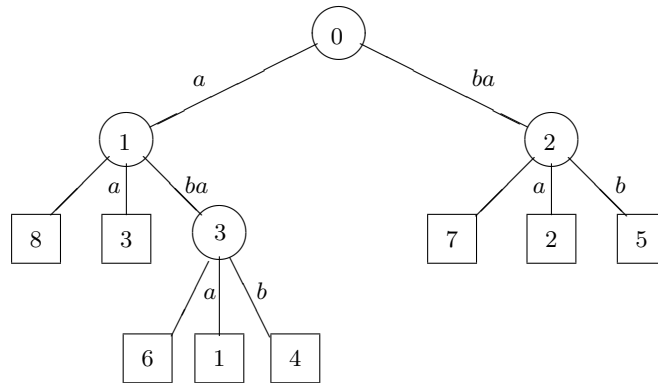


Fig. 1. $ST_{abaababa}$

Repetitions

A **repetition** in a string $x = x[1..n]$ is a substring $x[i..i+pe-1] = u^e$, where $|u| = p$ and $e \geq 2$. If moreover u itself is not a repetition, then u^e is said to be **irreducible**; and if neither $x[i-p..i-1]$ nor $x[i+pe..i+p(e+1)-1]$ equals u , then u^e is said to be **maximal**. Unless specifically stated otherwise, all repetitions referred to in this paper are both irreducible and maximal.

We call u the **generator**, p the **period**, and e the **exponent** of the repetition u^e . Note that a repetition is completely specified by the triple (i, p, e) . In the string

$$\begin{array}{cccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ \mathbf{x} & = & a & b & a & a & b & a & b & a, \end{array}$$

the repetitions are $(1, 3, 2) = (aba)^2$, $(3, 1, 2) = a^2$, $(4, 2, 2) = (ab)^2$, and $(5, 2, 2) = (ba)^2$. Since in each case $e = 2$, all of these are **squares**.

About a quarter-century ago, three algorithms were discovered [4, 2, 21] that employed widely different approaches to compute all the repetitions in a given string $\mathbf{x}[1..n]$ in $O(n \log n)$ time; of these algorithms, two were based on a form of suffix tree calculation ([2] explicitly, [4] implicitly), while the third used a divide-and-conquer technique. In [4] it was moreover shown that the Fibonacci string \mathbf{f}_K :

$$\mathbf{f}_0 = b, \mathbf{f}_1 = a; \quad \mathbf{f}_k = \mathbf{f}_{k-1}\mathbf{f}_{k-2}, \quad k = 2, 3, \dots, K$$

actually contains $\Theta(|\mathbf{f}_K| \log |\mathbf{f}_K|)$ repetitions. Hence these algorithms were regarded as asymptotically optimal, a concept that as we shall see depends heavily on what is accepted as a sufficient specification of a repetition.

Runs

It was mentioned above that the maximum number of repetitions in a string $\mathbf{x} = \mathbf{x}[1..n]$ is $\Theta(n \log n)$. But this is a count of repetitions that are both maximal and irreducible. If instead we were asked to output the distinct squares \mathbf{u}^2 without these restrictions, we would find that $\mathbf{x} = a^n$, for example, would require $\lceil n^2/4 \rceil$ — that is, $\Theta(n^2)$ — outputs to specify squares $\mathbf{x}[1..2], \mathbf{x}[2..3], \dots, \mathbf{x}[n-1..n], \mathbf{x}[1..4], \mathbf{x}[2..5], \dots, \mathbf{x}[n-3..n]$, and so on. Thus in restricting the output to maximal irreducible repetitions, we *encode* the output, by tacit agreement with the user, so as to reduce its quantity, hence the asymptotic complexity of the algorithm. For $\mathbf{x} = a^n$, this encoding dramatically reduces the output to a single repetition $(1, 1, n)$.

We now describe another encoding of repetitions that further reduces the quantity of output required to $\Theta(n)$. We say that a repetition $(i, p, e) = \mathbf{u}^e$ is *left-extendible* (LE) if there exists a repetition at position $i-1$ of \mathbf{x} that is also of period p . If no such repetition exists, we say that (i, p, e) is NLE. Given an NLE repetition (i, p, e) , denote by t the greatest integer such that, for every $j \in 0..t$, $(i+j, p, e)$ is a repetition. Note that since (i, p, e) is maximal, therefore $t \in 0..p-1$. We call t the *tail* of (i, p, e) . Then a *run* (*maximal periodicity*) is a 4-tuple (i, p, e, t) , where (i, p, e) is an NLE repetition of tail t .

The idea of a run was first introduced by Main in [20], where also an algorithm was proposed to compute the leftmost occurrence of every distinct run in $\mathbf{x}[1..n]$. Given the suffix tree $\text{ST}_{\mathbf{x}}$ and the Lempel-Ziv decomposition [18] of \mathbf{x} (computable in linear time from $\text{ST}_{\mathbf{x}}$), Main's algorithm computes all the leftmost runs in $\Theta(n)$ time. In [17] Kolpakov

& Kucherov showed that the maximum number $\rho(n)$ of runs in any string of length n satisfies

$$\rho(n) < k_1 n - k_2 \sqrt{n} \log_2 n \quad (1)$$

for some pair of universal positive constants k_1 and k_2 . They also extended Main's algorithm to compute *all* the runs in \mathbf{x} in time proportional to their number; thus by (1), given $\text{ST}_{\mathbf{x}}$, all the runs in \mathbf{x} , and so in effect all the repetitions, could be computed in $\Theta(n)$ time.

Repeats

A *repeat* in \mathbf{x} is a tuple

$$M_{\mathbf{x},\mathbf{u}} = (p; i_1, i_2, \dots, i_e),$$

where $e \geq 2$, $1 \leq i_1 < i_2 < \dots < i_e \leq n$, and

$$\mathbf{u} = \mathbf{x}[i_1..i_1+p-1] = \mathbf{x}[i_2..i_2+p-1] = \dots = \mathbf{x}[i_e..i_e+p-1].$$

Note that it may happen, for some $j \in 1..e-1$, that $i_{j+1} - i_j = p$ or that $i_{j+1} - i_j < p$ — that is, the substrings of a repeat may form a repetition or even overlap. Analogous to a repetition, we call \mathbf{u} the *generator*, p the *period*, and e the *exponent* of $M_{\mathbf{x},\mathbf{u}}$; also, $M_{\mathbf{x},\mathbf{u}}$ is called a *square* if $e = 2$. We say that $M_{\mathbf{x},\mathbf{u}}$ is *maximal* if for every

$$i \in 1..n \text{ and } i \notin \{i_1, i_2, \dots, i_e\},$$

we are assured that $\mathbf{x}[i..i+p-1] \neq \mathbf{u}$. Again analogous to runs, we say that $M_{\mathbf{x},\mathbf{u}}$ is *left-extendible* (LE) if

$$(p; i_1 - 1, i_2 - 1, \dots, i_e - 1)$$

is a repeat; in this case, $(p+1; i_1-1, i_2-1, \dots, i_e-1)$ is a repeat whose suffixes of length p are specified by $M_{\mathbf{x},\mathbf{u}}$. Similarly, $M_{\mathbf{x},\mathbf{u}}$ is *right-extendible* (RE) if

$$(p; i_1 + 1, i_2 + 1, \dots, i_e + 1)$$

is a repeat; in this case, $(p+1; i_1, i_2, \dots, i_e)$ is a repeat whose prefixes of length p are specified by $M_{\mathbf{x},\mathbf{u}}$. If $M_{\mathbf{x},\mathbf{u}}$ is neither LE nor RE, we say that it is *nonextendible* (NE). Unless explicitly stated otherwise, all repeats discussed in this paper are both maximal and NE.

In $\mathbf{x} = \text{abaababa}$, the repeats are

$$M_{\mathbf{x},a} = (1; 1, 3, 4, 6, 8) \text{ and } M_{\mathbf{x},aba} = (3; 1, 4, 6);$$

since every occurrence of b is both preceded and followed by a , there are no others.

Of particular interest are repeating substrings \mathbf{u} such that $M_{\mathbf{x}, \mathbf{v}_1 \mathbf{u} \mathbf{v}_2}$ is a repeat if and only if $\mathbf{v}_1 = \mathbf{v}_2 = \varepsilon$, the empty string — in other words, \mathbf{u} is not a proper substring of any other repeating substring. We call such repeats *supernonextendible* (SNE). The repeating substring \mathbf{u} in an SNE repeat $M_{\mathbf{x}, \mathbf{u}}$ is in some sense the longest in a class of nested repeating substrings that are substrings of \mathbf{u} . In the above example, $(3; 1, 4, 6)$ is the unique SNE repeat.

In [11, p. 147] an algorithm is described that, given the suffix tree $\text{ST}_{\mathbf{x}}$ of \mathbf{x} , computes all the NE squares (whether maximal or not) in \mathbf{x} in time $O(n+q)$, where q is the number of squares output. [3] uses similar methods to compute all NE squares $(p; i_1, i_2)$ such that $i_2 - i_1 \leq g$ for some user-defined *gap* g . [1] shows how to use the suffix array $\text{SA}_{\mathbf{x}}$ of \mathbf{x} to compute the NE squares and the SNE squares, both in time $O(n+q)$. [10] uses either the suffix trees of both \mathbf{x} and its reversed string $\bar{\mathbf{x}} = \mathbf{x}[n]\mathbf{x}[n-1] \cdots \mathbf{x}[1]$, or alternatively the suffix arrays of both, to compute all the (maximal and NE) repeats in \mathbf{x} in $\Theta(n)$ time. Thus the number of repeats in \mathbf{x} is linear in string length.

There is a noteworthy distinction between the reporting of repetitions/runs on the one hand and repeats on the other. Suppose

$$\mathbf{x} = \cdots \overset{i_1}{abababab} \cdots \overset{i_2}{abababab} \cdots ,$$

where neither a nor b occurs elsewhere in \mathbf{x} . Then two separate repetitions $(i_1, 2, 4) = (i_2, 2, 4) = (ab)^4$ will be reported in \mathbf{x} , whereas a single SNE repeat $(8; i_1, i_2) = abababab$ will be reported. The period of the repetition or run is *reduced* as much as possible, while the period of the SNE repeat is *extended* as much as possible.

Summary

We have seen in this section that the cardinality of the basic periodicities in a string (repetitions/runs and repeats) is in some sense linear in string length; further that, using suffix trees, these periodicities can be computed in linear time. Since every computation using a suffix tree $\text{ST}_{\mathbf{x}}$ can be replaced by an equivalent computation using a suffix array $\text{SA}_{\mathbf{x}}$ with no sacrifice in asymptotic complexity [1], and since $\text{SA}_{\mathbf{x}}$ can be computed in guaranteed $\Theta(n)$ time [14–16], it follows that the total time requirement for the computation of periodicities is $\Theta(n)$. This is a

significant achievement, but, as we discuss in the next section, there is some reason to believe that even “better” (that is, faster and simpler) algorithms can be found.

3 A Theory of Periodicities?

It turns out that the linear-time suffix array construction algorithms (SACAs) are not the fastest in practice! As shown in [23], there are several other SACAs with supralinear worst-case behaviour that, over a wide range of frequently-occurring strings, require less than half of both the time and the space used by their linear-time cousins. Even so, the space requirements in particular are substantial: at least 5 bytes per symbol. This means that for a large string ($n \geq 10^9$, say), most computers would need to use secondary storage for an additional $5n$ bytes, a requirement that would have a disastrous effect on the speed of the algorithm. Thus it becomes of interest to seek ways of avoiding suffix array construction altogether.

Since as we have seen runs and repeats can both be reported in linear time, and are therefore in some sense infrequent in strings, an alternate approach to their calculation is not *a priori* absurd. For repeats no such approach as to my knowledge been studied, but for runs theoretical questions arise whose resolution could plausibly lead to new methods, as we now explain.

The result (1) is flawed in a certain sense: its proof is nonconstructive, and so the magnitude of the constants k_1 and k_2 is completely unknown. At the same time, [17] provides convincing experimental evidence that the following statements are true:

- (a) $\rho(n) < n$;
- (b) $\rho(n+1) - \rho(n) \leq 2$;
- (c) there exists a cube-free string $\mathbf{x}[1..n]$ on $\{a, b\}$ that contains $\rho(n)$ runs.

These conjectures are fundamental ones about periodicity in strings, yet they have scarcely been studied.

Conjecture (a) relates to the well-known conjecture that

$$\sigma(n) < n, \tag{2}$$

where $\sigma(n)$ is the maximum number of *distinct* squares that can occur in any string of length n . Since every run must begin with a square, it follows that $\sigma(n) \leq \rho(n)$, so that a proof of (a) would establish (2). In

fact, it has been shown that $\sigma(n) \leq 2n - 2$ [8, 12], more recently that $\sigma(n) \leq 2n - \Theta(\log n)$ [13].

[9] identifies an infinite family of strings $\mathbf{x}[1..n]$ containing a number $r(n)$ of runs, where

$$\lim_{n \rightarrow \infty} \frac{r(n)}{n} = \frac{3}{2\phi}$$

and $\phi = (1 + \sqrt{5})/2$ is the **golden mean**. It is conjectured that this limit is a maximum over all infinite families of strings. In the same paper the following theorem is proved (a string is said to be **run-maximal** if it contains $\rho(n)$ runs):

Theorem 1 *Let $\mathbf{x} = \mathbf{x}[1..n]$ be a run-maximal string that contains $\alpha \geq 3$ distinct letters. Suppose that one of these letters λ occurs fewer than three times. Then there exists a run-maximal string of length n that contains $\alpha - 1$ distinct letters. \square*

This result, not easy to prove, succeeds only in establishing conjecture (c) in the very special case that all letters except two occur at most twice.

In order to establish conjecture (a), one needs to somehow limit to less than one the average number of runs that begin at the positions of \mathbf{x} . This requirement draws attention to positions i where two or more runs begin: one hopes to be able to show that at positions “neighbouring” to i , no runs can begin. To date, the most famous and important theorem restricting periodicity is the “periodicity lemma” [7]:

Theorem 2 *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} . \square*

Unfortunately this theorem is silent on the case of most interest to us: two squares occurring at the same position. A more relevant result is the “three squares lemma” [5, 19] which we now state using the convention that the length of string \mathbf{x} is denoted by x :

Theorem 3 *Suppose \mathbf{u}^2 is irreducible, and suppose $\mathbf{w} \neq \mathbf{u}^k$ for any $k \geq 1$. If \mathbf{u}^2 is a prefix of \mathbf{w}^2 , in turn a proper prefix of \mathbf{v}^2 , then $w \leq v - u$. \square*

This result tells us that, given two squares \mathbf{u}^2 and \mathbf{v}^2 of periods $u < v$ occurring at some position i in \mathbf{x} , any third square \mathbf{w}^2 whose period $w \in u + 1..v - 1$ must satisfy $w \leq v - u$. Thus the existence of two squares at i imposes restrictions on the period of any third square at i . Since as already observed a run begins with a square, this imposes corresponding restrictions on runs. A recent result, that we call the “new periodicity

lemma” [6, 23] substantially generalizes Theorem 3. Let us say that a square \mathbf{u}^2 is *regular* if no prefix of \mathbf{u} is a square.

Theorem 4 *Suppose that \mathbf{x} has regular prefix \mathbf{u}^2 and irreducible prefix \mathbf{v}^2 , $u < v$.*

- (a) $v > \max\{u + 1, 3u/2\}$;
- (b) if $v < 2u$, then for every $w \in u+1..v-1$ and for every $k \in 0..v-u-1$, $\mathbf{x}[k+1..k+2w]$ is not a square;
- (c) if $v < 2u$, then for every $w \in v-u+1..u-1$ and for every $k \in 2u-v..v-u-1$, $\mathbf{x}[k+1..k+2w]$ is not a square. \square

This result tells us that if there exists “small” $v < 2u$, then strict restrictions apply to the period w of any other square that might occur close to the start (at position $k+1$, in fact) of \mathbf{x} . Theorem 3 actually states the special case $k = 0$ of Theorem 4(b). Figure 2 shows the excluded values of w corresponding to the overlapping zones specified by k in Theorem 4(b)-(c).

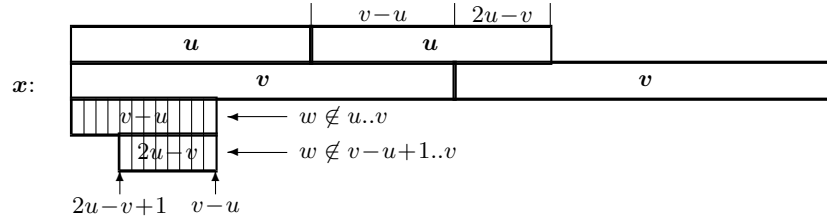


Fig. 2. Exclusion zones for \mathbf{u}^2 and \mathbf{v}^2

It should be remarked that Theorem 4 holds only trivially for $u < 3$, since for $u \leq 2$, we must have $v > 2u$:

- for $u = 1$, $\mathbf{u}^2 = \lambda^2$ for some letter λ , and the minimum $v = 3$ corresponds to $\mathbf{v}^2 = (\lambda^2\mu)^2$, $\mu \neq \lambda$;
- for $u = 2$, $\mathbf{u}^2 = (\lambda\mu)^2$, and the minimum $v = 5$ corresponds to $\mathbf{v}^2 = ((\lambda\mu)^2\nu)^2$.

Thus the occurrence of small regular squares \mathbf{u}^2 can be analyzed as special cases giving rise to “large” v , while for $u \geq 3$, squares \mathbf{v}^2 satisfying $v < 2u$ (“small” v) can always be formed and are covered by Theorem 4.

4 Future Work

In this paper I have described three important kinds of periodicity in strings — repetitions, runs, repeats — and the algorithms currently available for computing them. Generally speaking, these algorithms depend on the time- and space-consuming construction of suffix arrays/trees.

I have put forward conjectures and theoretical results that may facilitate — at least for the calculation of runs — a more direct approach.

This is very much work-in-progress. For example, it is not yet clear whether Theorem 4 is optimal: it is possible that the values of k and w for which $\mathbf{x}[k+1..k+w]^2$ is impossible can be extended. It is also not clear how Theorem 4 can be used to prove conjecture (a), still less how it might be used algorithmically, so as to avoid testing for the existence of certain squares in a certain neighbourhood of a position where two squares are already known to exist.

Proof of conjecture (c) is also an important and largely untouched objective: if (c) holds, then $\rho(n)$ must depend also on alphabet size α , and so the properties of a function $\rho(n, \alpha)$ become an interesting object of research, both from theoretical and algorithmic points of view.

References

1. Mohamed Ibrahim Abouelhoda, Stefan Kurtz & Enno Ohlebusch, **The enhanced suffix array and its applications to genome analysis**, *Proc. Second Workshop on Algorithms in Bioinformatics*, LNCS 2452 (2002) 449–463.
2. Alberto Apostolico & Franco P. Preparata, **Optimal off-line detection of repetitions in a string**, *Theoret. Comput. Sci.* 22 (1983) 297–315.
3. Gerth S. Brodal, Rune B. Lyngso, Christian N. S. Pedersen & Jens Stoye, **Finding maximal pairs with bounded gap**, *J. Discrete Algorithms* 1 (2000) 77–103.
4. Maxime Crochemore, **An optimal algorithm for computing the repetitions in a word**, *Inform. Process. Lett.* 12–5 (1981) 244–250.
5. Maxime Crochemore & Wojciech Rytter, **Squares, cubes and time-space efficient strings searching**, *Algorithmica* 13 (1995) 405–425.
6. Kangmin Fan, R. J. Simpson & W. F. Smyth, **A new periodicity lemma**, *Proc. 16th Annual Symp. Combin. Pattern Matching*, Alberto Apostolico, Maxime Crochemore & Kunsoo Park (eds.), LNCS 3537, Springer-Verlag (2005) 257–265.
7. N. J. Fine & Herbert S. Wilf, **Uniqueness theorems for periodic functions**, *Proc. Amer. Math. Soc.* 16 (1965) 109–114.
8. Aviezri Fraenkel & R. J. Simpson, **How many squares can a string contain?**, *J. Combin. Theory Ser. A* 82 (1998) 112–120.
9. Frantisek Franek, R. J. Simpson & W. F. Smyth, **The maximum number of runs in a string**, *Proc. 14th Australasian Workshop on Combin. Algorithms*, Mirka Miller & Kunsoo Park (eds.) (2003) 26–35.
10. Frantisek Franek, W. F. Smyth & Yudong Tang, **Computing all repeats using suffix arrays**, *JALC* 8–4 (2003) 579–591.

11. Dan Gusfield, *Algorithms on Strings, Trees & Sequences*, Cambridge University Press (1997) 534 pp.
12. Lucian Ilie, **A simple proof that a word of length n has at most $2n$ distinct squares**, *J. Combin. Theory Ser. A* (2005) to appear.
13. Lucian Ilie, **A note on the number of distinct squares in a word** (2005) manuscript.
14. Juha Kärkkäinen & Peter Sanders, **Simple linear work suffix array construction**, *Proc. 30th Internat. Colloq. Automata, Languages & Programming* (2003) 943–955.
15. Dong Kyue Kim, Jeong Seop Sim, Heejin Park & Kunsoo Park, **Linear-time construction of suffix arrays**, *Proc. 14th Annual Symp. Combin. Pattern Matching*, Ricardo Baeza-Yates, Edgar Chávez & Maxime Crochemore (eds.), LNCS 2676, Springer-Verlag (2003) 186–199.
16. Pang Ko & Srinivas Aluru, **Space efficient linear time construction of suffix arrays**, *Proc. 14th Annual Symp. Combin. Pattern Matching*, Ricardo Baeza-Yates, Edgar Chávez & Maxime Crochemore (eds.), LNCS 2676, Springer-Verlag (2003) 200–210.
17. Roman Kolpakov & Gregory Kucherov, **On maximal repetitions in words**, *J. Discrete Algorithms 1* (2000) 159–186.
18. Abraham Lempel & Jacob Ziv, **On the complexity of finite sequences**, *IEEE Trans. Information Theory 22* (1976) 75–81.
19. M. Lothaire, *Algebraic Combinatorics on Words*, Cambridge University Press (2002) 504 pp.
20. Michael G. Main, **Detecting leftmost maximal periodicities**, *Discrete Applied Maths. 25* (1989) 145–153.
21. 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.
22. Edward M. McCreight, **A space-economical suffix tree construction algorithm**, *J. Assoc. Comput. Mach. 23-2* (1976) 262–272.
23. Simon J. Puglisi, W. F. Smyth & Andrew Turpin, **A taxonomy of suffix array construction algorithms**, *Proc. Prague Stringology Conf. '05* (2005) to appear.
24. Bill Smyth, *Computing Patterns in Strings*, Pearson Addison-Wesley (2003) 423 pp.
25. Axel Thue, **Über unendliche zeichenreihen**, *Norske Vid. Selsk. Skr. I. Mat. Nat. Kl. Christiana 7* (1906) 1–22.
26. Peter Weiner, **Linear pattern matching algorithms**, *Proc. 14th Annual IEEE Symp. Switching & Automata Theory* (1973) 1–11.