



Longest previous overlapping factor array

Hideo Bannai^a, Shunsuke Inenaga^{b,c}, Neerja Mhaskar^{d,*}

^a M&D Data Science Center, Tokyo Medical and Dental University, Japan

^b Department of Informatics, Kyushu University, Japan

^c PRESTO, Japan Science and Technology Agency, Japan

^d Department of Computing and Software, McMaster University, Canada



ARTICLE INFO

Article history:

Received 31 December 2019

Received in revised form 18 November 2020

Accepted 13 January 2021

Available online 26 January 2021

Communicated by Abhi Shelat

Keywords:

Algorithms

Longest previous overlapping factor

Overlapping factor

Manhattan skyline problem

ABSTRACT

In this paper, we introduce the *longest previous overlapping factor array* of a string – a variant of the longest previous factor array. We show that it can be computed in linear time in the length of the input string, via a reduction to the Max-variant of the Manhattan skyline problem Crochemore et al. (2014) [5].

© 2021 Elsevier B.V. All rights reserved.

1. Introduction

Crochemore and Ilie in [3] introduced the *longest previous factor array* (LPF): given a string w of length n , for any position i in w , $\text{LPF}[i]$ stores the length of the longest prefix of $w[i..n]$ that has a previous occurrence which begins in the position range $[1..i-1]$ of w . Several linear time algorithms have been proposed to compute the LPF array [3,4]. Of the many applications of the LPF array, one of its most important application is its use in the computation of the famous Lempel-Ziv factorization [10] used for data compression.

Later, an important variant of the LPF array, called the *longest previous non-overlapping factor array* (LPnF) was proposed [6,7]: $\text{LPnF}[i]$ stores the length of the longest prefix of $w[i..n]$ that has a previous occurrence which ends in position range $[1..i-1]$ of w . Thus, for any $1 \leq i \leq n$, the corresponding factor $w[i..i + \text{LPnF}[i] - 1]$ must have a

previous occurrence that *does not* overlap with position i . It is known that the LPnF array of a string of length n can be computed in $\mathcal{O}(n)$ time [6,7]. The main application of LPnF is the Lempel-Ziv factorization *without self-references* (a.k.a. *f-factorization*) [2], which plays a central role in the solutions to various string problems such as approximation of the smallest grammar-based compression [9], and computation of all maximal repetitions (a.k.a. runs) [8] in a given string, just to mention a few. Other variants of the LPF array have been proposed in the literature, for example the longest reverse factor array [6], which can also be computed in linear time in the length of the input string.

In this paper, we discuss another variant of LPF, called the *longest previous overlapping factor array* (LPoF) of a given string. Namely, $\text{LPoF}[i]$ stores the length of the longest prefix of $w[i..n]$ that has a previous occurrence which ends in the position range $[i..i + \text{LPoF}[i] - 2]$ of w . Thus, for any $1 \leq i \leq n$, the corresponding factor $w[i..i + \text{LPoF}[i] - 1]$ must have a previous occurrence that *does* overlap with position i . The LPoF array is a natural extension to the LPF array and can be seen as a complement to LPnF, since $\text{LPF}[i] = \max\{\text{LPnF}[i], \text{LPoF}[i]\}$ always holds. While there exists a simple $\mathcal{O}(n)$ -time algo-

* Corresponding author.

E-mail addresses: hdbn.dsc@tmd.ac.jp (H. Bannai), inenaga@inf.kyushu-u.ac.jp (S. Inenaga), pophlin@mcmaster.ca (N. Mhaskar).

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
\mathbf{w}	a	b	a	a	a	b	a	b	a	a	a	b	a	b	a	a	b
L _{PF}	0	0	1	2	3	2	10	9	8	7	6	5	4	3	3	2	1
L _{PnF}	0	0	1	1	3	2	6	6	6	6	6	5	4	3	3	2	1
L _{PoF}	0	0	0	2	0	0	10	9	8	7	0	0	3	0	0	0	0

Fig. 1. The L_{PF}, L_{PnF}, and L_{PoF} arrays for string $\mathbf{w} = \text{abaaababaaababaab}$.

rithm which computes the L_{PF} array from the Suffix Array and the Longest Common Prefix Array [3,4], the above equation suggests yet another approach for computing the L_{PF} array.

On the other hand, we remark that the L_{PoF} array cannot be computed trivially from either or both the L_{PF} and L_{PnF} arrays. This poses an interesting theoretical question: Can we compute the L_{PoF} array in $\mathcal{O}(n)$ time? We answer this question affirmatively, by presenting a linear-time algorithm which computes the L_{PoF} array using a linear-time algorithm for solving a variant of the Manhattan skyline problem [5].

The outline of the paper is as follows: in Section 2, we give the mathematical background and briefly discuss a variant of the Manhattan skyline problem. In Section 3, we present our algorithm to compute the L_{PoF} array.

2. Preliminaries

2.1. Strings

A string \mathbf{w} is an ordered sequence of letters or symbols chosen from a finite totally ordered set of Σ , called an alphabet. The length of \mathbf{w} is written $|\mathbf{w}| = n$. Throughout this paper, we assume that Σ is the integer alphabet of polynomial size in n .

We also write a string \mathbf{w} of length n in an array form as $\mathbf{w}[1..n]$. We denote by $\mathbf{w}[i]$ the i -th letter of \mathbf{w} . A string $\mathbf{w}[i..j]$, where $1 \leq i \leq j \leq n$, is called a substring or a factor of $\mathbf{w}[1..n]$. A substring $\mathbf{w}[i..j]$ is a proper substring of \mathbf{w} if $j - i + 1 < n$. A substring $\mathbf{w}[1..j]$, where $1 \leq j \leq n$ is called a prefix of \mathbf{w} . A substring $\mathbf{w}[i..n]$, where $1 \leq i \leq n$ is called a suffix of \mathbf{w} .

A positive integer p is said to be a period of a string \mathbf{w} of length n if $\mathbf{w}[i] = \mathbf{w}[i + p]$ for any $1 \leq i \leq n - p$. A non-empty string \mathbf{w} of even length is said to be a square if \mathbf{w} has a period $\frac{|\mathbf{w}|}{2}$ (or equivalently if \mathbf{w} is of the form \mathbf{xx}). A substring $\mathbf{u} = \mathbf{w}[i..j]$ of \mathbf{w} is called a repetition in \mathbf{w} if the smallest period p of \mathbf{u} satisfies $p \leq \frac{|\mathbf{u}|}{2} = \frac{j-i+1}{2}$. A repetition $\mathbf{u} = \mathbf{w}[i..j]$ of \mathbf{w} is said to be a maximal repetition or a run in \mathbf{w} if the period p cannot be extended to the left or to the right, i.e., $i = 1$ or $\mathbf{w}[i - 1] \neq \mathbf{w}[i + p - 1]$, and $j = n$ or $\mathbf{w}[j + 1] \neq \mathbf{w}[j - p + 1]$. Each run \mathbf{r} in \mathbf{w} is represented by a triple (i, j, p) such that $\mathbf{r} = \mathbf{w}[i..j]$ and p is the smallest period of \mathbf{r} .

Theorem 1 ([8,1]). For any string \mathbf{w} of length n , the number of runs in \mathbf{w} is $\mathcal{O}(n)$. Also, all runs in \mathbf{w} can be computed in $\mathcal{O}(n)$ time and space if \mathbf{w} is drawn from an integer alphabet of polynomial size in n .

The paper [1] presents a new and much simpler algorithm for calculating runs.

2.2. Longest previous overlapping factor array

Given a string \mathbf{u} of the form $\mathbf{u} = \mathbf{vx} = \mathbf{yv}$, where $|\mathbf{u}| > |\mathbf{v}| > |\mathbf{x}|$, \mathbf{v} is called an overlapping factor. For example, given the string $\mathbf{u} = \text{abacabacaba} = (\text{abac})^2\text{aba}$, abacaba is an overlapping factor of \mathbf{u} .

Given a string \mathbf{w} of length n , the Longest Previous Overlapping Factor array (L_{PoF}) is an integer array of length n where each element

$$\text{L}_{\text{PoF}}[i] = \max\{\ell \geq 2 \mid \mathbf{w}[k..k + \ell - 1] = \mathbf{w}[i..i + \ell - 1], \\ k < i \leq k + \ell - 1\}$$

if such ℓ exists, otherwise it is 0.

Fig. 1 gives the L_{PF}, L_{PnF}, and L_{PoF} arrays corresponding to the string $\mathbf{w} = \text{abaaababaaababaab}$. It is easy to see that the L_{PF} array can be trivially computed from both the L_{PnF} and L_{PoF} arrays, as the i -th element of the L_{PF} array is simply the maximum of the i -th elements of the L_{PnF} and L_{PoF} arrays; that is, $\text{L}_{\text{PF}}[i] = \max\{\text{L}_{\text{PnF}}[i], \text{L}_{\text{PoF}}[i]\}$ for any $1 \leq i \leq n$.

Note however that the L_{PF} array can also be computed in linear time using the Suffix Array and the Longest Common Prefix Array, which is conceptually simpler than the approach mentioned above.

2.3. Manhattan skyline problem

Crochemore et al. [5] presented the Min-variant of the Manhattan Skyline Problem, and showed that it can be computed in linear time. Here we consider the Max-variant of the Manhattan Skyline Problem (MSP) as shown below.

Input: A set \mathcal{S} of $\mathcal{O}(n)$ subintervals of $[1..n]$, where each subinterval $[i..j]$ is associated with a natural weight, $\text{height}([i..j])$, of size $\mathcal{O}(n)$.

Output: A table $f[t] = \max\{\text{height}([i..j]) \mid t \in [i..j], [i..j] \in \mathcal{S}, t \in [1..n]\}$.

The Max-variant of MSP can be solved in $\mathcal{O}(n)$ time using the same strategy given in [5] to solve the Min-variant of MSP with the exception of sorting the intervals in non-increasing order instead of non-decreasing order. Consequently, we obtain the following lemma.

Lemma 2. The Max-variant of MSP can be solved in $\mathcal{O}(n)$ time.

3. Computing the L_{PoF} array

We follow a similar schema presented in [5] for computing internal local periods of a string, to compute the

LPoF array values; that is, we infer the LPoF array values from the run structure of \mathbf{w} . More specifically, we observe that any overlapping string in \mathbf{w} corresponds to one of the runs in \mathbf{w} . If $\text{LPoF}[i] = \ell$, then for some k such that $\ell > i - k$, $\mathbf{w}[k..i - 1] = \mathbf{w}[i..i + (i - k) - 1]$, and $\mathbf{w}[k..i + (i - k) - 1]$ is a square of period $i - k$. This also implies that there must exist a run (b, e, p) with period $p = i - k$ containing position i . Furthermore, it must satisfy $b + p \leq i \leq e - p$. The last condition is due to the fact that a run (b, e, p) with $\min\{1, i - p + 1\} \leq b$ or $e \leq \min\{i + p - 1, n\}$ does not contribute to the value of $\text{LPoF}[i]$. For example in Fig. 1, although the run $(1, 16, 6)$ contains the position $i = 11$, it does not satisfy the condition $b + p \leq i \leq e - p$, and therefore it does not contribute to the value of $\text{LPoF}[11]$.

The rationale for the last condition is as follows. For two occurrences of a substring \mathbf{u} , where $|\mathbf{u}| = \ell$, to overlap in a run, ℓ must be larger than the period p of the run. Consequently, these overlapping occurrences, say at indices i_1 and $i_2 > i_1$, are separated by p positions in the run; that is, $i_2 - i_1 = p$. However, for a substring occurring at an index $b \leq i < b + p$, its previous overlapping occurrence must be at $i - p < b$. This case is not possible as the run would be left extensible. Further, for a substring starting at index $e - p < i \leq e$, it has a length $< p$ in the run, and as result does not overlap with any of its occurrences in the run. Therefore, both these cases do not contribute to the $\text{LPoF}[i]$ values.

There can be more than one such run that contains position i , but the one that gives the largest ℓ is the one with the largest endpoint e . For example, in Fig. 1 the position $i = 7$ is contained in two runs $(1, 16, 6)$, and $(5, 9, 2)$ which satisfy the above conditions. Since we are interested in the largest ℓ , we consider the run $(1, 16, 6)$ with the largest endpoint $e = 16$.

To this end, we reduce the computation of each element of the LPoF array to an instance of the Max-variant of the Manhattan Skyline problem as follows:

$$\text{LPoF}[i] = \max\{\text{height}(t) \mid i \in t, t \in I\} - i + 1,$$

where I is the set of intervals corresponding to the runs in \mathbf{w} such that

$$I = \{[b + p..e - p] \mid (b, e, p) \text{ is a run in } \mathbf{w}, \\ \text{and } b + p \leq e - p\},$$

and $\text{height}(t) = e$ is the weight of the interval $t = [b + p..e - p]$ for each run (b, e, p) . Then, for any i , $\text{LPoF}[i]$ is the largest weight of all the intervals containing i , if such intervals exist; otherwise $\text{LPoF}[i] = 0$. This is an instance of the Max-variant of the Manhattan Skyline Problem shown above. By Lemma 2, it can be computed in $\mathcal{O}(n)$ time. Since we can compute all runs in \mathbf{w} in $\mathcal{O}(n)$ time by Theorem 1, LPoF array can be computed in $\mathcal{O}(n)$ time. Therefore, we have shown the following:

Theorem 3. For any string \mathbf{w} of length n , over an integer alphabet of polynomial size in n , the LPoF array of \mathbf{w} can be computed in $\mathcal{O}(n)$ time and space.

For example, the string \mathbf{w} given in Fig. 1 contains the following runs: $(1, 16, 6)$, $(3, 5, 1)$, $(5, 9, 2)$, $(9, 11, 1)$, $(11, 15, 2)$ and $(15, 16, 1)$. Then we get,

$$I = \{[7..10], [4..4], [7..7], [10..10], [13..13]\}$$

with weights 16, 5, 9, 11, and 15, respectively. Consider an index $i = 7$. This index is contained in intervals $[7..10]$ and $[7..7]$ with weights 16 and 9 respectively. Therefore, $\text{LPoF}[7] = 16 - 7 + 1 = 10$.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

This work is based on the authors' discussions during the StringMasters workshop held in Prague, August 29-31, 2018.

This work was supported by JSPS KAKENHI [Grant Numbers JP16H02783, 20H04141 (HB), JP17H01697 (SI)], and by JST PRESTO [Grant Number JPMJPR1922 (SI)].

References

- [1] H. Bannai, T. I. S. Inenaga, Y. Nakashima, M. Takeda, K. Tsuruta, The "Runs" theorem, *SIAM J. Comput.* 46 (5) (2017) 1501–1514.
- [2] M. Crochemore, C. Hancart, T. Lecroq, *Algorithms on Strings*, Cambridge University Press, 2007.
- [3] M. Crochemore, L. Ilie, Computing longest previous factor in linear time and applications, *Inf. Process. Lett.* 106 (2008) 75–80.
- [4] M. Crochemore, L. Ilie, C.S. Iliopoulos, M. Kubica, W. Rytter, T. Waleń, Computing the longest previous factor, *Electron. J. Comb.* 34 (2013) 15–26.
- [5] M. Crochemore, C.S. Iliopoulos, M. Kubica, J. Radoszewski, W. Rytter, T. Waleń, Extracting powers and periods in a word from its runs structure, *Theor. Comput. Sci.* 521 (2014) 29–41.
- [6] M. Crochemore, C.S. Iliopoulos, M. Kubica, W. Rytter, T. Waleń, Efficient algorithms for three variants of the LPF table, *J. Discret. Algorithms* 11 (2012) 51–61.
- [7] M. Crochemore, G. Tischler, Computing longest previous non-overlapping factors, *Inf. Process. Lett.* 111 (6) (2011) 291–295.
- [8] R.M. Kolpakov, G. Kucherov, Finding maximal repetitions in a word in linear time, in: *FOCS 1999*, 1999, pp. 596–604.
- [9] W. Rytter, Application of Lempel-Ziv factorization to the approximation of grammar-based compression, *Theor. Comput. Sci.* 302 (1–3) (2003) 211–222.
- [10] J. Ziv, A. Lempel, A universal algorithm for sequential data compression, *IEEE Trans. Inf. Theory* IT-23 (3) (1977) 337–349.