

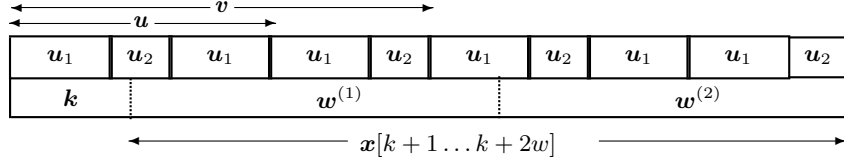
# Overlapping squares in strings (01/2011)

Jenya Kopylova<sup>1</sup> and W. F. Smyth<sup>2</sup>

<sup>1</sup> SEQUOIA, Laboratoire d'Informatique Fondamentale de Lille, France  
evguenia.kopylova@inria.fr

<sup>2</sup> Algorithms Research Group, Department of Computing & Software  
McMaster University, Hamilton, Ontario, Canada L8S 4K1  
smyth@mcmaster.ca

**Question:** What happens to the periodicity of a string when three squares  $(u^2, v^2, w^2)$  begin at neighbouring positions separated from each other by at most  $k$ ?



**Analysis:** In all the cases we know about so far [7, 1, 2], the string breaks down locally into a repetition of small period. This document draws its focus on the 14 subcases defined in Table 1.

**Table 1.** 14 subcases of  $3u/2 < v < 2u, 0 \leq k < v - u$

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

**Solution:** Conjectures for most of the 14 subcases were formulated with the aid of an implemented algorithm outlined in [2]:

**Table 2.** Generated Conjectures

Subcases $S$	Conditions	Breakdown of $\mathbf{x}/v^2$
1,2,5,6,8–10	$(\forall \mathbf{x}, \sigma = d)$	$\mathbf{x} = \mathbf{d}^{(x/d)}$
3,4,7	$\sigma = d$	$\mathbf{x} = \mathbf{d}^{(x/d)}$
	$\sigma > d$	$\mathbf{x} = \mathbf{s}^\alpha \mathbf{s}[1 \dots u_1 \bmod s] \mathbf{s}^\gamma \mathbf{s}[1 \dots u_1 \bmod s] \mathbf{s}^\epsilon$
11–14	$\sigma = d$	$\mathbf{x} = \mathbf{d}^{(x/d)}$
	$\sigma > d$	?

where,

$$d = \gcd(u_1, u_2, w); s = \gcd(u - w, w - u_1)$$

$$\alpha = \lfloor u/s \rfloor; \gamma = \lfloor v/s \rfloor; \epsilon = (u_1 + u_2)/s$$

For subcases  $\{1, 2, 5, 6, 8 - 10\}$ , the proofs have been derived in [2]. For subcases  $\{10, 11 - 14\}$  where  $\sigma = d$ , the proofs are outlined in [1].

**The remaining subcases  $\{3, 4, 7\}$  for  $\sigma \geq d$  and  $\{11 - 14\}$  for  $\sigma > d$  remain as open problems.**

**Why do we want to know?** To be able to use *combinatorial knowledge* about the occurrence of multiple squares at neighbouring positions to:

- \* provide a more precise and also computation-free analysis of the occurrence of runs in a string;
- \* compute repetitions (and perhaps other periodicities) directly rather than using all the heavy machinery of suffix arrays, etc.

*For further background information and examples, please see:*

<http://www.cas.mcmaster.ca/~bill/cv.shtml>  
and  
<http://www.cas.mcmaster.ca/~bj>

*and the accompanying document on “The Maximum Number of Runs in a String” by Bill Smyth.*

# Problem

## The Maximum Number of Runs in a String (2008)

Bill Smyth<sup>1,2</sup>

<sup>1</sup> Algorithms Research Group, Department of Computing & Software  
McMaster University, Hamilton, Ontario, Canada L8S 4K1  
smyth@mcmaster.ca  
www.cas.mcmaster.ca/cas/research/algorithms.htm

<sup>2</sup> Digital Ecosystems & Business Intelligence Institute  
and Department of Computing, Curtin University, GPO Box U1987  
Perth WA 6845, Australia  
smyth@computing.edu.au

Given a nonempty string  $u$  and an integer  $e \geq 2$ , we call  $u^e$  a *repetition*; if  $u$  itself is not a repetition, then  $u^e$  is a *proper repetition*. Given a string  $x$ , a *repetition in*  $x$  is a substring

$$x[i..i+e|u|-1] = u^e,$$

where  $u^e$  is a proper repetition and neither  $x[i+e|u|..i+(e+1)|u|-1]$  nor  $x[i-|u|..i-1]$  equals  $u$ . We say the repetition has *period*  $|u|$  and *exponent*  $e$ ; it can be specified by the integer triple  $(i, |u|, e)$ . It is well known [4] that the maximum number of repetitions in a string  $x = x[1..n]$  is  $\Theta(n \log n)$ , and that the number of repetitions in  $x$  can be computed in  $\Theta(n \log n)$  time [4, 3, 13].

A string  $u$  is a *run* iff it is periodic of (minimum) period  $p \leq |u|/2$ . Thus  $x = abaabaabaab = (aba)^4ab$  is a run of period  $|aba| = 3$ . A substring  $u = x[i..j]$  of  $x$  is a *run* or *maximal periodicity in*  $x$  iff it is a run of period  $p$  and neither  $x[i-1..j]$  nor  $x[i..j+1]$  is a run of period  $p$ . The run  $u$  has *exponent*  $e = \lfloor |u|/p \rfloor$  and possibly empty *tail*  $t = x[i+ep..j]$  (proper prefix of  $x[i..i+p-1]$ ). Thus

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

contains a run  $x[3..12]$  of period  $p = 3$  and exponent  $e = 3$  with tail  $t = a$  of length  $t = |t| = 1$ . It can be specified by a 4-tuple  $(i, p, e, t) = (3, 3, 3, 1)$ . and it includes the repetitions  $(aab)^3$ ,  $(aba)^3$  and  $(baa)^2$  of period  $p = 3$ . In general it is easy to see that for  $e = 2$  a run *encodes*  $t+1$  repetitions; for  $e > 2$ ,  $p$  repetitions. Clearly, computing all the runs in  $x$  specifies all the repetitions in  $x$ . The idea of a run was introduced in [12].

Let  $r_x$  denote the number of runs that actually occur in a given string  $x$ , and let  $\rho(n)$  denote the maximum number of runs that can possibly occur in any string  $x$  of given length  $n$ . A string  $x = x[1..n]$  such that  $r_x = \rho(n)$  is said to be *run-maximal*.

In [10, 11] it was shown that there exist universal positive constants  $k_1$  and  $k_2$  such that

$$\rho(n)/n < k_1 - k_2 \log_2 n / \sqrt{n},$$

but the proof was nonconstructive and provided no way of estimating the magnitude of  $k_1$  and  $k_2$ . In [10], using a brute force algorithm, a table of  $\rho(n)$  was computed for  $n = 5, 6, \dots, 31$ , giving also for each  $n$  an example of a run-maximal string; for every  $n$  in this range,  $\rho(n)/n < 1$  and  $\rho(n) \leq \rho(n-1) + 2$ . In [8] an infinite sequence  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots\}$  of strings was described, with  $|\mathbf{x}_{i+1}| > |\mathbf{x}_i|$  for every  $i \geq 1$ , such that

$$\lim_{i \rightarrow \infty} r\mathbf{x}_i/|\mathbf{x}_i| = \frac{3}{2\phi},$$

where  $\phi = \frac{1+\sqrt{5}}{2}$  is the golden mean. Moreover, it was conjectured that in fact

$$\lim_{n \rightarrow \infty} \rho(n)/n = \frac{3}{2\phi}. \quad (1)$$

Recently a different and simpler construction was found [9] to yield another infinite sequence  $X$  of strings for which the ratio  $r\mathbf{x}_i/|\mathbf{x}_i|$  approached the same limit; in addition, it was shown that for every  $\epsilon > 0$  and for every sufficiently large  $n = n(\epsilon)$ ,  $\frac{3}{2\phi} - \epsilon$  provides an asymptotic lower bound on  $\rho(n)/n$ .

In 2006 considerable progress was made on the estimation of an upper bound on  $\rho(n)/n$ :

- \*  $\rho(n)/n \leq 5.0$  [15];
- \*  $\rho(n)/n \leq 3.48$  [14];
- \*  $\rho(n)/n \leq 3.44$  [16]<sup>1</sup>;
- \*  $\rho(n)/n \leq 1.6$  [5].

Thus the problems may be stated as follows:

**Is conjecture (1) true?**  
**In any case, characterize the function  $\rho(n)/n$ .**

Help may be found in recent work studying the limitations imposed on the existence and length of runs in neighbourhoods of positions where two runs are known to exist [7, 17].

## Additional results in 2008

In [6] new perspectives on the problem are discussed. Based on further computational work, Lucian Ilie's website

<http://www.csd.uwo.ca/faculty/ilie/>

claims  $\rho(n)/n \leq 1.048n$ .

## References

1. R. J. SIMPSON, **Intersecting periodic words**, *Theoret. Comput. Sci.* 374 (2007) 58–65.

---

<sup>1</sup> Unfortunately, this bound has recently been found to be invalid, due to an error in a proof (2008)

2. EVGUENIA KOPYLOV & W. F. SMYTH, **The three squares lemma revisited**, *J. Discrete Algorithms* (submitted for publication).
3. Alberto Apostolico & Franco P. Preparata, **Optimal off-line detection of repetitions in a string**, *Theoret. Comput. Sci.* 22 (1983) 297–315.
4. Maxime Crochemore, **An optimal algorithm for computing the repetitions in a word**, *Inform. Process. Lett.* 12–5 (1981) 244–250.
5. Maxime Crochemore & Lucian Ilie, **Maximal repetitions in strings**, *J. Comput. Sys. Sci.* (2007) to appear.
6. Maxime Crochemore, Lucian Ilie & Liviu Tinta, **Towards a solution to the “runs” conjecture**, *Proc. 19<sup>th</sup> Annual Symp. Combinatorial Pattern Matching* (2008) to appear.
7. Kangmin Fan, Simon J. Puglisi, W. F. Smyth & Andrew Turpin, **A new periodicity lemma**, *SIAM J. Discrete Math.* 20–3 (2006) 656–668.
8. Frantisek Franek, R. J. Simpson & W. F. Smyth, **The maximum number of runs in a string**, *Proc. 14<sup>th</sup> Australasian Workshop on Combinatorial Algs.*, M. Miller & K. Park (eds.) (2003) 26–35.
9. Frantisek Franek & Qian Yang, **An asymptotic lower bound for the maximum-number-of-runs function**, *Proc. Prague Stringology Conference '06*, Jan Holub & Jan Žd'árek (eds.) (2006) 3–8.
10. Roman Kolpakov & Gregory Kucherov, *Maximal Repetitions in Words or How to Find all Squares in Linear Time*, Rapport LORIA 98-R-227, Laboratoire Lorrain de Recherche en Informatique et ses Applications (1998) 22 pp.
11. Roman Kolpakov & Gregory Kucherov, **On maximal repetitions in words**, *J. Discrete Algorithms* 1 (2000) 159–186.
12. Michael G. Main, **Detecting leftmost maximal periodicities**, *Discrete Applied Maths.* 25 (1989) 145–153.
13. Michael G. Main & Richard J. Lorentz, **An  $O(n \log n)$  algorithm for finding all repetitions in a string**, *J. Algs.* 5 (1984) 422–432.
14. Simon J. Puglisi, R. J. Simpson & W. F. Smyth, **How many runs can a string contain?**, *Theoret. Comput. Sci.* (2008) to appear.
15. 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.
16. Wojciech Rytter, **The number of runs in a string**, *Information & Computation* 205–9 (2007) 1459–1469.
17. R. J. Simpson, **Intersecting periodic words**, *Theoret. Comput. Sci.* (2007) 58–65.