

Algorithms to Compute the Lyndon Array^{*}

Frantisek Franek¹, A. S. M. Sohidull Islam², M. Sohel Rahman³, and

W. F. Smyth^{1,3,4}

¹ Algorithms Research Group
Department of Computing & Software
McMaster University, Hamilton, Canada
{franek/smyth}@mcmaster.ca

² School of Computational Science & Engineering
McMaster University, Hamilton, Canada
sohansayed@gmail.com

³ Department of Computer Science & Engineering
Bangladesh University of Engineering & Technology
msrahman@cse.buet.ac.bd

⁴ School of Engineering & Information Technology
Murdoch University, Perth, Australia

Abstract. In the Lyndon array $\lambda = \lambda_{\mathbf{x}}[1..n]$ of a string $\mathbf{x} = \mathbf{x}[1..n]$, $\lambda[i]$ is the length of the longest Lyndon word starting at position i of \mathbf{x} . The computation of λ has recently become of great interest, since it was shown (Bannai *et al.*, **The “Runs” Theorem**) that the runs in \mathbf{x} are computable in linear time from $\lambda_{\mathbf{x}}$. Here we describe two algorithms for computing $\lambda_{\mathbf{x}}$ based on previous results known in different context, but for which no explicit exposition in this context had been given. These two algorithms execute in $\mathcal{O}(n^2)$ time in the worst case. The third algorithm presented that executes in $\Theta(n)$ time had been suggested and discussed previously, and we provide a more substantial discussion and prove of correctness for one of its steps. This algorithm achieves its linearity at the expense of prior computation of both the suffix array and the inverse suffix array of \mathbf{x} . We then go on to sketch a new algorithm and its two variants that avoids prior computation of global data structures and indicate that in worst-case these algorithms perform in $\mathcal{O}(n \log n)$ time.

1 Introduction

If $\mathbf{x} = \mathbf{u}\mathbf{v}$ for some \mathbf{u} and nonempty \mathbf{v} , then $\mathbf{v}\mathbf{u}$ is said to be the $|\mathbf{u}|^{\text{th}}$ *rotation* of \mathbf{x} , written $\mathbf{v}\mathbf{u} = R_{|\mathbf{u}|}(\mathbf{x})$. If there exists a string \mathbf{u} and an integer $e > 1$ such

^{*} This work was supported in part by the Natural Sciences & Engineering Research Council of Canada. The authors wish to thank Maxime Crochemore and Hideo Bannai for helpful discussions.

that $\mathbf{x} = \mathbf{u}^e$, then \mathbf{x} is said to be a *repetition*; otherwise \mathbf{x} is *primitive*. A primitive string \mathbf{x} that is lexicographically strictly least among all its rotations $R_k(\mathbf{x})$, $k = 0, 1, \dots, |\mathbf{x}|-1$, is said to be a *Lyndon word*.

The *Lyndon array* $\boldsymbol{\lambda} = \boldsymbol{\lambda}_{\mathbf{x}}[1..n]$ of a given nonempty string $\mathbf{x} = \mathbf{x}[1..n]$ gives at each position i the length of the longest Lyndon word starting at i . Note that equivalently we could store in the i th position of the Lyndon array the end position of the longest Lyndon word starting in i . We will use the notation $\mathcal{L}[i]$ to indicate the end position for the longest Lyndon word starting at i .

$$\begin{array}{cccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & & \\ \mathbf{x} & = & a & b & a & a & b & a & b & a & a & b & \\ \boldsymbol{\lambda} & = & 2 & 1 & 5 & 2 & 1 & 2 & 1 & 3 & 2 & 1 & \\ \mathcal{L} & = & 2 & 2 & 7 & 5 & 5 & 7 & 7 & 10 & 10 & 10 & \end{array} \tag{1}$$

Since being Lyndon really depends on the order of the underlying alphabet of the string, the Lyndon array of a string will change when we change the order of the alphabet. The Lyndon array has recently become of interest since Bannai *et al.* [2] showed that the two Lyndon arrays, one with respect to a given order of the alphabet and the other with respect to the inverse of that order, can be used to compute all the maximal periodicities (“runs”) in a string in linear time.

In this paper we describe four algorithms to compute $\boldsymbol{\lambda}_{\mathbf{x}}$. Section 2 makes various observations that apply generally to the Lyndon array and its computation. In Section 3 we describe two algorithms that are based on previous results known in a different context, and we present them here explicitly in the context of computing Lyndon arrays. These two algorithms perform in $\mathcal{O}(n^2)$ time in the worst case, where n is the length of the input string. Despite the high worst case complexity, in practice these algorithms perform very well as they are simple and straightforward to implement and do not require any complicated data structures; they could be characterized almost as in-place. The third algorithm discussed in this section had been described previously and we provide a more substantial discussion and prove correctness of one of its steps that we could not find anywhere in the literature. This algorithm is simple and worst-case linear-time, but requires suffix array construction and so is a little slower. Section 4 describes two variants of an algorithm we designed that uses only elementary data structures (no suffix arrays). One variant is $\mathcal{O}(n^2)$ in the worst case, the other indicates $\mathcal{O}(n \log n)$ time, but with no clear advantage in processing time. Section 5 describes the results of preliminary experiments on the algorithms; Section 6 outlines future work.

2 Preliminaries

Here we make various observations that apply to the algorithms described below.

Observation 1 *Let $\mathbf{x} = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_k$ be the Lyndon decomposition [5, 9] of \mathbf{x} , with Lyndon words $\mathbf{w}_1 \geq \mathbf{w}_2 \geq \cdots \geq \mathbf{w}_k$. Then every Lyndon word $\mathbf{x}[i.. \mathcal{L}[i]]$ of length $\boldsymbol{\lambda}[i]$ is a substring of some \mathbf{w}_h , $h \in 1..k$.*

Proof. For some $h \in 1..k-1$, consider \mathbf{w}_h with a nonempty proper suffix \mathbf{v}_h , and for some $t \in 1..k-h$, consider \mathbf{w}_{h+t} with a nonempty prefix \mathbf{u}_{h+t} . Since \mathbf{w}_h is a Lyndon word, $\mathbf{w}_h < \mathbf{v}_h$, and by lexorder, $\mathbf{u}_{h+t} \leq \mathbf{w}_{h+t}$. Thus $\mathbf{v}_h > \mathbf{w}_h \geq \mathbf{w}_{h+t} \geq \mathbf{u}_{h+t}$, and so $\mathbf{v}_h \mathbf{w}_{h+1} \cdots \mathbf{w}_{h+t-1} \mathbf{u}_{h+t}$ cannot be a Lyndon word for any choice of h or t . \square

Therefore to compute $\mathcal{L}\mathbf{x}$ it suffices to consider separately each distinct element \mathbf{w}_h in the Lyndon decomposition of \mathbf{x} . Hence, without loss of generality we suppose that \mathbf{x} is a Lyndon word and write it in the form $\mathbf{x}_1 \mathbf{x}_2 \cdots \mathbf{x}_m$, where for each $r \in 1..m$, $|\mathbf{x}_r| = \ell_r$ and

$$\mathbf{x}_r[1] \leq \mathbf{x}_r[2] \leq \cdots \leq \mathbf{x}_r[\ell_r], \quad (2)$$

while for $1 \leq r < m$,

$$\mathbf{x}_r[\ell_r] > \mathbf{x}_{r+1}[1]. \quad (3)$$

We call \mathbf{x}_r a **range** in \mathbf{x} and the boundary between \mathbf{x}_r and \mathbf{x}_{r+1} a **drop**. We identify a position j in range \mathbf{x}_r , $1 \leq j \leq \ell_r$, with its equivalent position i in \mathbf{x} by writing $i = S_{r,j} = \sum_{r'=1}^{r-1} \ell_{r'} + j$.

Observation 2 *Let $i = S_{r,j}$ be a position in \mathbf{x} that corresponds to position j in range \mathbf{x}_r .*

- (a) *If $\mathbf{x}_r[j] = \mathbf{x}_r[\ell_r]$, then $\mathcal{L}[i] = i$.*
- (b) *Otherwise, $\mathcal{L}[i] = i'$, where i' is the final position in some range $\mathbf{x}_{r'}$, $r' \geq r$; that is, $i' = \sum_{s=1}^{r'} \ell_s$.*

Proof. (a) is an immediate consequence of (2) and (3). To prove (b), suppose that $\mathbf{x}[i.. \mathcal{L}[i]]$ is a maximum-length Lyndon word, where $\mathcal{L}[i]$ falls within range r' but $\mathcal{L}[i] < i'$. Since by (2) $\mathbf{x}[\mathcal{L}[i]] \leq \mathbf{x}[\mathcal{L}[i]+1]$, there are two consecutive Lyndon words $\mathbf{x}[i.. \mathcal{L}[i]]$, $\mathbf{x}[\mathcal{L}[i]+1]$ that by the Lyndon decomposition theorem [5] can be merged into a single Lyndon word $\mathbf{x}[i.. \mathcal{L}[i]+1]$. Thus $\mathbf{x}[i.. \mathcal{L}[i]]$ is not maximum-length, a contradiction. \square

We see then that if $\mathbf{x}_r[j] < \mathbf{x}_r[\ell_r]$, then $\mathbf{x}_r[j.. \ell_r]$ is a (not necessarily maximum-length) Lyndon word, and for $i = S_{r,j}$, $\mathcal{L}[i] \geq S_{r,\ell_r}$:

$$\begin{array}{ccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ \mathbf{x} & = & a & a & a & b & | & a & a & b & | & a & b & | & a & a & b & b \\ \mathcal{L} & = & 13 & 13 & 4 & 4 & 9 & 7 & 7 & 9 & 9 & 13 & 13 & 12 & 13 \end{array} \quad (4)$$

More generally, the integer interval $\langle i, \mathcal{L}[i] \rangle = i.. \mathcal{L}[i]$ satisfy a ‘‘Monge’’ property that is exploited by Algorithm NSV* (Section 4):

Observation 3 *Suppose positions i, j in $\mathbf{x}[1..n]$ satisfy $1 \leq i < j \leq n$. Then either $\mathcal{L}[i] \leq j$ or $\mathcal{L}[i] \geq \mathcal{L}[j]$: the intervals $\langle i, \mathcal{L}[i] \rangle$ and $\langle j, \mathcal{L}[j] \rangle$ are not overlapping.*

Proof. Suppose two such intervals do overlap. Then the maximum-length Lyndon words $\mathbf{w}_1 = \mathbf{x}[i..\mathcal{L}[i]]$ and $\mathbf{w}_2 = \mathbf{x}[j..\mathcal{L}[j]]$ have a nonempty overlap, so that we can write $\mathbf{w}_1 = \mathbf{uv}$, $\mathbf{w}_2 = \mathbf{vv}'$ for some nonempty \mathbf{v} . But then, by well-known properties of Lyndon words, $\mathbf{w}_1 < \mathbf{v} < \mathbf{w}_2 < \mathbf{v}'$, implying that $\mathbf{w}_1\mathbf{v}'$ is a Lyndon word, contradicting the assumption that \mathbf{w}_1 is of a maximal length. \square

Expressing a string in terms of its ranges has the same useful lexorder property that writing it in terms of its letters does:

Observation 4 *Suppose strings \mathbf{x} and \mathbf{y} are expressed in terms of their ranges: $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \cdots \mathbf{x}_m$, $\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \cdots \mathbf{y}_n$. Suppose further that for some least integer $r \in 1..\min(m, n)$, $\mathbf{x}_r \neq \mathbf{y}_r$. Then $\mathbf{x} < \mathbf{y}$ (respectively, $\mathbf{x} > \mathbf{y}$) according as $\mathbf{x}_r < \mathbf{y}_r$ (respectively, $\mathbf{x}_r > \mathbf{y}_r$).*

Proof. If $\mathbf{x}_r < \mathbf{y}_r$, then either

- (a) \mathbf{x}_r is a nonempty proper prefix of \mathbf{y}_r ; or
- (b) there is some least position j such that $\mathbf{x}_r[j] < \mathbf{y}_r[j]$.

In case (a), if $r = m$, then \mathbf{x} is actually a prefix of \mathbf{y} , so that $\mathbf{x} < \mathbf{y}$, while if $r < m$, then by (3), $\mathbf{x}_{r+1}[1] < \mathbf{y}_r[|\mathbf{x}_r|+1]$, and again $\mathbf{x} < \mathbf{y}$. In case (b) the result is immediate. The proof for $\mathbf{x}_r > \mathbf{y}_r$ is similar. \square

3 Basic Algorithms

Here we outline three algorithms for which no clear exposition in the context of Lyndon arrays is available in the literature. We remark that the Lyndon array computation is equivalent to “Lyndon bracketing”, for which an $\mathcal{O}(n^2)$ algorithm was described in [18].

3.1 Folklore — Iterated MaxLyn

This algorithm, see Figure 1, is based on Duval’s linear time algorithm for Lyndon factorization, [10] – it is the application of its first step which we refer to as **MaxLyn** since it returns the size of the longest Lyndon word starting at that position. This process is iterated for all positions in the input string and this thus gives immediately $\mathcal{O}(|x|^2)$ worst case complexity for an input string x . Since Duval’s algorithm is in-place, this algorithm is simple and almost in-place, except the space for the Lyndon array. Below, we sketch the reasons the algorithm provides the correct answer.

For a string \mathbf{x} of length n , recall that the *prefix table* $\pi[1..n]$ is an integer array in which for every $i \in 1..n$, $\pi[i]$ is the length of the longest substring beginning at position i of \mathbf{x} that matches a prefix of \mathbf{x} . Given a nonempty string \mathbf{x} on alphabet Σ , let us define $\mathbf{x}' = \mathbf{x}\$,$ where the sentinel $\$ < \mu$ for every letter $\mu \in \Sigma$.

Observation 5 \mathbf{x} is a Lyndon word if and only if for every $i \in 2..n$, $\mathbf{x}'[1+k] < \mathbf{x}'[i+k]$, where $k = \pi[i]$.

This result forms the basis of the algorithm given in Figure 1 that computes the length $\max \in 1..n - j + 1$ of the longest Lyndon factor at a given position j in $\mathbf{x}[1..n]$. Its efficiency is a consequence of the instruction $i \leftarrow i + k + 1$ that skips over positions in the range $i + 1..i + k - 1$, effectively assuming that for every position i^* in that range, $i^* + \pi[i^*] \leq i + k$. Lemma 11, given in Appendix 1, justifies this assumption. Simply repeating MaxLyn at every position j of \mathbf{x} gives a simple, fast $\mathcal{O}(n^2)$ time and $\mathcal{O}(1)$ additional space algorithm to compute $\lambda_{\mathbf{x}}$.

```

procedure MaxLyn( $\mathbf{x}[1..n], j, \Sigma, \prec$ ) : integer
 $i \leftarrow j + 1$ ;  $max \leftarrow 1$ 
while  $i \leq n$  do
   $k \leftarrow 0$ 
  while  $\mathbf{x}'[j+k] = \mathbf{x}'[i+k]$  do
     $k \leftarrow k + 1$ 
  if  $\mathbf{x}'[j+k] \prec \mathbf{x}'[i+k]$  then
     $i \leftarrow i + k + 1$ ;  $max \leftarrow i - 1$ 
  else
    return max

```

Fig. 1. Algorithm MaxLyn

Recent work on the prefix table [4, 6] has confirmed its importance as a data structure for string algorithms. In this context it is interesting to find that Lyndon words \mathbf{x} can be characterized in terms of $\pi_{\mathbf{x}}$:

Observation 6 Suppose $\mathbf{x} = \mathbf{x}[1..n]$ is a string on alphabet Σ such that $\mathbf{x}[1]$ is the least letter in \mathbf{x} . Then \mathbf{x} is a Lyndon word over Σ if and only if for every $i \in 2..n$,

- (a) $i + \pi_{\mathbf{x}}[i] < n + 1$; and
- (b) for every $j \in i + 1..i + \pi_{\mathbf{x}}[i] - 1$, $j + \pi_{\mathbf{x}}[j] \leq i + \pi_{\mathbf{x}}[i]$.

In Appendix 1, the reader can find an additional result that justifies the strategy employed by **MaxLyn** (Figure 1).

3.2 Recursive Duval Factorization: Algorithm RDuval

Rather than independently computing the maximum-length Lyndon factor at each position i , as MaxLyn does, Algorithm RDuval recursively computes the Lyndon decomposition, [10], into maximum factors, at each step taking advantage of the fact that $\mathcal{L}[i]$ is known for the first position i in each factor, then re-computing with the first letters removed. This again gives immediate worst case

complexity of $O(n^2)$. We consider it only because it allows for a more refined discussion of the complexity in special cases for strings over binary alphabets giving an average case complexity of $O(n \log n)$, see below.

By Observation 1, whenever $\mathbf{x} = \mathbf{x}[1..n]$ is a Lyndon word, we know that $\mathcal{L}[1] = n$. Thus computing the Lyndon decomposition $\mathbf{x} = \mathbf{w}_1 \mathbf{w}_2 \cdots \mathbf{w}_k$, $\mathbf{w}_1 \geq \mathbf{w}_2 \geq \cdots \geq \mathbf{w}_k$, allows us to assign $\lambda[i_j] = |\mathbf{w}_j|$, where i_j is the first position of \mathbf{w}_j , $j = 1, 2, \dots, k$.

Algorithm RDUval applies this strategy recursively, by assigning $\lambda[i_j] \leftarrow |\mathbf{w}_j|$, then removing the first letter i_j from each \mathbf{w}_j to form \mathbf{w}'_j , to which the Lyndon decomposition is applied in the next recursive step. This process continues until each Lyndon word is reduced to a single letter.

The asymptotic time required for RDUval is bounded above by n times the maximum depth of the recursion, thus $O(n^2)$ in the worst case — consider, for example, the string $\mathbf{x} = a^{n-1}b$. However, to estimate expected behaviour, we can make use of a result of Bassino *et al.* [3]. Given a Lyndon word \mathbf{w} , they call $\mathbf{w} = \mathbf{u}\mathbf{v}$ the **standard factorization** of \mathbf{w} if \mathbf{u} and \mathbf{v} are both Lyndon words and \mathbf{v} is of maximum size. They then show that if \mathbf{w} is a binary string ($\Sigma = \{a, b\}$), the average length of \mathbf{v} is asymptotically $3|\mathbf{w}|/4$. Thus each recursive application of RDUval yields a left Lyndon factor of expected length $|\mathbf{w}|/4$ and a remainder of length $3|\mathbf{w}|/4$ to be factored further. It follows that the expected number of recursive calls of RDUval is $O(\log_{4/3} n)$. Hence

Lemma 7 *On binary strings RDUval executes in $O(n \log_{4/3} n)$ time on average.*

Example 8 *For*

$$\begin{array}{cccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ \mathbf{x} & = & a & a & b & a & a & b & b & a & b & b & a & b \\ \boldsymbol{\lambda} & = & 12 & 2 & 1 & 9 & 3 & 1 & 1 & 3 & 1 & 1 & 2 & 1 \end{array}$$

the factors considered are first 1–12, then

- 2–3 and 4–12 in the first level of recursion;
- 3, 5–7, 8–10 and 11–12 in the second level;
- 6, 7, 9, 10, 12 in the third level.

Positions are assigned as follows: $\lambda[1] \leftarrow 12$; $\lambda[2] \leftarrow 2$, $\lambda[4] \leftarrow 9$; $\lambda[3] \leftarrow 1$, $\lambda[5] \leftarrow 3$, $\lambda[8] \leftarrow 3$, $\lambda[11] \leftarrow 2$; $\lambda[6] \leftarrow 1$, $\lambda[7] \leftarrow 1$, $\lambda[9] \leftarrow 1$, $\lambda[10] \leftarrow 1$, $\lambda[12] \leftarrow 1$.

3.3 NSV Applied to the Inverse Suffix Array

The idea of the “next smaller value” (NSV) array for a given array \mathbf{x} had been proposed in various forms and under various names [1, 11, 12, 16].

Definition 9 (Next Smaller Value) *Given an array $\mathbf{x}[1..n]$ of ordered values, $NSV = NSV_{\mathbf{x}[1..n]}$ is the **next smaller value array** of \mathbf{x} if and only if for every $i \in 1..n$, $NSV[i] = j$, where*

- (a) *for every $h \in 1..j-1$, $\mathbf{x}[i] \leq \mathbf{x}[i+h]$; and*

(b) either $i+j = n+1$ or $\mathbf{x}[i] > \mathbf{x}[i+j]$.

Example 10

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \mathbf{x} & = & 3 & 8 & 7 & 10 & 2 & 1 & 4 & 9 & 6 & 5 \\ \text{NSV}_{\mathbf{x}} & = & 4 & 1 & 2 & 1 & 1 & 5 & 4 & 1 & 1 & 1 \end{array}$$

As shown in various contexts in [12], $\text{NSV}_{\mathbf{x}}$ can be computed in $\Theta(n)$ time using a stack. Our main observation here, also mentioned in [13], is that $\lambda_{\mathbf{x}}$ can be computed merely by applying NSV to the inverse suffix array $\text{ISA}_{\mathbf{x}}$. Proof of this claim can be found in Appendix 2; here we present the very simple $\Theta(n)$ -time, $\Theta(n)$ -space algorithm for this calculation:

```

procedure NSVISA( $\mathbf{x}[1..n]$ ) :  $\lambda_{\mathbf{x}}[1..n]$ 
  Compute  $\text{SA}_{\mathbf{x}}$  (see [15, 17])
  Compute  $\text{ISA}_{\mathbf{x}}$  from  $\text{SA}_{\mathbf{x}}$  in place (see [17])
   $\lambda_{\mathbf{x}} \leftarrow \text{NSV}(\text{ISA}_{\mathbf{x}})$  (in place)

```

Fig. 2. Apply NSV to $\text{ISA}_{\mathbf{x}}$

4 Elementary Computation of $\lambda_{\mathbf{x}}$ Using Ranges

In this section we describe an approach to the computation of $\lambda_{\mathbf{x}}$ that applies a variant of the NSV idea to the ranges of \mathbf{x} . Figure 3 gives pseudocode for Algorithm NSV* that uses the NSV stack ACTIVE to compute λ . The processing identifies ranges in a single left-to-right scan of \mathbf{x} , making use of two range comparison routines, COMP and MATCH. COMP compares adjacent individual ranges \mathbf{x}_r and \mathbf{x}_{r+1} , returning $\delta_1 = -1, 0, +1$ according as $\mathbf{x}_r < \mathbf{x}_{r+1}$, $\mathbf{x}_r = \mathbf{x}_{r+1}$, $\mathbf{x}_r > \mathbf{x}_{r+1}$. MATCH similarly returns δ_2 for adjacent *sequences* of ranges; that is,

$$\begin{aligned} \mathbf{X}_r &= \mathbf{x}_r \mathbf{x}_{r+1} \cdots \mathbf{x}_{r+s}, \text{ for some } s \geq 1; \\ \mathbf{X}_{r+s+1} &= \mathbf{x}_{r+s+1} \mathbf{x}_{r+s+2} \cdots \mathbf{x}_{r+s+t}, \text{ for some } t \geq 1. \end{aligned}$$

Algorithm NSV* is based on the idea encapsulated in Lemma 15 of Appendix 2, the main basis of the correctness of Algorithm NSVISA (see Figure 2). We process \mathbf{x} from left to right, using a stack ACTIVE initialized with index 1. At each iteration, the top of the stack (say, j) is compared with the current index (say, i). In particular, we need to compare $\mathbf{s}_{\mathbf{x}}(i)$ with $\mathbf{s}_{\mathbf{x}}(j)$, where $\mathbf{s}_{\mathbf{x}}(i) \equiv \mathbf{x}[i..n]$. As long as $\mathbf{s}_{\mathbf{x}}(i) \succeq \mathbf{s}_{\mathbf{x}}(j)$, NSV* pushes the current index and continues to the next. When $\mathbf{s}_{\mathbf{x}}(i) \prec \mathbf{s}_{\mathbf{x}}(j)$, it pops the stack and puts appropriate values in the corresponding indices of $\lambda_{\mathbf{x}}$. As noted above, especially Observations 1–3, ranges are employed to expedite these suffix comparisons.

```

procedure NSV* ( $\mathbf{x}, \boldsymbol{\lambda}$ )
  nextequal  $\leftarrow 0^n$ ; period  $\leftarrow 0^n$ 
  push(ACTIVE)  $\leftarrow 1$ 
   $\triangleright \mathbf{x}[n+1] = \$$ , a letter smaller than any in  $\Sigma$ .
  for  $i \leftarrow 2$  to  $n+1$  do
     $prev \leftarrow 0$ ;  $j \leftarrow$  peek(ACTIVE)
     $\triangleright$  COMP compares suffixes specified by  $i, j$  of two ranges.
     $\delta_1 \leftarrow$  COMP( $\mathbf{x}[j], \mathbf{x}[i]$ );  $\delta_2 \leftarrow 1$ 
    while ( $\delta_1 \geq 0$  and  $\delta_2 > 0$ ) do
      if  $\delta_1 = 0$  then  $\delta_2 \leftarrow$  MATCH( $\mathbf{x}[j], \mathbf{x}[i]$ )
      if  $\delta_2 > 0$  then
        if  $prev = 0$  or nextequal[ $j$ ]  $\neq prev$  then  $\boldsymbol{\lambda}[j] \leftarrow i-j$ 
        else
           $\boldsymbol{\lambda}[j] \leftarrow$  offset  $\leftarrow prev-j$ 
          if period[ $prev$ ] = 0 then
            if  $\boldsymbol{\lambda}[prev] > offset$  then
               $\boldsymbol{\lambda}[j] \leftarrow \boldsymbol{\lambda}[j] + \boldsymbol{\lambda}[prev]$ 
            else
              if nextequal[ $j$ ] =  $prev$  and offset  $\neq \boldsymbol{\lambda}[prev]$  then
                 $\boldsymbol{\lambda}[j] \leftarrow \boldsymbol{\lambda}[j] + period[prev]$ 
              if  $\boldsymbol{\lambda}[prev] = offset$  then
                 $\triangleright$  Current position is a part of periodic substring
                if period[ $prev$ ] = 0 then
                  period[ $j$ ]  $\leftarrow$  period[ $prev$ ] +  $2 \times offset$ 
                else
                  period[ $j$ ]  $\leftarrow$  period[ $prev$ ] + offset
          pop(ACTIVE)
           $prev \leftarrow j$ ;  $j \leftarrow$  peek(ACTIVE)
           $\triangleright$  Empty stack implies termination.
          if  $j = 0$  then EXIT
           $\delta_1 \leftarrow$  COMP( $\mathbf{x}[j], \mathbf{x}[i]$ )
     $\triangleright$  Finished processing  $i$  — it goes to stack.
    if  $\delta_2 = 0$  then nextequal[ $j$ ]  $\leftarrow i$ 
    push(ACTIVE)  $\leftarrow i$ 

```

Fig. 3. Computing $\boldsymbol{\lambda}_x$ using modified NSV

Two auxiliary arrays, **nextequal** and **period**, are required to handle situations in which MATCH finds that a suffix of a previous range at position j equals the current range at position i . Thus, when $\delta_2 = 0$, the algorithm assigns nextequal[j] $\leftarrow i$ before i is pushed onto ACTIVE. Then when a later MATCH yields $\delta_2 = 0$, the value of **period** — that is, the extent of the following periodicity — may need to be set or adjusted, as shown in the following example:

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
\mathbf{x}	=	a	a	a	b	a	a	b	a	a	b	a	a	b	a	b
nextequal	=	0	5	0	0	8	0	0	11	0	0	0	14	0	0	0
period	=	0	12	0	0	9	0	0	6	0	0	0	4	0	0	0

A straightforward implementation of COMP and MATCH could require a number of letter comparisons equal to the length of the shorter of the two sequences of ranges being matched. However, by performing $\Theta(n)$ -time preprocessing, we can compare two ranges in $\mathcal{O}(\sigma)$ time, where $\sigma = |\Sigma|$ is the alphabet size. Given $\Sigma = \{\mu_1, \mu_2, \dots, \mu_\sigma\}$, we define Parikh vectors $P_r[1..\sigma]$, where $P_r[j]$ is the number of occurrences of μ_j in range \mathbf{x}_r . Since ranges are monotone nondecreasing in the letters of the alphabet, it is easy to compute all the $P_r, r = 1, 2, \dots, m$, in linear time in a single scan of \mathbf{x} . Similarly, during the processing of each range \mathbf{x}_r , any value $P_{r,j}$, the Parikh vector of the suffix $\mathbf{x}_r[j..l_r]$, can be computed in constant time for each position considered. Thus we can determine the lexicographical order of any two ranges (or part ranges) \mathbf{x}_r and $\mathbf{x}_{r'}$ in $\mathcal{O}(\sigma)$ time rather than time $\mathcal{O}(\max(\ell_r, \ell_{r'}))$. The variant of NSV* that uses Parikh vectors is called PNSV*; otherwise NPNSV* for Not Parikh.

In Appendix 3 we describe briefly another approach to this suffix comparison problem, which we believe achieves run time $\mathcal{O}(n \log n)$ by maintaining a simple data structure requiring $\mathcal{O}(n \log n)$ space.

Now consider the worst case behaviour of Algorithm NSV*. Given the initial string $\mathbf{x}_0 = a^h b a^h c_0$, $h \geq 1$, $c_0 > b > a$, let $\mathbf{x}_k^{(h)} = \mathbf{x}_k = \mathbf{x}_{k-1} \mathbf{x}_{k-1}^*$, $k = 1, 2, \dots$, with \mathbf{x}_{k-1}^* identical to \mathbf{x}_{k-1} except in the last position, where the letter $c_k > c_{k-1}$ replaces c_{k-1} . Then \mathbf{x}_k has length $n = (h+1)m$, where $m = 2^{k+1}$ is the number of ranges in \mathbf{x}_k . We believe and are working towards a proof that \mathbf{x}_k is a worst-case input for Algorithm NSV*, which requires $\mathcal{O}(n \log n)$ range matches in such cases. Since PNSV* compares two ranges in $\mathcal{O}(\sigma)$ time, it therefore would require $\mathcal{O}(\sigma n \log n)$ time in the worst case, thus $\mathcal{O}(n \log n)$ for constant σ .

5 Preliminary Experimental Results

We have done some preliminary tests on the implementations of the two variants of NSV*, with and without employing Parikh vectors. The equipment used was an Intel(R) Core i3 at 1.8GHz and 4GB main memory under a 64-bit Windows 7 operating system. For each length 10000,20000...100000 we generated 500 random strings for alphabets of sizes $\sigma = 2, 4$ and 8. The results indicate, that at least for random strings, the processing time seems linear. The processing time for “with Parikh vectors” is greater because of the initial pre-processing. The data and the corresponding graphs are in Figures 8..11 below.

Processing time versus length of the string for the program without Parikh vectors

length of the input strings	10000	20000	30000	40000	50000	60000	70000	80000	90000	100000
alphabet size 2	0.001122	0.002265	0.003335	0.004444	0.005552	0.006682	0.007813	0.008936	0.010117	0.011167
alphabet size 4	0.001103	0.002233	0.003345	0.004464	0.0056	0.006677	0.007792	0.008903	0.010036	0.011157
alphabet size 8	0.001067	0.002112	0.003181	0.004255	0.005337	0.006372	0.007418	0.008522	0.009646	0.010653

Fig. 4. Processing times in seconds for the implementation without Parikh vectors

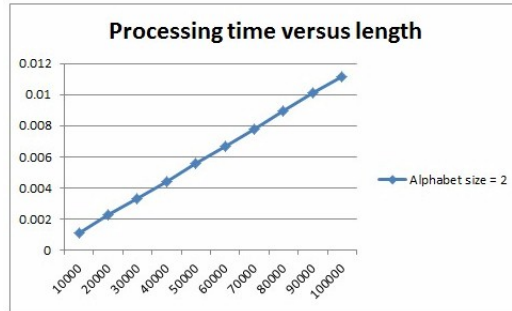


Fig. 5. Processing times for random strings over the binary alphabet; without Parikh vectors

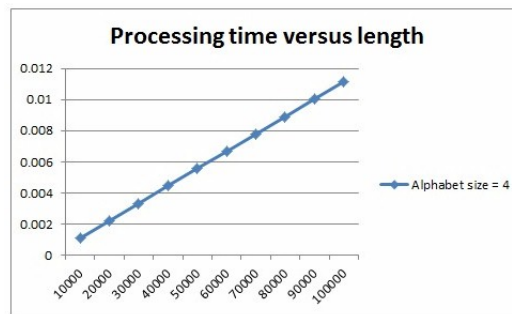


Fig. 6. Processing times for random strings over the alphabet of size 4; without Parikh vectors

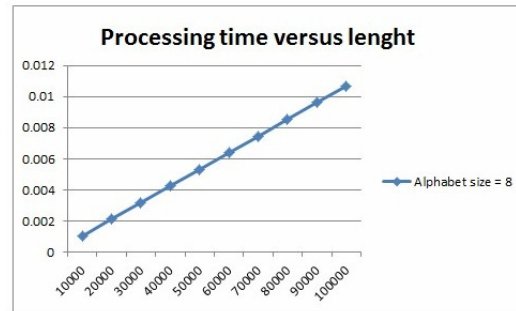


Fig. 7. Processing times for random strings over the alphabet of size 8; without Parikh vectors

Processing time versus length of the string for the program with Parikh vectors

length of the input strings	10000	20000	30000	40000	50000	60000	70000	80000	90000	100000
alphabet size 2	0.001368	0.002724	0.004072	0.005413	0.006782	0.008304	0.009465	0.010868	0.012201	0.013669
alphabet size 4	0.001898	0.003782	0.005694	0.007516	0.009498	0.011421	0.01319	0.015226	0.017124	0.019068
alphabet size 8	0.002385	0.004721	0.007053	0.009409	0.011874	0.014219	0.016546	0.018967	0.021337	0.024008

Fig. 8. Processing times in seconds for the implementation with Parikh vectors

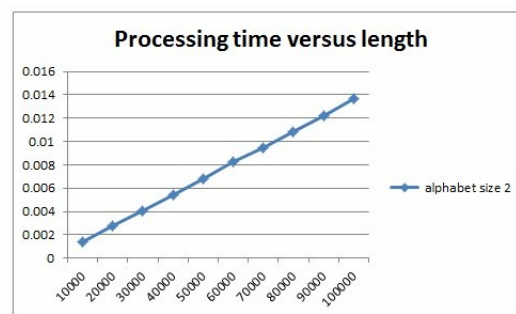


Fig. 9. Processing times for random strings over the binary alphabet; with Parikh vectors

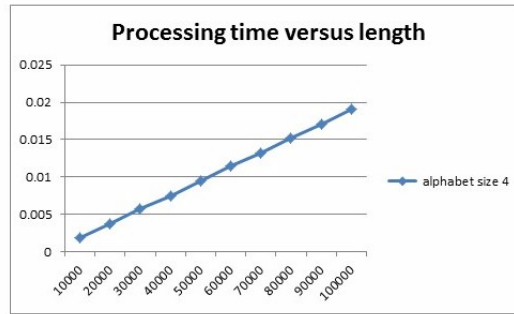


Fig. 10. Processing times for random strings over the alphabet of size 4; with Parikh vectors

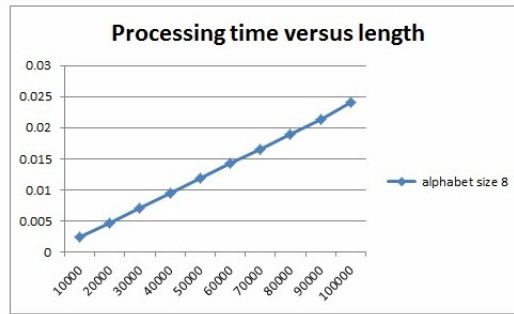


Fig. 11. Processing times for random strings over the alphabet of size 8; with Parikh vectors

6 Future Work

There is reason to believe [14] that the Lyndon array computation is less hard than suffix array construction. Thus the authors conjecture that there is a linear-time elementary algorithm (no suffix arrays) to compute the Lyndon array.

References

1. Stephen Alstrup, Cyril Gavoille, Haim Kaplan, and Theis Rauhe. Nearest common ancestors: a survey and new distributed algorithm. In *Proc. 11th Annual ACM Symp. on Parallel Algorithms & Architectures*, pages 258–264, 2002.
2. Hideo Bannai, Tomohiro I, Shunsuke Inenaga, Yuto Nakashima, Masayuki Takeda, and Kazuya Tsuruta. The “runs” theorem. 2014. *arXiv:1406.0263v6*.
3. Frédérique Bassino, Julien Clément, and Cyril Nicaud. The standard factorization of Lyndon words: an average point of view. *Discrete Mathematics*, 290(1):1–25, 2005.

4. Widmer Bland, Gregory Kucherov, and W. F. Smyth. Prefix table construction and conversion. *Proc. 24th Internat. Workshop on Combinatorial Algs. (IWOCOA)*, pages 41–53, 2013.
5. K. T. Chen, R. H. Fox, and R. C. Lyndon. Free differential calculus. iv. the quotient groups of the lower central series. *Annals of Mathematics*, 68(1):81–95, 1958.
6. Manolis Christodoulakis, P. J. Ryan, W. F. Smyth, and Shu Wang. Indeterminate strings, prefix arrays and undirected graphs. *Theoretical Comput. Sci.*, 600:34–48, 2015.
7. Maxime Crochemore, Christophe Hancart, and Thierry Lecroq. *Algorithms on Strings*. Cambridge University Press, New York, NY, USA, 2007.
8. Maxime Crochemore and Wojciech Rytter. *Jewels of stringology*. World Scientific, 2002.
9. Jean-Pierre Duval. Factorizing words over an ordered alphabet. *J. Algorithms*, 4:363–381, 1983.
10. Jean-Pierre Duval. Factorizing words over an ordered alphabet. *J. Algorithms*, 4(4):363–381, 1983.
11. Johannes Fischer, Veli Mäkinen, and Gonzalo Navarro. An(other) entropy-based compressed suffix tree. In *19th Annual Symp. on Combinatorial Pattern Matching*, volume 5029 of *Lecture Notes in Computer Science*, pages 152–165. Springer, 2008.
12. Keisuke Goto and Hideo Bannai. Simpler and faster Lempel-Ziv factorization. In *Data Compression Conference*, pages 133–142, 2013.
13. Christophe Hohlweg and Christophe Reutenauer. Lyndon words, permutations and trees. *Theor. Comput. Sci.*, 307(1):173–178, 2003.
14. Dmitry Kosolobov. Lempel-Ziv factorization may be harder than computing all runs. In *Proc. 32nd Symp. on Theoretical Aspects of Computer Science*, 2015. *arXiv:1409.5641*.
15. Ge Nong, Sen Zhang, and Wai H. Chan. Linear suffix array construction by almost pure induced-sorting. *Data Compression Conference*, 0:193–202, 2009.
16. Enno Ohlebusch and Simon Gog. Lempel-Ziv factorization revisited. In *22nd Annual Symp. on Combinatorial Pattern Matching*, volume 6661 of *Lecture Notes in Computer Science*, pages 15–26. Springer, 2011.
17. Simon J. Puglisi, W. F. Smyth, and Andrew H. Turpin. A taxonomy of suffix array construction algorithms. *ACM Comput. Surv.*, 39(2):1–31, July 2007.
18. Joe Sawada and Frank Ruskey. Generating Lyndon brackets: an addendum to “Fast algorithms to generate necklaces, unlabeled necklaces and irreducible polynomials over GF(2)”. *J. Algorithms*, 46:21–26, 2003.

Appendix 1

The following result justifies the strategy employed in Algorithm MaxLyn (Figure 1):

Lemma 11 *Suppose that for some position i in a Lyndon word $\mathbf{x}[1..n]$, $k = \pi[i] \geq 2$. Then for every $j \in i + 1 .. i + k - 1$, $\pi[j] \leq i + k - j$.*

Proof. The result certainly holds for $i + k = n + 1$, so we consider $i + k \leq n$. Assume that for some $j \in i + 1 .. i + k - 1$, $\pi[j] > i + k - j$. It follows that

$$\mathbf{x}[1 .. i + k - j + 1] = \mathbf{x}[j .. i + k], \quad (5)$$

while $\mathbf{x}[j-i+1..k] = \mathbf{x}[j..i+k-1]$. Since \mathbf{x} is Lyndon, therefore $\mathbf{x}[1+k] \prec \mathbf{x}[i+k]$, and so we find that

$$\mathbf{x}[j-i+1..1+k] \prec \mathbf{x}[j..i+k]. \quad (6)$$

From (5) and (6) we see that $\mathbf{x}[1..k+1]$ has suffix $\mathbf{x}[j-i+1..k+1]$ satisfying $\mathbf{x}[j-i+1..k+1] \prec \mathbf{x}[1..i+k-j+1]$, contradicting the assumption that \mathbf{x} is Lyndon. \square

Appendix 2

Here we prove Theorem 12 that justifies the algorithm given in Figure 2:

Theorem 12 *For a given string $\mathbf{x} = \mathbf{x}[1..n]$ on alphabet Σ , totally order by \prec , let $ISA = ISA_{\mathbf{x}}^{\prec}$. Then for every $i \in 1..n$, the substring $\mathbf{x}[i..j]$ is a longest Lyndon factor with respect to \prec if and only if*

- (a) *for every $h \in i+1..j$, $ISA[j] < ISA[h]$; and*
- (b) *either $j = n$ or $ISA[j+1] < ISA[i]$.*

The following well-known result is needed to prove Lemma 14:

Lemma 13 (Duval, Lemma 1.6, [9]) *Suppose $\mathbf{x} \in \Sigma^+$, where Σ is an alphabet totally ordered by \prec . Let $\mathbf{x} = \mathbf{u}^r \mathbf{u}_1 b$, where \mathbf{u} is nonempty, $r \geq 1$, \mathbf{u}_1 a possibly empty proper prefix of \mathbf{u} , and the letter $b \neq \mathbf{u}[|\mathbf{u}_1|+1]$.*

- (a) *If $b \prec \mathbf{u}[|\mathbf{u}_1|+1]$, then \mathbf{u} is a longest Lyndon prefix of $\mathbf{x}\mathbf{y}$ for any \mathbf{y} ;*
- (b) *if $b \succ \mathbf{u}[|\mathbf{u}_1|+1]$, then \mathbf{x} is Lyndon with respect to \prec .*

For a given string $\mathbf{x}[1..n]$, let $\mathbf{s}\mathbf{x}(i) = \mathbf{x}[i..n]$ denote the suffix of \mathbf{x} beginning at position i . When clear from context we write just $\mathbf{s}(i)$.

Lemma 14 *Consider a string $\mathbf{x} = \mathbf{x}[1..n]$ over alphabet Σ totally ordered by \prec . Let $\mathbf{x}[i..j]$ be the longest Lyndon factor of \mathbf{x} starting at i . Then $\mathbf{s}\mathbf{x}(i) \prec \mathbf{s}\mathbf{x}(k)$ for every $k \in i+1..j$ and either $j = n$ or $\mathbf{s}\mathbf{x}(j+1) \prec \mathbf{s}\mathbf{x}(i)$.*

Proof. Because $\mathbf{x}[i..j]$ is Lyndon, therefore for any $i < k \leq j$, $\mathbf{x}[i..j] \prec \mathbf{x}[k..j]$ and so $\mathbf{s}(i) \prec \mathbf{s}(k)$. If $j = n$, we are done. So we may assume $j < n$, and we want to show that $\mathbf{s}(j+1) \prec \mathbf{s}(i)$. Suppose then that $\mathbf{s}(j+1) \not\prec \mathbf{s}(i)$. Since $\mathbf{s}(i)$ and $\mathbf{s}(j+1)$ are distinct, it follows that $\mathbf{s}(i) \prec \mathbf{s}(j+1)$. If we let $d = \text{lcp}(\mathbf{s}(i), \mathbf{s}(j+1)) + 1$, two cases arise:

- (a) $0 \leq d \leq j - i$.

Here $i \leq i + d \leq j$. Thus $\mathbf{x}[i..i+d-1] = \mathbf{x}[j+1..j+d]$ and $\mathbf{x}[i+d] \prec \mathbf{x}[j+1+d]$, and so for $j < k \leq j+1+d$, $\mathbf{x}[i..j+1+d] \prec \mathbf{x}[k..j+1+d]$. Since $\mathbf{x}[i..j]$ is Lyndon, $\mathbf{x}[i..j] \prec \mathbf{x}[k..j]$ and so $\mathbf{x}[i..j+1+d] \prec \mathbf{x}[k..j+1+d]$ for any $i < k \leq j$. Thus $\mathbf{x}[i..j+1+d]$ is Lyndon, contradicting the assumption that $\mathbf{x}[i..j]$ is the longest Lyndon factor starting at i .

(b) $0 < j - i \leq d$.

Let $d = r(j - i) + d_1$, where $0 \leq d_1 < j - i$. Then $r \geq 1$ and $\mathbf{x}[i..j+1+d] = \mathbf{u}^r \mathbf{u}_1 b$ where $\mathbf{u} = \mathbf{x}[i..j]$,

$$\mathbf{u}_1 = \mathbf{x}[j+r(j-i)+1..j+r(j-i)+d_1-1] = \mathbf{x}[j+r(j-i)+1..j+d-1]$$

is a prefix of $\mathbf{x}[i..j]$, and $\mathbf{x}[i+d] \prec \mathbf{x}[j+1+d]$, so that by Lemma 13 (b), $\mathbf{x}[i..j+1+d]$ is Lyndon, contradicting the assumption that $\mathbf{x}[i..j]$ is the longest Lyndon factor starting at i .

Thus $\mathbf{s}(j+1) \prec \mathbf{s}(i)$, as required. \square

Lemma 15 describes the property of being a longest Lyndon factor of a string \mathbf{x} in terms of relationships between corresponding suffixes.

Lemma 15 *Consider a string $\mathbf{x} = \mathbf{x}[1..n]$ over an alphabet Σ with an ordering \prec . A substring $\mathbf{x}[i..j]$ is a longest Lyndon factor of \mathbf{x} with respect to \prec if and only if $\mathbf{s}\mathbf{x}(i) \prec \mathbf{s}\mathbf{x}(k)$ for every $k \in i+1..j$ and either $j = n$ or $\mathbf{s}\mathbf{x}(j+1) \prec \mathbf{s}\mathbf{x}(i)$.*

Proof. Let (A) denote $\{\mathbf{x}[i..j] \text{ is a longest Lyndon factor of } \mathbf{x}\}$ and let (B) denote $\{\mathbf{s}(i) \prec \mathbf{s}(k) \text{ for any } 1 \leq k \leq j \text{ and } \mathbf{s}(j+1) \prec \mathbf{s}(i)\}$. Then (A) \Rightarrow (B) follows from Lemma 14, so we need to prove that (B) \Rightarrow (A).

Suppose then that (B) holds, and let $\mathbf{x}[i..k]$ be a longest Lyndon factor of \mathbf{x} starting at position i . If $k < j$, then by Lemma 14, $\mathbf{s}(k+1) \prec \mathbf{s}(i)$, a contradiction since $k+1 \leq j$. If $k > j$, then by Lemma 14, $\mathbf{s}(i) \prec \mathbf{s}(j+1)$ because $j+1 \leq k$, which again gives us a contradiction. Thus $k = j$ and $\mathbf{x}[i..j]$ is a longest Lyndon factor of \mathbf{x} . \square

Now we reformulate Lemma 15 in terms of the inverse suffix array ISA of \mathbf{x} using the relationship that $\mathbf{s}(i) \prec \mathbf{s}(j) \iff \text{ISA}[i] < \text{ISA}[j]$, thus yielding Theorem 12, as required. Hence the Lyndon array can be computed in a simple three-step algorithm, as shown in Figure 2, that executes in $\theta(n)$ time and uses only one additional array of integers.

Appendix 3

Here we describe a simple data structure that yields an alternative approach to Algorithm NSV*, based on the comparison of longest Lyndon factors as described in Lemma 15. The **dictionary of basic factors** [7, 8] of string $\mathbf{x}[1..n]$ consists of a sequence of arrays \mathcal{D}_t , $0 \leq t \leq \log n$. The array \mathcal{D}_t records information about factors of \mathbf{x} of length 2^t — that is, the basic factors. In particular, $\mathcal{D}_t[i]$ stores the rank of $\mathbf{x}[i..i+2^t-1]$, so that

$$\mathbf{x}[i..i+2^t-1] \preceq \mathbf{x}[i..i+2^t-1] \iff \mathcal{D}_t[i] \leq \mathcal{D}_t[i].$$

This dictionary requires $O(n \log n)$ space and can be constructed in $O(n \log n)$ time as follows. \mathcal{D}_0 contains information about consecutive symbols of \mathbf{x} and

hence can be computed in $O(n \log n)$ time by sorting all the symbols appearing in \mathbf{x} and mapping them to numbers from 1 and onward. Once \mathcal{D}_t is computed, we can easily compute \mathcal{D}_{t+1} by spending $O(n)$ time on a radix sort, because $u[i..i + 2^{t+1} - 1]$ is in fact a concatenation of the factors $u[i..i + 2^t - 1]$ and $u[i + 2^t..i + 2^{t+1} - 1]$.

Once this dictionary is computed, we can compare any two factors by comparing two appropriate overlapping basic factors (i.e., factors having length power of two), which is done by checking the corresponding \mathcal{D} array from the dictionary. This will require constant time and hence each suffix-suffix comparison can be done in constant time.