

## Computing the Longest Previous Factor

Maxime Crochemore, Lucian Ilie, Costas Iliopoulos, Marcin Kubica, Wojciech Rytter, Tomasz Waleń

► **To cite this version:**

Maxime Crochemore, Lucian Ilie, Costas Iliopoulos, Marcin Kubica, Wojciech Rytter, et al.. Computing the Longest Previous Factor. *European Journal of Combinatorics*, Elsevier, 2013, 34 (1), pp.15 - 26. <10.1016/j.ejc.2012.07.011>. <hal-01616480>

**HAL Id: hal-01616480**

**<https://hal-upec-upem.archives-ouvertes.fr/hal-01616480>**

Submitted on 13 Oct 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Contents lists available at SciVerse ScienceDirect

European Journal of Combinatorics

journal homepage: [www.elsevier.com/locate/ejc](http://www.elsevier.com/locate/ejc)



## Computing the Longest Previous Factor<sup>☆,☆☆</sup>

Maxime Crochemore<sup>a,d</sup>, Lucian Ilie<sup>b</sup>, Costas S. Iliopoulos<sup>a,e</sup>, Marcin Kubica<sup>c</sup>,  
Wojciech Rytter<sup>c,f</sup>, Tomasz Waleń<sup>c</sup>

<sup>a</sup> Department of Informatics, King's College London, London WC2R 2LS, UK

<sup>b</sup> Department of Computer Science, University of Western Ontario, London, Ontario, N6A 5B7, Canada

<sup>c</sup> Institute of Informatics, Warsaw University, ul. Banacha 2, 02-097 Warszawa, Poland

<sup>d</sup> Université Paris-Est, France

<sup>e</sup> Digital Ecosystems & Business Intelligence Institute, Curtin University of Technology, Perth WA 6845, Australia

<sup>f</sup> Department of Math. and Informatics, Copernicus University, Torun, Poland

### ARTICLE INFO

Article history:

Available online 24 August 2012

### ABSTRACT

The *Longest Previous Factor* array gives, for each position  $i$  in a string  $y$ , the length of the longest factor (substring) of  $y$  that occurs both at  $i$  and to the left of  $i$  in  $y$ . The Longest Previous Factor array is central in many text compression techniques as well as in the most efficient algorithms for detecting motifs and repetitions occurring in a text. Computing the Longest Previous Factor array requires usually the Suffix Array and the Longest Common Prefix array. We give the first time–space optimal algorithm that computes the Longest Previous Factor array, given the Suffix Array and the Longest Common Prefix array. We also give the first linear-time algorithm that computes the permutation that applied to the Longest Common Prefix array produces the Longest Previous Factor array.

© 2012 Elsevier Ltd. All rights reserved.

### 1. Introduction

The *Longest Previous Factor* array gives, for each position  $i$  in a string  $y$ , the length of the longest factor (substring) of  $y$  that occurs both at  $i$  and to the left of  $i$  in  $y$ . For example, for the string  $y = \text{abbaba}$ , the Longest Previous Factor corresponding to position 3 is 2 because  $ab$  is the longest factor at position 3 that appears before (at position 0). The Longest Previous Factor array is central in many text compression techniques as well as in the most efficient algorithms for detecting motifs and

<sup>☆</sup> A preliminary version of this paper has been presented at IWOCA'09; see [5].

<sup>☆☆</sup> Research supported in part by the Royal Society, UK.

E-mail address: [ilie@csd.uwo.ca](mailto:ilie@csd.uwo.ca) (L. Ilie).

repetitions occurring in a text. The problem of computing the Longest Previous Factor array may be regarded as an extension of the Ziv–Lempel factorization (LZ77) of a string as defined in [22].

The LZ77 factorization is used in several adaptive compression methods which carefully encode re-occurrences of phrases by pointers or integers (see [20,21]). It yields more powerful compressors than the factorization in [23] (called LZ78) where a phrase is an extension by a single letter of a previous phrase. But the LZ77 factorization is more difficult to compute. One remarkable application of the Longest Previous Factor array is that the LZ77 factorization can be easily computed from it (see [4]). This implies a linear-time solution for LZ77 factorization on integer alphabets. Previous solutions, using a Suffix Tree [19] or a Suffix Automaton [2] of the string, not only run in time  $O(n \log |A|)$  ( $n$  is the length of the string and  $|A|$  is the size of the alphabet) but also the data structures they use are more space-expensive than the Suffix Array which is used for computing the Longest Previous Factor array.

A slight variant of string parsing, whose relation with the LZ77 factorization is analyzed in [1], plays an important role in String Algorithms. The intuitive reason is that, when processing a string on-line, the work done on an element of the factorization can usually be skipped because it was done already on one of its previous occurrences. A typical application of this idea resides in algorithms to compute repetitions in strings (see [2,15,14]). For example, the algorithm in [14] reports all maximal repetitions (called runs) occurring in a string of length  $n$  in  $O(n \log |A|)$  time. It runs in linear time if a Suffix Array is used instead of a Suffix Tree [4]. Indeed, this seems to be the only technique that leads to linear-time algorithms, independently of the alphabet size, for this type of question.

The Suffix Arrays provide an ideal data structure to solve many questions requiring an index on all the factors of a string. Introduced by Manber and Myers [16], the structure can be built in linear-time by different methods [10,12,13,18] for sorting the suffixes of the string, possibly adding the method of [11] to compute the Longest Common Prefix array. The result holds if the string is drawn from an integer alphabet, that is, if the alphabet of the string can be sorted in linear time (otherwise the  $\Omega(n \log n)$  lower bound for sorting applies).

The notion of Longest Previous Factor has been introduced in [4], however an array with the same definition already appeared in McCreight's suffix tree construction algorithm [17] (the *head* array) and recently in [8] (the  $\pi$  array). In [8], the authors presented a direct computation running in  $O(n \log n)$  average time. A naive computation of the Longest Previous Factor array, either on the string itself or on its Suffix Array, as done by the algorithm LPF-NAIVE in Section 5, leads to quadratic effective running time on many inputs.

The first linear-time algorithm for computing the Longest Previous Factor array appears in [4] where the application to computing the LZ77 factorization (in linear time) is given as well. An improved version, running on-line on the Suffix Array of the string and requiring only  $O(\sqrt{n})$  extra memory space, is shown in [6].

In this article we intensively use the Suffix Array of the string to be processed, and we consider only linear-time solutions. A graphical representation of the Suffix Array structure helps understand the design of the algorithms.

We improve on the previous results and show that the computation of the Longest Previous Factor array of a string from its Suffix Array can be implemented to run in linear time with only a constant amount of additional memory space. Thus, the method is time–space optimal.

The next section introduces the necessary material for the design of Longest Previous Factor computations. An algorithm similar to the one of [4] is described in Section 3. Section 4 shows how the computation can be done on-line on the Suffix Array and Section 5 describes the time–space optimal algorithm for constructing the Longest Previous Factor array. Finally, the Longest Previous Factor array can be obtained as a permutation of the Longest Common Prefix array and we give in Section 6 a linear-time algorithm for that. Alternatively, this can be viewed as a sorting network transforming the Longest Common Prefix array into the Longest Previous Factor array. We conclude by discussing a couple of important research directions.

## 2. Basic definitions

We consider a string  $y$  of length  $n$  over an alphabet  $A$ . The  $i$ th letter of  $y$  is denoted by  $y[i]$  and  $y[i..j] = y[i]y[i+1] \cdots y[j]$ . If, for a string  $x$  and two positions  $i$  and  $j$  in  $y$  we have  $x = y[i..j]$ , then we

**Table 1**

The LPF, prev, and next arrays for the string  $y = \text{abaabababbabb}$ .

Position $i$	0	1	2	3	4	5	6	7	8	9	10	11	12	13
$y[i]$	a	b	a	a	b	a	b	a	b	b	a	b	b	b
LPF[ $i$ ]	0	0	1	3	2	4	3	2	1	4	3	2	2	1
prev[ $i$ ]	-1	-1	1	2	3	4	5	1	7	8	9	10	9	12
next[ $i$ ]	1	14	7	7	7	7	7	14	14	14	12	12	14	14

say that  $x$  occurs at position  $i$  in  $y$ . We index  $y$ 's letters from 0 to  $n - 1$ , that is,  $y = y[0..n - 1]$ . The suffix of  $y$  that starts at position  $i$  is  $\text{suf}_i(y) = y[i..n - 1]$ . The length of  $y$  is denoted by  $|y| = n$ .

In the LZ77 factorization, a string  $y$  is decomposed into factors, called phrases,  $u_0, u_1, \dots, u_k$ , such that

- (i)  $y = u_0u_1 \dots u_k$ ,
- (ii)  $u_0 = y[0]$  and,
- (iii) for  $i \geq 1$ ,  $u_i$  is either a new letter or, otherwise, the longest prefix of  $u_iu_{i+1} \dots u_k$  occurring both at position  $|u_0u_1 \dots u_{i-1}|$  and before it (that is, to its left) in  $y$ ; see [Example 1](#).

The Longest Previous Factor array is defined by  $\text{LPF}[0] = 0$  and, for  $1 \leq i \leq n - 1$ , by

$$\text{LPF}[i] = \max\{k \mid y[i..i + k - 1] = y[j..j + k - 1], \text{ for some } 0 \leq j < i\}.$$

Given the LPF array, the LZ77 factorization is easily computed by the following algorithm (see [4]):

```

LZ77(y, LPF)
1   $u_0 \leftarrow y[0]$ 
2   $k \leftarrow 0$ 
3  while ( $(\ell \leftarrow |u_0u_1 \dots u_k|) < n$ ) do
4       $u_{k+1} \leftarrow y[\ell \dots \ell + \max\{1, \text{LPF}[\ell]\} - 1]$ 
5       $k \leftarrow k + 1$ 
6  return  $u_0, u_1, \dots, u_k$ 
    
```

**Example 1.** For the string  $y = \text{abaabababbabb}$ , the LZ77 factorization is

$$y = \text{a.b.a.aba.bab.babb.b.}$$

The LPF array corresponding to  $y$  is given in the third row of [Table 1](#); the other two arrays, prev and next, will be defined in [Section 3](#).

The Suffix Array of the string  $y$  is a data structure used for indexing its content. It comprises two arrays that we denote SA and LCP (from Longest Common Prefix) and that are defined as follows.

The array SA stores the list of positions of  $y$  associated with the sorted list of its suffixes in increasing lexicographic order. That is, the array is such that

$$\text{suf}_{\text{SA}[0]}(y) < \text{suf}_{\text{SA}[1]}(y) < \dots < \text{suf}_{\text{SA}[n-1]}(y).$$

Thus, SA is indexed by the ranks of the suffixes in their sorted list.

The second array, LCP, is also indexed by the ranks of suffixes and stores the longest common prefixes between consecutive suffixes in the sorted list. Let  $\text{lcp}(i, j)$  be the longest common prefix of  $\text{suf}_i(y)$  and  $\text{suf}_j(y)$ , for two positions  $i$  and  $j$  of  $y$ . Then,  $\text{LCP}[0] = 0$  and, for  $0 < r < n$ , we have

$$\text{LCP}[r] = |\text{lcp}(\text{SA}[r - 1], \text{SA}[r])|.$$

We shall use sometimes the array ISA, the inverse of SA, that provides the rank of a position.

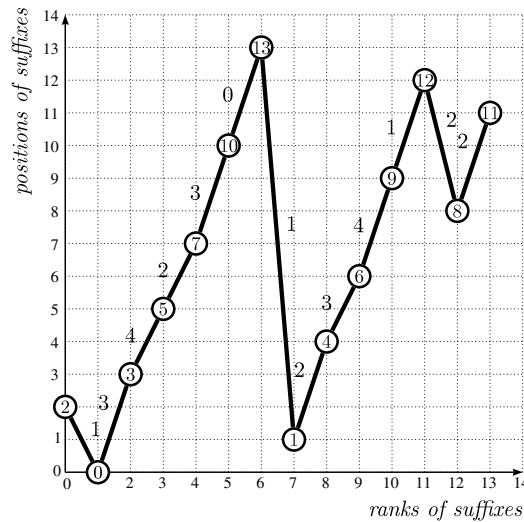
**Example 2.** For our running example  $y = \text{abaabababbabb}$ , the SA and LCP arrays are shown in [Table 2](#).

It is natural to define the LPF array indexed on positions in the string and the SA and LCP arrays indexed on ranks rather than positions. We show in the last column of [Table 2](#) the LPF array indexed also by rank, that is,  $\text{LPF}[\text{SA}[\cdot]]$ . This duality, of indexing by position or by rank, will appear a number of times throughout the paper.

**Table 2**

The SA and LCP arrays for the string  $y = \text{abaabababbabb}$ . The last column gives the LPF array indexed by ranks of suffixes instead of positions.

Rank $r$	SA[ $r$ ]	$\text{suf}_{\text{SA}[r]}(y)$	LCP[ $r$ ]	LPF[SA[ $r$ ]]
0	2	aabababbabb	0	1
1	0	abaabababbabb	1	0
2	3	abababbabb	3	3
3	5	ababbabb	4	4
4	7	abbabb	2	2
5	10	abbb	3	3
6	13	b	0	1
7	1	baabababbabb	1	0
8	4	bababbabb	2	2
9	6	babbabb	3	3
10	9	babb	4	4
11	12	bb	1	2
12	8	bbabb	2	1
13	11	bbb	2	2



**Fig. 1.** Graph of the permutation of positions of  $\text{abaabababbabb}$  (Example 2) in the lexicographic order of the corresponding suffixes. Labels of edges are Longest Common Prefix lengths between consecutive suffixes.

The computation of the Suffix Array can be done in time  $O(n \log n)$  in the comparison model [16] (see [3,7,9]). If the string is over an alphabet of integers in the range  $[0, n^c]$  for some constant  $c$ , then the Suffix Array can be built in time  $O(n)$  [10,12,13,18] (see also [3]). An elegant algorithm for building the LCP array in linear time is given in [11].

2.1. Graphic representation

The Suffix Array of the string  $y$  has a nice graphic representation that helps understand the algorithms computing the LPF array.

Fig. 1 shows the Suffix Array representation for the string  $y = \text{abaabababbabb}$ . The abscissa axis refers to ranks of suffixes and the ordinate axis refers to their positions. The sorted list of suffixes is plotted by their positions, and consecutive ranks are linked by an edge whose label is the associated LCP value.

### 3. Computing LPF using SA and LCP

The LCP array satisfies the following simple property, a consequence of the lexicographic ordering of suffixes in the SA: the length of the lcp value between two positions at ranks  $r$  and  $t$ ,  $r < t$ , that is,  $|\text{lcp}(\text{SA}[r], \text{SA}[t])|$  is the minimal value in  $\text{LCP}[r + 1..t]$ .

For a rank  $r$ , let us define  $\text{prev}[r]$  as the largest rank  $s < r$ , for which  $\text{SA}[s] < \text{SA}[r]$ , if it exists, and as  $-1$  otherwise. Let us also define the dual notion,  $\text{next}[r]$ , as the smallest rank  $t > r$  for which  $\text{SA}[t] < \text{SA}[r]$ , if it exists, and as  $n$  otherwise. The two arrays  $\text{prev}$  and  $\text{next}$  for our running example are shown in Table 1.

The above property of the LCP array implies that to compute  $\text{LPF}[\text{SA}[r]]$  there is no need to look at ranks smaller than  $\text{prev}[r]$  or greater than  $\text{next}[r]$ , that is,

$$\text{LPF}[\text{SA}[r]] = \max\{|\text{lcp}(\text{SA}[r], \text{SA}[\text{prev}[r]])|, |\text{lcp}(\text{SA}[r], \text{SA}[\text{next}[r]])|\},$$

where undefined LCP values are assumed to be 0. In particular, if a position  $i$  is a “peak” at rank  $r$  in the graphic representation of the Suffix Array ( $\text{SA}[r] = i$ ) we get

$$\text{LPF}[i] = \max\{\text{LCP}[r], \text{LCP}[r + 1]\}.$$

Also, the minimum of the two values is the length of the lcp between positions at ranks  $r - 1$  and  $r + 1$ , if defined, that is,  $|\text{LCP}(\text{SA}[r - 1], \text{SA}[r + 1])|$ . This gives the idea underlying the algorithm LPF-SIMPLE, which is similar to the one of [4].

```

LPF-SIMPLE(SA, LCP, n)
1  LCP-store  $\leftarrow$  LCP
2  LCP[n]  $\leftarrow$  0
3  ISA  $\leftarrow$  inverse SA
4  for  $r \leftarrow 0$  to  $n - 1$  do
5      prev[r]  $\leftarrow r - 1$ 
6      next[r]  $\leftarrow r + 1$ 
7  for  $i \leftarrow n - 1$  downto 0 do
8       $r \leftarrow$  ISA[i]
9      LPF[i]  $\leftarrow$  max{LCP[r], LCP[next[r]]}
10     LCP[next[r]]  $\leftarrow$  min{LCP[r], LCP[next[r]]}
11     if prev[r]  $\geq$  0 then
12         next[prev[r]]  $\leftarrow$  next[r]
13     if next[r] < n then
14         prev[next[r]]  $\leftarrow$  prev[r]
15  return LPF

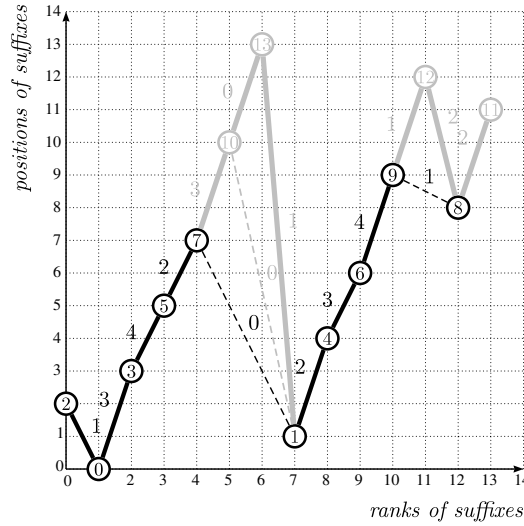
```

The algorithm computes the two arrays  $\text{prev}$  and  $\text{next}$  at the same time as computing the LPF array. Steps 9 and 10 implement the earlier discussion and it is enough to explain why the  $\text{prev}$  and  $\text{next}$  arrays are correctly computed. The initialization in steps 4–6 gives the correct values for the peaks only. However, in the *for* cycle at step 7, the peaks are considered in decreasing order of their height (that is, position) and hence the updates in steps 11–14 will provide the correct values.

For instance, the algorithm starts with the highest peak in the graph. The LPF value can already be computed for this position. In Fig. 1, the highest peak has position 13 and rank 6, that is,  $\text{SA}[6] = 13$ . Neighboring positions are lower and therefore  $\text{prev}[6] = 5$  and  $\text{next}[6] = 7$  are correct. Using the above formula for LPF, we have  $\text{LPF}[13] = \max\{0, 1\} = 1$ . Once this value is set, we remove the node from the graph by adding an edge between its neighbors, 10 and 1, labeled by the minimum of the two LCP values,  $\min\{0, 1\} = 0$ . The values  $\text{next}[5] = 7$  and  $\text{prev}[7] = 5$  are also set, the former being already the correct value (due to the peak position 10 at rank 5).

The algorithm LPF-SIMPLE runs in linear time but requires the extra arrays  $\text{prev}$ ,  $\text{next}$  and ISA in addition to its input and output. Since the LCP array is modified during the algorithm, it can be stored (step 1) in case it is needed later.

Fig. 2 illustrates several steps of the algorithm when the input is our running example abaabababbabbb.



**Fig. 2.** Illustration of the run of the algorithm LPF-SIMPLE on the string abaabababbabbb (see Fig. 1) after processing positions 13, 12, 11, 10. Gray positions and gray edges are no longer considered. The next step is to process position 9 at rank 10. According to the weights of edges pending from 9, we have  $LPF[9] = \max\{4, 1\} = 4$ . Positions 6 and 8 will be linked, through the prev and next arrays on their ranks, with an edge of weight  $\min\{4, 1\} = 1$ , which is indeed  $LPF[8]$ .

**Proposition 1.** The algorithm LPF-SIMPLE computes the LPF array of a string of length  $n$  from its Suffix Array in linear time and space. It requires  $3n + O(1)$  integer cells in addition to its input and output.

### 3.1. An alternative algorithm

The array next can alternatively be pre-computed by the following procedure NEXT whose linear-time behavior is an interesting exercise left to the reader. Only the code in black is needed for the computation of next.

When all codes (black and gray) are used, the procedure computes also the array  $LPF_{>}$  (defined in [4]) that is similar to LPF but accounts for larger ranks only. It is defined, for  $r = 0, \dots, n - 1$ , by

$$LPF_{>}[r] = |\text{lcp}(SA[r], SA[\text{next}[r]])|.$$

NEXT(SA, LCP,  $n$ )

```

1 next[n - 1] ← n
2 LPF_{>}[n - 1] ← 0
3 for r ← n - 2 downto 0 do
4   t ← r + 1
5   ℓ ← LCP[t]
6   while t < n and SA[t] > SA[r] do
7     ℓ ← min{ℓ, LPF_{>}[t]}
8     t ← next[t]
9   next[r] ← t
10  LPF_{>}[r] ← ℓ
11 return next, LPF_{>}

```

The computation of the dual prev array is done symmetrically. Also, the array  $LPF_{<}$  is symmetrically defined and computed. When both arrays  $LPF_{<}$  and  $LPF_{>}$  are available, the computation of the LPF array is done using the following equality:

$$LPF[i] = \max\{LPF_{>}[r], LPF_{<}[r]\}.$$

#### 4. On-line computation

Techniques of the previous sections to compute the LPF array are simple but space consuming. In this section and the next one we address this issue. We show that an on-line computation on the Suffix Array using a stack reduces the memory space to only  $O(\sqrt{n})$  for a string of length  $n$ . This algorithm is similar to the one of [6] with some differences in view of our time–space optimal algorithm in the next section.

The design of the on-line computation still relies on the property used for the algorithm of Section 2 and related to “peaks” (see lines 6–8). It relies additionally on another straightforward property that we describe now. Assume that a position  $SA[r]$  at rank  $r$  satisfies  $LCP[r] \geq LCP[r + 1]$ . Then, no position after it in the list can provide a larger LCP value and therefore we get  $LPF[SA[r]] = LCP[r]$ . This is implemented in lines 10–11 of the algorithm LPF-ON-LINE.

Note that lines 15–18 may be removed from the algorithm LPF-ON-LINE if the Suffix Array can be extended to rank  $n$  and initialized to  $SA[n] = -1$  and  $LCP[n] = 0$ . But we prefer the present design that is compatible with the algorithm of the next section.

```

LPF-ON-LINE(SA, LCP, n)
1  EMPTYSTACK(S)
2  for  $r \leftarrow 0$  to  $n - 1$  do
3     $r\text{-lcp} \leftarrow LCP[r]$ 
4    while not EMPTY(S) do
5       $(t, t\text{-lcp}) \leftarrow \text{TOP}(S)$ 
6      if  $SA[r] < SA[t]$  then
7         $LPF[SA[t]] \leftarrow \max\{t\text{-lcp}, r\text{-lcp}\}$ 
8         $r\text{-lcp} \leftarrow \min\{t\text{-lcp}, r\text{-lcp}\}$ 
9        POP(S)
10     elseif  $(SA[r] > SA[t])$  and  $(r\text{-lcp} \leq t\text{-lcp})$  then
11        $LPF[SA[t]] \leftarrow t\text{-lcp}$ 
12       POP(S)
13     else break
14   PUSH(S, (r, r-lcp))
15   while not EMPTY(S) do
16      $(t, t\text{-lcp}) \leftarrow \text{TOP}(S)$ 
17      $LPF[SA[t]] \leftarrow t\text{-lcp}$ 
18     POP(S)
19   return LPF

```

The extra memory space used by the algorithm LPF-ON-LINE to compute the LPF array of a string is occupied by the stack and a constant number of integer variables. To evaluate the total size required by the algorithm it is then important to determine the maximal size of the stack for a string of length  $n$ . It is proved in [6] that this quantity is  $O(\sqrt{n})$ .

For most values of  $n$  there are plenty of strings for which the stack reaches its maximal size. But if  $n$  is of the form  $k(k + 1)/2$ , that is, if it is the sum of the first  $k$  positive integers, then there is a unique string on the alphabet  $\{a, b\}$  (with  $a < b$ ) giving the maximal size stack. This word is  $aabab^2 \dots ab^{k-1}$  and the maximal stack size is  $k$ .

The maximal stack sizes for strings of lengths 4–22 are given in Table 3.

**Proposition 2.** *The algorithm LPF-ON-LINE computes the LPF array of a string of length  $n$  from its Suffix Array in linear time and  $O(\sqrt{n})$  space. It requires less than  $2\sqrt{2n} + O(1)$  integer cells in addition to its input and output.*

#### 5. The time–space optimal algorithm

In this section we show that the computation of the LPF array of a string can be implemented with only constant memory space in addition to the SA, LCP, and LPF arrays. The underlying property



used for this purpose is the  $O(\sqrt{n})$  stack size reported in the previous section. The property allows an implementation of the stack inside the LCP array for a sufficiently large part of the string. The rest of the computation for the remaining positions is done in a more time-expensive manner but for a small part of the string. This preserves the linear running time of the whole computation.

LPF-OPTIMAL(SA, LCP,  $n$ )

```

1  ▷ it is assumed that  $n \geq 8$ 
2   $K \leftarrow \lfloor n - 2\sqrt{2n} \rfloor$ 
3  ▷ next three procedures share the same LCP array
4  LPF-ON-LINE(SA, LCP,  $K$ )
5  LPF-NAIVE(SA, LCP,  $K, n$ )
6  LPF-ANCHORED(SA, LCP,  $K, n$ )
7  return LPF

```

It is rather clear that only constant extra memory space is required to implement the strategy retained by the algorithm LPF-OPTIMAL. The choice of the parameter  $K$  is a result of the previous section and is done to let enough space in the LCP array to implement the stack used by the algorithm LPF-ON-LINE.

Although not done here, the choice of the parameter  $K$  can be dynamic and done during the algorithm LPF-ON-LINE as  $n - 2k - 1$ , where  $k$  is the size of the stack just before executing line 15. This certainly reduces the actual running time but does not improve its asymptotic evaluation.

In the algorithm LPF-OPTIMAL, the stack of the procedure LPF-ON-LINE is implemented in the LCP array. Access to the array is done via the SA array. Doing so, the stack is treated like a continuous space LPF[SA[ $K..n - 1$ ]]. Elements are stored sequentially so that the elementary stack operations (empty, top, push, pop) are all executed in constant time. Therefore, the running time of the first step is  $O(n)$  (indeed  $O(K)$ ) as for the algorithm LPF-ON-LINE.

**Example 3.** Table 4 shows the content of the LCP array at two stages of the run of the procedure LPF-ON-LINE on Example 2: (i) immediately after processing rank 5 and (ii) at the end of the first step.

The row (i) shows that LCP values of positions 2, 3, 5 at respective ranks 0, 2, 3 have already been computed. The part LPF[SA[8..13]] of the array stores the content of the stack: ((1, 0), (4, 2), (5, 3)). The row (ii) shows that values have already been computed for positions at ranks 0–6. The content of the stack is ((7, 0), (8, 2)). On this example,  $K$  could have been set dynamically to 5.

### 5.1. Naive computation

The second step of the algorithm LPF-OPTIMAL processes the Suffix Array from rank  $K$ . For each rank  $r$ , the values  $\text{prev}[r]$ ,  $\text{next}[r]$  and the corresponding LCP values are computed starting from  $r$  and going backward and forward, respectively. The code is given below for the sake of completeness.

LPF-NAIVE(SA, LCP,  $K, n$ )

```

1  LCP[ $n$ ]  $\leftarrow$  0
2  for  $r \leftarrow K$  to  $n - 1$  do
3       $left \leftarrow$  LCP[ $r$ ]
4       $s \leftarrow r - 1$ 
5      while  $s \geq K$  and SA[ $s$ ] > SA[ $r$ ] do
6           $left \leftarrow$  min{ $left$ , LCP[ $s$ ]}
7           $s \leftarrow s - 1$ 
8      if  $s = K - 1$  then
9           $left \leftarrow$  0
10      $t \leftarrow r + 1$ 
11      $right \leftarrow$  LCP[ $t$ ]
12     while  $t < n$  and SA[ $t$ ] > SA[ $r$ ] do
13          $t \leftarrow t + 1$ 
14          $right \leftarrow$  min{ $right$ , LCP[ $t$ ]}
15     LPF[SA[ $r$ ]]  $\leftarrow$  max{ $left$ ,  $right$ }
16 return LPF

```

**Table 3**

Maximum stack size for various string lengths.

Length $n$	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
Max-stack-size	2	2	3	3	3	3	4	4	4	4	4	5	5	5	5	5	5	6	6

**Table 4**

The content of the LPF table during the procedure LPF-ON-LINE.

	Rank $r$	0	1	2	3	4	5	6	7	8	9	10	11	12	13
(i)	LPF[SA[ $r$ ]]	1		3	4					5	3	4	2	1	0
(ii)	LPF[SA[ $r$ ]]	1	0	3	4	2	3	1				8	2	7	0

The process takes  $O(n-K)$  time per rank, and hence the total running time of the step is  $O((n-K)^2)$ , which is  $O(n)$  since  $K = n - 2\sqrt{2n}$ .

### 5.2. Completing the computation

The first two steps of the algorithm LPF-OPTIMAL process independently two segments of the Suffix Array. The last step consists in joining their results, which requires updating some LPF values. Indeed, for a rank  $r$ ,  $r < K$ ,  $\text{next}[r]$  can be in  $[K, n - 1]$ , which implies that the computation of LPF[SA[ $r$ ]] might not be achieved. The same phenomenon happens for a rank in the second part, whose associated prev rank is smaller than  $K$ .

The next algorithm updates all LPF values and completes the whole computation. It assumes that the LPF calculation has been done independently on parts LPF[SA[0.. $K - 1$ ]] and LPF[SA[ $K$ .. $n - 1$ ]] of the array, which is realized by the first two steps on the algorithm LPF-OPTIMAL.

LPF-ANCHORED(SA, LCP,  $K$ ,  $n$ )

```

1   $s \leftarrow K - 1$ 
2   $t \leftarrow K$ 
3   $\ell \leftarrow \text{LCP}[K]$ 
4  while ( $s \geq 0$ ) and ( $t < n$ ) do
5      ▷ invariant:  $\ell = \min \text{LCP}[s + 1, \dots, t]$ 
6      if SA[ $s$ ] < SA[ $t$ ] then
7          LPF[SA[ $t$ ]]  $\leftarrow \max\{\text{LPF}[\text{SA}[t]], \ell\}$ 
8           $t \leftarrow t + 1$ 
9           $\ell \leftarrow \min\{\ell, \text{LCP}[t]\}$ 
10     else LPF[SA[ $s$ ]]  $\leftarrow \max\{\text{LPF}[\text{SA}[s]], \ell\}$ 
11          $\ell \leftarrow \min\{\ell, \text{LCP}[s]\}$ 
12          $s \leftarrow s - 1$ 
13 return LPF
    
```

The running time of this last step LPF-ANCHORED is obviously linear,  $O(n)$ , as are the other steps of the algorithm LPF-OPTIMAL.

The first conclusion of the section is the following statement.

**Theorem 1.** *The Longest Previous Factor array of a string of length  $n$  on an integer alphabet can be built from the read–write Suffix Array and Longest Common Prefix array in time  $O(n)$  (independently of the alphabet size) with a constant amount of extra memory space.*

The algorithm LPF-OPTIMAL uses the Suffix Array of the input string in a read-only manner but does not use the LPF array in a write-only manner. If this last condition is to be satisfied, the question remains of whether there exists a linear-time LPF array construction running with constant extra space. We get this feature if the algorithm LPF-ANCHORED is applied recursively by dividing the Suffix Array into two equal parts. The running time becomes  $O(n \log n)$  in the model of computation allowing priority writes.

**Proposition 3.** *The Longest Previous Factor array of a string of length  $n$  on an integer alphabet can be built from its read-only Suffix Array and Longest Common Prefix array in time  $O(n \log n)$  (independently of the alphabet size) with a constant amount of extra memory space and with a write-only output.*

Despite the use of the output as auxiliary storage in the ultimate linear-time algorithm, the series of algorithms described in the article provide a large range of efficient solutions that meet many practical needs.

## 6. Computing LPF by permuting LCP

It has been noticed in [4] that the content of the LPF array is the same as that of the LCP array up to some permutation. Then the following questions arise:

- Can this permutation be computed efficiently?
- Does it depend on the string, its Suffix Array or maybe just its Longest Common Prefix array?

It turns out that, for fixed SA and LCP arrays, it does not depend on the actual string. It is possible to construct such a sorting network, whose shape depends only on the Suffix Array, that transforms the LCP array into the LPF array. This observation leads to the algorithm LPF-BY-PERMUTING-LCP, producing LPF by permuting the elements of LCP.

In the algorithm LPF-BY-PERMUTING-LCP below we assume that the array `next` has been pre-computed by the procedure NEXT from Section 3.1 (in which the gray instructions, related to  $\text{LPF}_{>}$ , are removed, as well as the parameter LCP).

The array `next` has been defined as the next rank with a smaller suffix. It helps defining another array, `nextPos`, that behaves similarly but that is indexed by positions instead of ranks. More precisely, it computes, for each position  $i$ , the next closest position `nextPos[i]` in the Suffix Array that is smaller than  $i$ , that is,

$$\text{nextPos}[\text{SA}[r]] = \text{SA}[\min\{t \mid t > r \text{ and } \text{SA}[t] < \text{SA}[r]\}] = \text{SA}[\text{next}[r]].$$

Equivalently, if we set the position  $i = \text{SA}[r]$ , then

$$\text{nextPos}[i] = \text{SA}[\text{next}[\text{ISA}[i]]].$$

Initially, the LPF array is a copy of the LCP array permuted according to the SA array. The algorithm permutes the elements of the array to get the correct LPF array. This involves a modification of the “peaks” idea used previously. If we look again at Fig. 1, the highest peak (position 13) has the two adjacent edges labeled 0 and 1. The higher of the two values becomes  $\text{LPF}[13]$  whereas the smaller one is the length of the lcp between suffixes at ranks 5 and 7. Another way to do this is by permuting the labels 0 and 1 such that each goes to the right place, that is, 1 is on the edge to the left of 13 (and will remain there as the correct LPF value) and 0 goes on the edge to the left of 1, as the new LCP value.

```

LPF-BY-PERMUTING-LCP(SA, LCP, n)
1  SA[n] ← -1
2  π ← SA
3  for r ← 0 to n - 1 do
4    LPF[SA[r]] ← LCP[r]
5    nextPos[SA[r]] ← SA[next[r]]
6  for i ← n - 1 downto 0 do
7    if (nextPos[i] ≥ 0) and (LPF[i] < LPF[nextPos[i]]) then
8      EXCHANGE(LPF[i], LPF[nextPos[i]])
9      π ← π ◦ (LPF[i], LPF[nextPos[i]])
10 return LPF, π

```

Note that the algorithm LPF-BY-PERMUTING-LCP computes also explicitly the permutation  $\pi$  to transform the LCP array into the LPF array, that is, the permutation that satisfies  $\text{LPF}[\pi[i]] = \text{LCP}[i]$ .

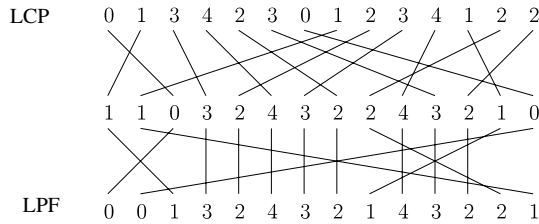


Fig. 3. The permutation of LCP that gives LPF, as computed by the algorithm LPF-BY-PERMUTING-LCP.

The final permutation is obtained by composing the SA (which is a permutation) with the transpositions produced in step 8. For our example, this process is shown in Fig. 3. The actual final permutation is

$$\begin{pmatrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ 0 & 2 & 3 & 5 & 7 & 10 & 1 & 13 & 4 & 6 & 9 & 8 & 12 & 11 \end{pmatrix}.$$

The algorithm LPF-BY-PERMUTING-LCP uses only one integer array in addition to its input and output since the next array it uses is substituted for the nextPos array. We have proved the following result.

**Theorem 2.** *The algorithm LPF-BY-PERMUTING-LCP computes the LPF array of a string of length  $n$  from its Suffix Array in linear time and space. It requires  $n + O(1)$  integer cells in addition to its input and output.*

Note that the elements of LPF are exchanged (lines 7–9) only if they are not in increasing order. Which elements are compared, depends on values in nextPos, and this in turn depends only on the Suffix Array. So, for a given Suffix Array, one can construct a sorting network implementing lines 6–9 of the algorithm LPF-BY-PERMUTING-LCP. Hence, the following proposition holds.

**Proposition 4.** *For a given Suffix Array of a string, there exists a sorting network processing a sequence of  $n$  numbers in such a way that it transforms the LCP array into the LPF array. Moreover, the shape of the sorting network depends only on the Suffix Array, but not on its LCP array nor on the actual string.*

## 7. Conclusion

A number of interesting problems remain open. First, it is easy to compute the LZ77 factorization using the LPF array. However, no way of doing the opposite computation is known. The LZ77 factorization contains less information and therefore, it may not be of great use in constructing the LPF array.

Second, our algorithm LPF-OPTIMAL is time–space optimal assuming the SA and LCP arrays are used for computing the LPF array. It is of great interest to find fast algorithms that compute the LPF array without such information or, more generally, without using any string indexes.

## Acknowledgments

The authors would like to warmly thank German Tischler for his careful inspection of the algorithms described in the article. Thanks are also due to the anonymous referee for pointing out the head array in McCreight’s paper.

The second author’s research was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC). The fifth author was supported by grant N206 566740 of the National Science Centre in Poland.

## References

- [1] J. Berstel, A. Savelli, Crochemore factorization of Sturmian and other infinite words, in: R. Kralovic, P. Urzyczyn (Eds.), Mathematical Foundations of Computer Science, in: Lecture Notes in Computer Science, vol. 4162, Springer, 2006, pp. 157–166.

- [2] M. Crochemore, Transducers and repetitions, *Theoret. Comput. Sci.* 45 (1) (1986) 63–86.
- [3] M. Crochemore, C. Hancart, T. Lecroq, *Algorithms on Strings*, Cambridge University Press, Cambridge, UK, 2007, p. 392.
- [4] M. Crochemore, L. Ilie, Computing longest previous factor in linear time and applications, *Inform. Process. Lett.* 106 (2) (2008) 75–80.
- [5] M. Crochemore, L. Ilie, C. Iliopoulos, M. Kubica, W. Rytter, T. Waleń, LPF computation revisited, in: J. Fiala, J. Kratochvíl, M. Miller (Eds.), *Proc. of IWOCA 2009*, in: *Lecture Notes in Comput. Sci.*, vol. 5874, Springer, Heidelberg, 2009, pp. 158–169.
- [6] M. Crochemore, L. Ilie, W.F. Smyth, A simple algorithm for computing the Lempel–Ziv factorization, in: J.A. Storer, M.W. Marcellin (Eds.), *18th Data Compression Conference*, IEEE Computer Society, Los Alamitos, CA, 2008, pp. 482–488.
- [7] M. Crochemore, W. Rytter, *Jewels of Stringology*, World Scientific Publishing, Hong-Kong, 2002, p. 310.
- [8] F. Franek, J. Holub, W.F. Smyth, X. Xiao, Computing quasi suffix arrays, *J. Autom. Lang. Comb.* 8 (4) (2003) 593–606.
- [9] D. Gusfield, *Algorithms on Strings, Trees and Sequences: Computer Science and Computational Biology*, Cambridge University Press, Cambridge, UK, 1997, p. 534.
- [10] J. Kärkkäinen, P. Sanders, Simple linear work suffix array construction, in: J.C.M. Baeten, J.K. Lenstra, J. Parrow, G.J. Woeginger (Eds.), *Automata, Languages and Programming*, in: *Lecture Notes in Computer Science*, vol. 2719, Springer, 2003, pp. 943–955.
- [11] T. Kasai, G. Lee, H. Arimura, S. Arikawa, K. Park, Linear-time longest-common-prefix computation in suffix arrays and its applications, in: *Combinatorial Pattern Matching*, in: *Lecture Notes in Computer Science*, vol. 2089, Springer, 2001, pp. 181–192.
- [12] D.K. Kim, J.S. Sim, H. Park, K. Park, Linear-time construction of suffix arrays, in: *Combinatorial Pattern Matching*, in: *Lecture Notes in Computer Science*, vol. 2676, Springer, 2003, pp. 186–199.
- [13] P. Ko, S. Aluru, Space efficient linear time construction of suffix arrays, in: *Combinatorial Pattern Matching*, in: *Lecture Notes in Computer Science*, vol. 2676, Springer, 2003, pp. 200–210.
- [14] R.M. Kolpakov, G. Kucherov, Finding maximal repetitions in a word in linear time, in: *FOCS*, 1999, pp. 596–604.
- [15] M.G. Main, Detecting leftmost maximal periodicities, *Discrete Appl. Math.* 25 (1989) 145–153.
- [16] U. Manber, G. Myers, Suffix arrays: a new method for on-line string searches, *SIAM J. Comput.* 22 (5) (1993) 935–948.
- [17] E.M. McCreight, A space-economical suffix tree construction algorithm, *J. ACM* 23 (2) (1976) 262–272.
- [18] G. Nong, S. Zhang, W.H. Chan, Linear time suffix array construction using  $D$ -critical substrings, in: G. Kucherov, E. Ukkonen (Eds.), *Combinatorial Pattern Matching*, in: *Lecture Notes in Computer Science*, vol. 5577, Springer, 2009, pp. 54–67.
- [19] M. Rodeh, V.R. Pratt, S. Even, Linear algorithm for data compression via string matching, *J. ACM* 28 (1) (1981) 16–24.
- [20] J. Storer, T. Szymanski, Data compression via textual substitution, *J. ACM* 29 (4) (1982) 928–951.
- [21] I. Witten, A. Moffat, T. Bell, *Managing Gigabytes*, Van Nostrand Reinhold, New York, 1994.
- [22] J. Ziv, A. Lempel, A universal algorithm for sequential data compression, *IEEE Trans. Inform. Theory* 23 (3) (1977) 337–343.
- [23] J. Ziv, A. Lempel, Compression of individual sequences via variable-rate coding, *IEEE Trans. Inform. Theory* 24 (5) (1978) 530–536.