

Linear searching for a square in a word.

Bul. Euro. Assoc. Theor. Comput. Sci. 24:66-72, 1984

Max CROCHEMORE

Laboratoire d'Informatique

Université de Haute-Normandie

BP 67

76130 MONT-SAINT-AIGNAN

FRANCE

## Abstract

Searching for a square in a word may be implemented in time proportional to the length of the word on a random access machine provided the alphabet is fixed.

Several methods have been developed to find a square in a word  $x$  [AP83, Cr81, ML83]. They share in common an  $O(|x| \cdot \log|x|)$  time complexity. In fact the three algorithms are able to produce all the repetitions inside the given word with the same complexity. On an infinite alphabet or if the alphabet is not known this complexity cannot be improved even when only one square is searched. We show here that the situation changes when the alphabet is fixed. In this case the search for a square can be implemented in linear time. Our algorithm makes use of a result on concatenation of squarefree words and the fact that suffixes of a word, on a fixed alphabet, may be lexicographically ordered in linear time [Mc76, We73].

## 1 Definitions

The words that are considered are elements of the free monoid  $A^*$  generated by some finite alphabet  $A$ . The empty word is denoted by  $1$ ,  $A^+$  is  $A^* - \{1\}$  and  $|x|$  means the length of a word  $x$  in  $A^*$ .

A word  $x$  of  $A^*$  is said to contain a square  $yy$  if  $y \in A^+$  and  $x$  may be written  $x_1yyx_2$  for some  $x_1$  and  $x_2$  in  $A^*$ . Otherwise  $x$  is said to be square-free. For example, words of length  $\leq 1$  are square-free.

When  $x = yz$ ,  $y$  is called a prefix of  $x$  and  $z$  a suffix of  $x$ . To each such decomposition  $(y, z)$  of  $x$  is associated its separator denoted by  $\text{sep}(y, z)$  : it is the greatest couple of words  $(v, w)$  which satisfies :  $|v| < |y|$  if  $w \neq \lambda$ ,  $vw$  is a prefix of  $yz$  and  $w$  is a prefix of  $z$ . Couples of words are ordered by the relation :

$$(v', w') \leq (v, w) \text{ iff } |w'| < |w| \text{ or } (|w'| = |w| \text{ and } |v'| \geq |v|).$$

When  $(v, w)$  is the separator of  $(y, z)$ ,  $w$  is then the longest prefix of  $z$  which appears at least twice in the word  $yz$ . The word  $v$  marks the first occurrence of  $w$  in  $yz$ .

Example 1. If  $a$  is a letter which does not occur in  $y$  then  $\text{sep}(y, az) = (l, l)$ . For any word  $x \in A^*$ ,  $\text{sep}(l, x) = \text{sep}(x, l) = (l, l)$ . ■

Example 2. On  $A = \{a, b, c, d\}$ , for the word

$$x_0 = \text{abcdbabdcadbcbdbabcabc}, \text{ we have}$$

$$\text{sep}(\text{abcdbabdcadbcbdbabc}, \text{abcb}) = (l, abc). ■$$

Based on the separators of the decompositions of a word  $x \in A^*$  in prefix and suffix is defined the s-factorization of  $x$  which is a sequence  $(u_1, u_2, \dots, u_k)$  of non-empty words such that

$$x = u_1 u_2 \dots u_k$$

and which is built iteratively as follows : if  $u_1, \dots, u_{i-1}$  are already built and  $x = u_1 \dots u_{i-1} az$  for  $a \in A$  and  $z \in A^*$  and if  $\text{sep}(u_1 \dots u_{i-1}, az) = (v, w)$ , then  $u_i$  is  $a$  or  $w$  whether  $w$  is empty or not.

The s-factorization of a word is unique and the empty sequence is associated to the empty word. We may note that  $u_1$  the first term of the s-factorization of a non-empty word is its first letter.

Example 2 (continued). The s-factorization of  $x_0$  is  
 $(a, b, c, d, b, ab, d, c, a, db, c, bd, bab, ca, bcb)$ . ■

We finally introduce two properties on couples of square-free words which deals with squares in their concatenation.

When  $u, u' \in A^*$ , we say that  $(u, u')$  satisfies the property right, denoted by  $r(u, u')$  if  $u$  and  $u'$  are square-free and

$$\exists u_1, u_2, u_3, u_4 \in A^* \quad u_2 u_3 \neq 1, \quad u = u_1 u_2 \text{ and } u' = u_3 u_2 u_3 u_4.$$

The property left, denoted by  $l(u, u')$ , is true for  $(u, u')$  exactly when  $r(\bar{u}', \bar{u})$  is true ( $\bar{u}$  and  $\bar{u}'$  are the minor images of  $u$  and  $u'$  respectively).

Then,  $r(u, u') = \text{true}$  means that  $u$  and  $u'$  are squarefree but their concatenation  $uu'$  contains a square centered at the right of  $u$ .

Example 2 (continued)

$$r(\text{abcdbabdcadbcbdbab}, \text{cabcb}) = \text{true}$$

with  $u_1 = \text{abcdbabdcadbcbdb}$ ,  $u_2 = ab$ ,  $u_3 = c$ ,  $u_4 = b$ . The word  $x_0$  contains the square  $\text{abcabc}$ . ■

## 2 Main Theorem

Separators alone are not sufficient to characterize words which contain a square ; nevertheless they provide a sufficient condition which is given in the following lemma. With the help of s-factorizations we are able to give an 'if and only if' condition for words that contain a square. The result in the next theorem is the base of our linear test for squarefreeness of words.

Lemma. Let  $a_1, \dots, a_n$  be letters of  $A$ . If the word  $a_1 \dots a_n$  satisfies  
 $\exists i \in (1, 2, \dots, n) \quad \exists v, w \in A^*$   
 $(1) \quad (v, w) = \text{sep}(a_1 \dots a_{i-1}, a_i \dots a_n), \quad w \neq 1 \text{ and } a_1 \dots a_{i-1} \text{ prefix of } vw$   
then it contains a square.

Proof :  $a_1 \dots a_n$  contains two occurrences of  $w$  which overlap or are consecutive and then  $a_1 \dots a_n$  is not square-free. ■

Theorem. Let  $a_1, \dots, a_n$  be letters of  $A$  and  $(u_1, \dots, u_k)$  be the s-factorization of  $a_1 \dots a_n$ . The word  $a_1 \dots a_n$  contains a square if and only if it satisfies (1) or

$$(2) \quad \exists j \in (1, \dots, k-1) \quad l(u_i, u_{i+1}) \text{ or } r(u_i, u_{i+1}) \text{ or } r(u_1 \dots u_{i-1}, u_i u_{i+1})$$

Proof : We have only to prove that if  $a_1 \dots a_n$  contains a square and does not satisfy (1), it satisfies (2).

Let  $x_1yy$  be the shortest prefix of  $a_1 \dots a_n$  which contains a square ( $x_1 \in A^*$ ,  $y \in A^+$ ). Let  $i$  be the smallest integer such that  $x_1yy$  is prefix of  $u_1u_2 \dots u_iu_{i+1}$ . Note that  $i \geq 1$ , since  $|u_1| = 1$  and then  $u_1$  is square-free.

The definition of  $i$  and the fact that  $u_i$ 's are non empty words ensure that  $u_1 \dots u_i$  is square-free and is a proper prefix of  $x_1yy$ .

From this, we deduce that  $u_1 \dots u_{i-1}$  (which is empty if  $i=1$ ) is a proper prefix of  $x_1y$  : the contrary would contradict the definition of  $u_i$  and particularly the maximality of its length.

It then follows that the second occurrence of  $y$  is a factor (neither prefix nor suffix) of  $u_i u_{i+1}$ . We conclude, noting that  $l(u_i, u_{i+1}) = r(u_i, u_{i+1}) = \text{false}$  implies  $u_i u_{i+1}$  is squarefree and then  $d(u_1 \dots u_{i-1}, u_i u_{i+1})$  is true.  $\square$

Example 2 (continued) Using the methods of the function `SQUARE`.

The word  $x_0$  does not satisfy the condition (1) of theorem but (2) holds.  $\square$

### 3 Square-freeness test

Following the conditions in the previous theorem, we get a square-freeness test that we write as a function named `SQUARE` whose value is true exactly when the input word contains a square.

```

function SQUARE ( $a_1 \dots a_n$ ) ;
   $i \leftarrow 1$  ;
  while  $i \leq n$  do
     $(v, w) \leftarrow \text{sep}(a_1 \dots a_{i-1}, a_i \dots a_n)$  ;
    if  $w = 1$  then  $u' \leftarrow a_i$  ;  $j \leftarrow i$  ;  $i \leftarrow i+1$ 
    else  $u \leftarrow u'$  ;  $u' \leftarrow w$  ;
      if  $(|vw| \geq i-1 \text{ or } l(u, u') \text{ or } r(u, u') \text{ or } r(a_1 \dots a_{j-1}, uu'))$ 
      then return (true) ;
      end if ;
       $j \leftarrow i$  ;  $i \leftarrow i+|w|$ 
    end if
  end while

```

```
    return (false)
end function
```

The function SQUARE follows the theorem in 2 and takes into account that fact : if  $a_1 \dots a_{i-1}$  is square-free and  $\text{sep}(a_1 \dots a_{i-1}, a_i \dots a_n) = (l, l)$ , the letter  $a_i$  does not appear in  $a_1 \dots a_{i-1}$  and then  $a_1 \dots a_i$  is also square-free.

#### 4 Time complexity

The time complexity of the function SQUARE depends on the time needed to compute the values of functions  $\text{sep}$ ,  $l$  and  $r$ .

On any finite alphabet, it is possible to find the separators of all the decompositions  $(y, z)$  of a word  $x (=yz)$  in linear time adapting the construction of the suffix tree of  $x$  [Mc 76]. The suffix tree may be precomputed, or built during the execution of the function SQUARE. In the latter case the algorithm becomes 'almost on-line'.

It remains to see that all the evaluations of the functions  $l$  and  $r$  is also linear in the length of the input word. The Morris and Pratt's algorithm [KMP 77] includes the computation of a function called  $f$  which is useful for this purpose.

For any word  $u$  in  $A^+$ ,  $f(u)$  is defined to be the longest word distinct from  $u$  which is both a prefix and a suffix of  $u$ . Furthermore, the computations of the  $f(u)$ 's, for the non empty prefixes  $u$  of a word  $x$ , may be done in time proportionnal to  $|x|$ .

The next proposition gives a base to built algorithms which realize the functions  $l$  and  $r$ .

Proposition. Let  $u, u' \in A^*$  which satisfy  $r(u, u')$ , and let  $\text{psp}$  be the shortest prefix of  $u'$  such that  $s$  is a suffix of  $u$  ( $p \in A^*$ ,  $s \in A^+$ ). Then

i-  $p = f(\text{psp})$

ii- if  $p'vs'p'$  is a prefix of  $u'$  ( $p' \in A^*$ ,  $v, s' \in A^+$ ) such that  $|\text{psp}| < |p'vs'p'|$  and  $s'$  is the longest common suffix of both  $u$  and  $p'vs'$ , then  $|ps| < \max(|p'v|, |p'vs'|/2)$ .

To evaluate  $d(u, u')$  when  $u$  and  $u'$  are square-free we proceed as follows : reading  $u'$  from left to right  $f$  is first computed ; then from right to left for a prefix  $psp$  of  $u'$  such that  $p = f(psp)$  it is checked if  $s$  is a suffix of  $u$ . The total number of comparisons between letters during these tries is bounded by  $2|u'|$  as a consequence of part ii of the proposition which allows to make jumps.

Corollary. When  $u$  and  $u'$  are square-free words in  $A^*$ ,  $r(u, u')$  (resp.  $l(u, u')$ ) may be evaluated in time proportionnal to  $|u'|$  (resp.  $|u|$ ).

Turning back to the total cost of calls to functions  $l$  and  $r$  in the function **SQUARE**, the above corollary gives a complexity proportionnal to  $\sum_{i=1}^{k-1} |u_i u_{i+1}|$  if  $(u_1, \dots, u_k)$  is the s-factorization of the input word  $a_1 \dots a_n$ , and this is majorate by  $2n$ .

Theorem. On any finite alphabet, searching a word for a square may be realized in time linear in the length of the input word.

## 5 Conclusion

The above algorithm gives a positive answer to a question of Main & Lorentz which were the first to publish an  $O(|x| \log |x|)$  algorithm to test whether a word  $x$  is square-free or not. On arbitrary (or infinite) alphabets our algorithm is still useful : but, in that case, the implementation could be different using hash tables or search trees to deal with the distinct letters of the word. The worst case complexity becomes  $O(|x| \log m)$  where  $m$  is the number of distinct letters of  $x$ .

The method used here to find a square in a word cannot be directly applied to cubes or higher power but we conjecture that recognition of these other patterns may also be realized in linear time.

## 6 References

AP 83      A. Apostolico & F.P. Préparata.  
 Optimal off-line detection of repetitions in a string, Theor-  
 Compt. Sci. 22 (1983) 297-315.

Cr81      **M. Crochemore**  
An optimal algorithm for computing the repetitions in a word,  
Information Processing Letters, 12 (1981) 244-250.

KMP77      **D.E. Knuth, J.H. Morris & V.R. Pratt.**  
Fast pattern-matching in strings, SIAM J. Comput. 6 (1977)  
323-350.

ML83      **M. Main & R. Lorentz.**  
An  $O(n \log n)$  algorithm for finding all repetitions in a string,  
J. Of algorithms (1983) to appear.

Mc76      **E.M. Mc Creight.**  
A space-economical suffix tree construction algorithm, J. of  
the ACM, 23, 2(1976) 262-272.

We73      **P. Weiner.**  
Linear pattern matching algorithms, in : (proceedings of the  
14th annual symposium on switching and automata theory, 1973)  
1-11.

---

```
@ARTICLE (MCrochemore84eatcs,
  author = {Maxime Crochemore},
  title = {Linear searching for a square in a word},
  journal = {Bul. Euro. Assoc. Theor. Comput. Sci.},
  volume = {24},
  year = {1984},
  pages = {66--72},
  note = {Presented at ICALP'84.  
Abstract in \textit{Automata, Languages and Programming},  
LNCS 172, pp. 137. Springer, Berlin, 1984},
)
```

See also:

```
@ARTICLE (Cro83cras,
  author = {Maxime Crochemore},
  title = {Recherche lin\'aire d'un carr\'e dans un mot},
  journal = {C. R. Acad. Sc. Paris S\'er. I Math.},
  volume = {296},
  number = {18},
  year = {1983},
  pages = {781--784},
```

---