



RESEARCH REPOSITORY

This is the author's final version of the work, as accepted for publication following peer review but without the publisher's layout or pagination. The definitive version is available at:

<https://doi.org/10.1016/j.tcs.2017.04.008>

Daykin, J.W., Franěk, F., Holub, J., Islam, A.S.M.S. and Smyth, W.F. (2017) Reconstructing a string from its Lyndon arrays. Theoretical Computer Science.

<http://researchrepository.murdoch.edu.au/id/eprint/37033/>

Copyright: © 2017 Elsevier B.V.
It is posted here for your personal use. No further distribution is permitted.

Accepted Manuscript

Reconstructing a string from its Lyndon arrays

Frantisek Franek, Jacqueline W. Daykin, Jan Holub, A.S.M. Sohidull Islam, W.F. Smyth

PII: S0304-3975(17)30330-4
DOI: <http://dx.doi.org/10.1016/j.tcs.2017.04.008>
Reference: TCS 11161

To appear in: *Theoretical Computer Science*

Received date: 13 June 2016
Revised date: 9 April 2017
Accepted date: 17 April 2017

Please cite this article in press as: F. Franek et al., Reconstructing a string from its Lyndon arrays, *Theoret. Comput. Sci.* (2017), <http://dx.doi.org/10.1016/j.tcs.2017.04.008>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



Highlights

- An algorithm to determine whether or not a given array is a Lyndon array.
- Two algorithms to compute a string from a valid Lyndon array.
- An algorithm to compute a unique string from “rotated” Lyndon arrays.
- Plus theoretical background.

Reconstructing a String from its Lyndon Arrays^{*}

Frantisek Franek^{**1,4}, Jacqueline W. Daykin^{2,5}, Jan Holub³,
A. S. M. Sohidull Islam⁴, and W. F. Smyth^{**1,4,5,6}

¹ Algorithms Research Group, Department of Computing & Software
McMaster University

{franek,smyth}@mcmaster.ca

² Department of Computer Science
Aberystwyth University (Mauritius Branch Campus)
jackie.daykin@gmail.com

³ Department of Theoretical Computer Science
Czech Technical University in Prague

Jan.Holub@fit.cvut.cz

⁴ School of Computational Science & Engineering
McMaster University
sohansayed@gmail.com

⁵ Department of Informatics, King's College London

⁶ School of Engineering & Information Technology
Murdoch University

Abstract. Given a string $\mathbf{x} = \mathbf{x}[1..n]$ on an ordered alphabet Σ of size σ , the Lyndon array $\lambda = \lambda_{\mathbf{x}}[1..n]$ of \mathbf{x} is an array of positive integers such that $\lambda[i]$, $1 \leq i \leq n$, is the length of the maximal Lyndon word over the ordering of Σ that begins at position i in \mathbf{x} . The Lyndon array has recently attracted considerable attention due to its pivotal role in establishing the long-standing conjecture that $\rho(n) < n$, where $\rho(n)$ is the maximum number of maximal periodicities (runs) in any string of length n . Here we first describe two linear-time algorithms that, given a valid Lyndon array λ , compute a corresponding string — one for an alphabet of size n , the other for a smaller alphabet. We go on to describe another linear-time algorithm that determines whether or not a given integer array is a Lyndon array of some string. Finally we show how σ Lyndon arrays $\lambda_{\Sigma} = \{\lambda_1 = \lambda, \lambda_2, \dots, \lambda_{\sigma}\}$ corresponding to σ “rotations” of the alphabet can be used to determine uniquely the string \mathbf{x} on Σ such that $\lambda_{\mathbf{x}} = \lambda$.

1 Introduction

We suppose throughout that \mathbf{x} is a nonempty string. If $\mathbf{x} = \mathbf{uv}$ for some \mathbf{u} and nonempty \mathbf{v} , then \mathbf{vu} is said to be the $|\mathbf{u}|^{\text{th}}$ *rotation* of \mathbf{x} , written

^{*} The authors express their thanks to two anonymous referees whose comments have materially improved the quality of this paper.

^{**} Supported in part by grants from the Natural Sciences & Engineering Research Council (NSERC) of Canada.

$vu = R_{|u|}(x)$. A string x is said to be a *repetition* if $x = u^e$ for some nonempty string u and some integer $e \geq 2$; otherwise, x is said to be *primitive*. A primitive string x that is lexicographically least among all its rotations $R_k(x)$, $k = 0, 1, \dots, |x|-1$, is said to be a *Lyndon word* [26]. As a consequence of their interesting properties, Lyndon words have been much studied: the existence of a unique factorization $x = w_1 w_2 \cdots w_s$ of a string into Lyndon words $w_1 \geq w_2 \geq \cdots \geq w_s$ was demonstrated some 60 years ago [7] and a simple linear-time algorithm to compute the *Lyndon factorization* was proposed a quarter-century later [18].

In fact, Lyndon words are a special case of Unique Maximal Factorization Families (UMFFs), that over the last 15 years have also been studied extensively [16, 17, 13, 12]. When every factor w_j , $1 \leq j \leq s$, of a (not necessarily Lyndon) factorization of x belongs to a specified set \mathcal{W} , we say that it is a *factorization of x over \mathcal{W}* , denoted by $F_{\mathcal{W}}(x)$. Then a subset $\mathcal{W} \subseteq \Sigma^+$ is a *factorization family* (FF) if and only if for every nonempty string x on Σ there exists a factorization $F_{\mathcal{W}}(x)$. If \mathcal{W} is an FF on an alphabet Σ , then \mathcal{W} is said to be a *unique maximal factorization family* (UMFF) if and only if there exists a unique factorization $F_{\mathcal{W}}(x)$ for every string $x \in \Sigma^+$. We expect that the results given here for Lyndon arrays can be generalized to UMFFs.

The *Lyndon array* $\lambda = \lambda_x[1..n]$ (equivalently, $L = L_x[1..n]$) of a given $x = x[1..n]$ gives at each position i the length (equivalently, the end position) of the longest (or *maximal*) Lyndon word starting at i . Thus $L_x[i] = \lambda_x[i] + i - 1$. Apparently first introduced as “Lyndon bracketing” [25], the Lyndon array has recently become of interest because of the central role it plays in the surprising and simple proof [6] that the maximum number $\rho(n)$ of maximal periodicities (runs) in any string of length n satisfies $\rho(n) < n$. A recent paper [22] studies algorithms to compute λ_x , exhibiting several that apparently execute in linear expected time, while conjecturing that there exists a worst-case linear-time algorithm to compute λ_x that is “elementary” — not a precise term, but we intend by it an algorithm that computes local features of a string while avoiding prior computation of global data structures such as the suffix array. Indeed, such an algorithm has recently been found [4, 5] as a first step in a two-step non-recursive linear-time suffix array construction algorithm. Here is an example of a Lyndon array, taken from [22]:

$$\begin{array}{rcccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\
 x & = & a & b & a & a & b & a & b & a & a & b \\
 \lambda_x & = & 2 & 1 & 5 & 2 & 1 & 2 & 1 & 3 & 2 & 1 \\
 L_x & = & 2 & 2 & 7 & 5 & 5 & 7 & 7 & 10 & 10 & 10
 \end{array} \tag{1}$$

Since λ and L are arrays of positive integers, it is natural to ask under what conditions a given integer array is a Lyndon array. In Section 2 we give necessary and sufficient conditions that a given integer array L^* is a Lyndon array $L_{\mathbf{x}}$ of some string \mathbf{x} on some alphabet Σ . We then describe linear-time algorithms that compute a string \mathbf{x} corresponding to a given Lyndon array L^* — the problem of computing a lexicographically least such string on a minimum-size alphabet appears to be computationally difficult. Finally we describe a linear-time algorithm to determine whether or not a given integer array is a Lyndon array of some string.

In Section 3 we go on to establish a “reverse engineering” result for Lyndon arrays; that is, given certain Lyndon arrays $L_{\mathbf{x}}$ based on orderings of a given alphabet Σ of size σ , what can be said about the corresponding string \mathbf{x} ? This kind of problem was first introduced in [21, 19] for the border array, then later considered for various common string data structures; for example, prefix tables [9, 3, 8], KMP arrays [20, 23, 24], cover arrays [10], and many others. Section 3 also presents an $\mathcal{O}(\sigma n)$ -time algorithm to compute the unique string \mathbf{x} determined by the Lyndon arrays computed for σ rotations of the alphabet. In Section 4 we discuss a variety of open problems arising.

2 When is L^* a Valid Lyndon Array of Some String?

We begin with a result from [22]:

Observation 1 *Suppose positions i, j in $\mathbf{x}[1..n]$ satisfy $1 \leq i < j \leq n$. Then either $L_{\mathbf{x}}[i] < j$ or $L_{\mathbf{x}}[i] \geq L_{\mathbf{x}}[j]$; that is, the vectors $(i, L_{\mathbf{x}}[i])$ and $(j, L_{\mathbf{x}}[j])$ are nonintersecting, and $L_{\mathbf{x}}[i] = j \implies j = L_{\mathbf{x}}[j]$.*

(We remark that a generalization of this result follows also from the **xyz** Lemma of [16] which states that an FF \mathcal{W} is an UMFF if and only if whenever $\mathbf{xy}, \mathbf{yz} \in \mathcal{W}$ for some nonempty \mathbf{y} , then $\mathbf{xyz} \in \mathcal{W}$. If we let \mathbf{xy} be the Lyndon word at $\mathbf{x}[i..L_{\mathbf{x}}[i]]$ and \mathbf{yz} the Lyndon word at $\mathbf{x}[j..L_{\mathbf{x}}[j]]$, it follows that the case $i < j < L_{\mathbf{x}}[i] < L_{\mathbf{x}}[j]$ cannot arise.)

Observation 1 tells us that a nonintersecting, or Monge-like, property necessarily holds for the arcs $(i, L[i])$ determined by the Lyndon array $L_{\mathbf{x}}$ of every string \mathbf{x} . To see that this property is also sufficient, consider an integer array $L^*[1..n]$ in which $i \leq L^*[i] \leq n$ for every $i \in 1..n$, and where either $L^*[i] < j$ or $L^*[i] \geq L^*[j]$ for every $1 \leq i < j \leq n$. Suppose an alphabet $\Sigma = \{\mu_1, \mu_2, \dots, \mu_n\}$ is given, with $\mu_1 < \mu_2 < \dots < \mu_n$. We now outline an algorithm (see Figure 1) that assigns the n letters of Σ to positions in \mathbf{x} in such a way that $L_{\mathbf{x}} = L^*$.

```

procedure SimpleAssign ( $L^*, n, \Sigma, \sigma; \mathbf{x}$ )
  Radix sort pairs  $(L^*[i], i)$ ,  $1 \leq i \leq n$ , in ascending
    order of  $i$  within descending order of  $L^*[i]$ 
    to form sorted positions  $I = I[1..n]$ .
  for  $j \leftarrow 1$  to  $n$  do
     $\mathbf{x}[I[j]] \leftarrow j$ 

```

Fig. 1. Given a valid Lyndon array L^* and an ordered alphabet $\Sigma = \{1, 2, \dots, n\}$, in $\mathcal{O}(n)$ time construct a string \mathbf{x} on Σ whose Lyndon array is L^* .

This algorithm ensures that arcs $(i, L^*[i])$ are processed in descending order of $L^*[i]$ — specifically, so that for all i_1, i_2, \dots, i_m such that

$$L^*[i_1] = L^*[i_2] = \dots = L^*[i_m], \quad i_1 < i_2 < \dots < i_m,$$

it follows that $\mathbf{x}[i_1] < \mathbf{x}[i_2] < \dots < \mathbf{x}[i_m]$. Thus for each choice of $L^*[i]$, we ensure that $L^*[i] = L_{\mathbf{x}}[i]$. The descending order, together with the nonintersecting property, then guarantees that this identity holds for all i , and so the nonintersecting property is sufficient to ensure that L^* is the Lyndon array of some string — in particular \mathbf{x} . Thus:

Lemma 2 *Suppose that $L^*[1..n]$ is an integer array such that $1 \leq L^*[i] \leq n$ for all $i \in 1..n$. Then L^* is a Lyndon array $L_{\mathbf{x}}$ of some string \mathbf{x} if and only if for all i, j such that $1 \leq i < j \leq n$, either $L^*[i] < j$ or $L^*[i] \geq L^*[j]$.*

For the example (1), Algorithm SimpleAssign yields the following:

$$\begin{array}{rcccccccccc}
 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\
 L^* & = & 2 & 2 & 7 & 5 & 5 & 7 & 7 & 10 & 10 & 10 \\
 I & = & 8 & 9 & 10 & 3 & 6 & 7 & 4 & 5 & 1 & 2 \\
 \mathbf{x} & = & 9 & 10 & 4 & 7 & 8 & 5 & 6 & 1 & 2 & 3
 \end{array} \tag{2}$$

In an effort to construct a string \mathbf{x} on a smaller alphabet, we employ a strategy (see Figure 2) that for each range of increasing values in I (that is, for each maximal Lyndon word of length at least two):

- chooses an initial letter one greater than the initial letter in the immediately following maximal Lyndon word;
- assigns the same letter to consecutive positions at the beginning of the current maximal Lyndon word — but excluding the final position;
- thereafter increments the letter by one at each successive position.

For example, in (2), after selecting $\mathbf{x}[8..10] = 112$, we choose

$$\mathbf{x}[3] = 2, \mathbf{x}[6] = 3, \mathbf{x}[7] = 4,$$

corresponding to $I[4..6] = 367$ and ensuring that $L\mathbf{x}[5] \leq 5$. Then, for $I[7..8] = 45$, we choose

$$\mathbf{x}[4] = 4, \mathbf{x}[5] = 5,$$

finally yielding $\mathbf{x} = 3424534112$, on an alphabet of size 5 rather than 10, but still far from the minimum of 2 ($\mathbf{x} = 1211212112$). Clearly Algorithm BetterAssign also executes in $\mathcal{O}(n)$ time; it yields the same worst-case result as SimpleAssign (when $I = n, n-1, \dots, 1$), but otherwise finds a string \mathbf{x} on a smaller alphabet.

```

procedure BetterAssign ( $L^*, n, \Sigma, \sigma; \mathbf{x}$ )
  Compute  $I[1..n]$  as in SimpleAssign
   $I[n+1] \leftarrow 0; h \leftarrow 1; i \leftarrow 1$ 
  while  $i \leq n$  do
     $\triangleright$  Assign letters to range of increasing values from  $I[i]$ .
    repeat
       $\triangleright$  Consecutive positions at start of range are identical.
       $\mathbf{x}[I[i]] \leftarrow h; i \leftarrow i+1$ 
    until  $I[i+1] \neq I[i]+1$ 
    if  $I[i+1] < I[i]$  then
       $\triangleright$  End position in range must be incremented.
       $\mathbf{x}[I[i]] \leftarrow h+1; i \leftarrow i+1$ 
    else
      while  $I[i+1] > I[i]$  do
         $\triangleright$  Elsewhere in range every position is incremented.
         $h \leftarrow h+1; \mathbf{x}[I[i]] \leftarrow h; i \leftarrow i+1$ 
     $\triangleright$  Reset  $h$  depending on the next range in  $\mathbf{x}$ .
    if  $i \leq n$  then
       $j \leftarrow i$ 
      while  $I[j+1] > I[j]$  do  $j \leftarrow j+1$ 
       $h \leftarrow \mathbf{x}[I[j]+1]+1$ 

```

Fig. 2. Construct a string \mathbf{x} on a subset of $\Sigma = \{1, 2, \dots, n\}$ with Lyndon array L^* .

We know of no approach other than brute force (trial and error) to the computation of \mathbf{x} on a minimum alphabet consistent with L^* . Hence


```

function CheckLyndon ( $L^*, n$ )
STACK  $\leftarrow \emptyset$ 
for  $i \leftarrow 1$  to  $n$  do
  if  $L^*[i] < i$  or  $L^*[i] > n$  return (FALSE)
  if  $L^*[i] > i$  then
     $end \leftarrow 0$ 
    while STACK  $\neq \emptyset$  and  $end < L^*[i]$  do
       $end \leftarrow peek(\text{STACK})$ 
      if  $end < i$  then  $pop(\text{STACK})$ 
      elseif  $end < L^*[i]$  then return (FALSE)
     $push(\text{STACK}, L^*[i])$ 
return (TRUE)

```

Fig. 3. Determine whether (TRUE) or not (FALSE) a given integer array L^* is a Lyndon array of some string.

Problem 3 Given a valid Lyndon array L^* , what is the complexity of the problem of constructing a string \mathbf{x} on a minimum alphabet consistent with L^* ?

We turn now to the problem of determining whether or not a given integer array L^* is valid; that is, whether or not it is a Lyndon array of some string. To solve this problem we introduce Algorithm CheckLyndon (see Figure 3), based on Lemma 2. It processes the segments $(i, L^*[i])$ in ascending order of position i and places the “end” of each nontrivial segment on the stack. Before doing so, it checks to see if any previous end lies within the current segment: if so, L^* cannot be a Lyndon array. If not, then either the previous entry ended before the current range and so can be deleted from the stack, or else it includes the current range and so must be kept in the stack to be tested against later segments. Note that entries in the stack have all been tested against *preceding* segments. We claim therefore that CheckLyndon is correct.

To see that the algorithm executes in linear time, observe that segment i is either wholly contained in a preceding segment, so that access to the stack is terminated, or else the current stack entry is deleted. Thus the total time requirement of the **while** loop is $\mathcal{O}(n)$. We have:

Lemma 4 Algorithm CheckLyndon correctly determines in $\mathcal{O}(n)$ time whether or not a given integer array $L^*[1..n]$ is a Lyndon array.

3 Reconstructing a String from its Lyndon Arrays

Suppose an alphabet $\Sigma = \{\ell_1, \ell_2, \dots, \ell_\sigma\}$, $\sigma \geq 2$, is given with initial global order $R_1: \ell_1 < \ell_2 < \dots < \ell_\sigma$. For $j = 2, 3, \dots, \sigma$, the j^{th} rotation R_j of R_1 is the order $\ell_j < \ell_{j+1} < \dots < \ell_\sigma < \ell_1 < \dots < \ell_{j-1}$. Thus ℓ_j is the least letter, and for $j > 1$ ℓ_{j-1} is the largest, in the rotation R_j . The collection of σ rotations is denoted by R_Σ .

In this section we deal with the problem of identifying a unique string on alphabet Σ corresponding to R_Σ . We begin with an observation from [22], that we can write \mathbf{x} in the form $\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_m$, where for each $r \in 1..m$, $|\mathbf{x}_r| = \text{len}_r$ and

$$\mathbf{x}_r[1] \leq \mathbf{x}_r[2] \leq \dots \leq \mathbf{x}_r[\text{len}_r], \quad (3)$$

while for $1 \leq r < m$,

$$\mathbf{x}_r[\text{len}_r] > \mathbf{x}_{r+1}[1]. \quad (4)$$

We call \mathbf{x}_r a *range* in \mathbf{x} , and we identify a position j in range \mathbf{x}_r , $1 \leq j \leq \text{len}_r$, with its equivalent position i in \mathbf{x} by writing $i = S_{r,j} = \sum_{r'=1}^{r-1} \text{len}_{r'} + j$. Then, again from [22], we have the following:

Observation 5 *Let $i = S_{r,j}$ be a position in \mathbf{x} that corresponds to position j in range \mathbf{x}_r .*

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

Based on these remarks, for the special case $\sigma = 2$, we can now prove:

Lemma 6 *Let $L_{\mathbf{x}}$ be the Lyndon array of a string $\mathbf{x}[1..n]$ on $\Sigma = \{a, b\}$, $a < b$. Then, provided that $\lambda_{\mathbf{x}} \neq 1^n$, \mathbf{x} is determined uniquely by $L_{\mathbf{x}}$.*

Proof. First observe that $\lambda_{\mathbf{x}} = 1^n$ if and only if $\mathbf{x} = b^m a^{n-m}$ for some $m \in 0..n$ — thus in this case the corresponding $L_{\mathbf{x}} = 12 \dots n$ corresponds to $n+1$ choices for \mathbf{x} . Otherwise, let n' be the smallest index such that for every $i \in n'..n$, there exists no $j < i$ such that $L_{\mathbf{x}}[j] = i$. If there is no such n' , then $L_{\mathbf{x}}[j] = n$ for some $j < n$; in this case, set $n' = n+1$. Then for every $i \in n'..n$, $L_{\mathbf{x}}[i] = i$, which, since $\mathbf{x} \neq b^n$ by hypothesis, implies that $\mathbf{x}[i] = a$; that is, $\mathbf{x}[n'..n] = a^{n-n'+1}$. Since, again by hypothesis, $\mathbf{x} \neq a^n$, it follows that $n' > 2$ and $\mathbf{x}[n'-1] = b$. More generally, by Observation 5, for every $i < n'$, $\mathbf{x}[i] = b$ if and only if $L_{\mathbf{x}}[i] = i$. Therefore, for all other i , $\mathbf{x}[i] = a$. Thus \mathbf{x} is uniquely determined by $L_{\mathbf{x}}$, as required. \square

Suppose then that $\sigma \geq 3$. In this case, corresponding to R_Σ , we assume that Lyndon arrays $\lambda_\Sigma = \{\lambda_1, \lambda_2, \dots, \lambda_\sigma\}$ are given, where λ_j , $1 \leq j \leq \sigma$, is based on rotation R_j of the alphabet and determined by some (unknown) string $\mathbf{x} = \mathbf{x}[1..n]$. Then for $j \in 1..\sigma$, $i \in 1..n$, $\lambda_j[i]$ is the length of the maximal Lyndon word at position i in \mathbf{x} based on rotation R_j of the alphabet. We say that position i in \mathbf{x} , $1 \leq i \leq n$, is **covered** by λ_j if and only if there exists a position $i' < i$ such that $i' + \lambda_j[i'] > i$ (alternatively, $L_j[i'] \geq i$).

We now prove the following result, enabling us to uniquely determine \mathbf{x} from R_Σ :

Theorem 7 *Suppose that \mathbf{x} is a string on an alphabet Σ of size $\sigma \geq 3$, whose Lyndon arrays λ_Σ are given based on rotations R_Σ . Let*

$$\lambda^+[i] = \max(\lambda_1[i], \lambda_2[i], \dots, \lambda_\sigma[i]), \quad 1 \leq i \leq n,$$

and let $P[i]$ be the nonempty ascending sequence $\{j_1, j_2, \dots, j_k\}$ of indices j that specify all the rotations R_j for which $\lambda^+[i] = \lambda_j[i]$; that is, that yield the maximum Lyndon word at position i . Then

1. $|P[i]| = 1 \implies \mathbf{x}[i] = \ell_{j_1}$.
2. $1 < |P[i]| < \sigma \implies \mathbf{x}[i] = \ell_{j_h}$, where j_h is the unique (leftmost) entry in $P[i]$ such that $j_{(h+1) \bmod \sigma} \notin P[i]$.
3. $|P[i]| = \sigma \iff \lambda^+[i] = 1$, and $\mathbf{x}[i..n] = \ell_{j_h}^{n-i+1}$, where i is the least integer such that $\lambda^+[i] = 1$ and j_h is the unique entry in $P[i]$ such that $\mathbf{x}[i]$ is not covered by λ_{j_h} .

We remark that, since the assignments to positions in \mathbf{x} made here under Cases 1–3 are unique, therefore \mathbf{x} must be the only string on Σ that satisfies the constraints given by λ_Σ . In Figure 4 the various cases of Theorem 7 are illustrated:

$$\begin{array}{l} i = 1 \ 2 \ 3 \ 4 \ 5 \\ \lambda_1 = 1 \ \underline{4} \ 3 \ 2 \ 1 \quad (a < b < c) \\ \lambda_2 = \underline{2} \ 1 \ \underline{3} \ \underline{2} \ 1 \quad (b < c < a) \\ \lambda_3 = 1 \ 3 \ 1 \ 1 \ \underline{1} \quad (c < a < b) \\ \lambda^+ = 2 \ 4 \ 3 \ 2 \ 1 \\ k = 1 \ 1 \ 2 \ 2 \ 3 \\ \mathbf{x} = b \ a \ b \ b \ c \end{array}$$

Fig. 4. Lyndon arrays based on rotated orders for $\sigma = 3$.

- in columns $i = 1, 2$, the respective maximum values $\lambda_2[1] = 2$, $\lambda_1[2] = 4$ occur only once, so that Case 1 applies, and $\mathbf{x}[1] = b$, $\mathbf{x}[2] = a$;
- in column $i = 3$, we find $k = 2 < \sigma$, so Case 2 applies, and since $2 \in P[3]$, $3 \notin P[3]$, we choose $\mathbf{x}[3] = \ell_2 = b$;
- similarly in column $i = 4$, Case 2 applies and again we choose $\mathbf{x}[4] = b$;
- of course in column $i = n = 5$, Case 3 applies, and we set $\mathbf{x}[5] = c$ because, while position 5 is covered by preceding entries in λ_1 and λ_2 , it is not covered by any preceding entry in λ_3 .

In order to prove Theorem 7, we first need the following:

Lemma 8 *Let $i \in 1..n$ be a position in \mathbf{x} such that $\mathbf{x}[i] = \ell_j \in \Sigma$ for some $j \in 1..\sigma$, $\sigma \geq 3$. Then $\lambda^+[i] = \lambda_j[i] \geq \lambda_{j'}[i]$ for every $\ell_{j'} \in \Sigma$.*

Proof. Assume the contrary. Then there exists $j' \neq j$ ($\ell_{j'} \neq \ell_j$) such that $\lambda_{j'}[i] > \lambda_j[i]$. Suppose now that for some position h satisfying $i < h < i + \lambda_j[i]$ (in the range of the Lyndon word corresponding to R_j that begins at i), $\mathbf{x}[h] = \ell_*$, where ℓ_* is a letter such that $\ell_* < \ell_j$ in $R_{j'}$. But this implies that $\lambda_{j'}[i] < \lambda_j[i]$, contradicting our original assumption. Thus every letter ℓ that occurs in the range $\mathbf{x}[i..i+\lambda_j-1]$ must satisfy $\ell \geq \ell_j$ in both R_j and $R_{j'}$. Since the same condition holds for both rotations, therefore $\lambda_{j'}[i]$ can be at most equal to $\lambda_j[i]$. \square

If now we suppose, as in Case 1 of Theorem 7, that there exists a single j_1 such that $\lambda_{j_1}[i]$ is maximum, then it follows immediately from Lemma 8 that, for every rotation R_t , $t \in 1..\sigma$, except $t = j_1$, $\mathbf{x}[i] \neq \ell_t$. Thus $\mathbf{x}[i] = \ell_{j_1}$, establishing Case 1.

The next result gives us a basis for establishing Case 2 by providing a simple characterization of the entries in $P[i]$ when $1 < |P[i]| < \sigma$:

Lemma 9 *Suppose that j_1 is the least value and $j_2 > j_1$ the greatest value (with $j_2 - j_1 < \sigma - 1$) such that for some $i \in 1..n$, $\lambda^+[i] = \lambda_{j_1}[i] = \lambda_{j_2}[i]$. Then $P[i]$ contains exactly one of*

$$\begin{aligned} P_1 &= j_1, j_1+1, \dots, j_2-1, \\ P_2 &= j_2, j_2+1, \dots, \sigma, 1, 2, \dots, j_1-1, \end{aligned}$$

where we suppose that the sequence $1, 2, \dots, j_1-1$ is empty if $j_1 = 1$.

Proof. By hypothesis $\mathbf{x}^* = \mathbf{x}[i..i + \lambda^+[i] - 1]$ is a Lyndon word in both orders R_{j_1} and R_{j_2} . Therefore every letter in \mathbf{x}^* must be greater than or equal to $\mathbf{x}[i]$ in *both* orders

$$\begin{aligned} J_1 &= \{j_1 < j_1 + 1 < \dots < j_2 - 1\}, \\ J_2 &= \{j_2 < j_2 + 1 < \dots < \sigma < 1 < 2 \dots < j_1 - 1\}, \end{aligned} \quad (5)$$

where again we must take account of the special case $j_1 = 1$. We can write $R_{j_1} \equiv J_1\{j_2 - 1 < j_2\}J_2$ and $R_{j_2} \equiv J_2\{j_1 - 1 < j_1\}J_1$. Now observe that if \mathbf{x}^* contains letters from *both* J_1 and J_2 , there must be at least one letter in one of the two orderings that is less than $\mathbf{x}[i]$, and so \mathbf{x}^* cannot be a Lyndon word in one of R_{j_1}, R_{j_2} . Thus \mathbf{x}^* contains letters from exactly one of P_1, P_2 , as required. \square

In the context of Lemma 9, consider a maximal sequence of entries

$$j', j'+1, \dots, j'+t, \quad t > 0, \quad (6)$$

in $\lambda^+[i]$, where $(j'+t+1) \bmod \sigma \notin \lambda^+[i]$ — as in Lemma 9, we suppose that the sequence is circular, so that 1 follows σ . Recall that $R_{j'+1}$ is the rotation of $R_{j'}$ that turns the least letter $\ell_{j'}$ of rotation $R_{j'}$ into the greatest letter of rotation $R_{j'+1}$. Thus the occurrence of $j'+1$ in the sequence $\lambda^+[i]$ ensures that the letter $\ell_{j'}$ cannot be the first letter of the Lyndon array at position i — if it were, then in $R_{j'+1}$, we could have only $\lambda_{j'}[i] = 1$, certainly not maximum. It follows that only the final letter $\ell_{j'+t}$ in the sequence (6) can be the first letter of the Lyndon array at i , because it is the only letter that is not rotated. Noting that in both of the two possible orders given in (5) $j'+t$ will be the *leftmost* occurrence in $P[i]$, we thus establish Case 2 of Theorem 7.

In order to deal with Case 3, we first need:

Lemma 10 $|P[i]| = \sigma \iff \lambda^+[i] = 1$.

Proof. Suppose $|P[i]| = \sigma$. Since every letter in Σ occurs on the right in some rotation R_j of the alphabet, and so is maximum, it follows that $\lambda_j[i] = 1$ for some j . But since every such j yields a maximum $\lambda_j[i]$, it follows that $\lambda^+[i] = 1$. Conversely, if $\lambda^+[i] = 1$, then every rotation R_j yields $\lambda_j[i] = 1$, so that $|P[i]| = \sigma$. \square

Now consider $i = n$, and note that for every rotation R_j , $\lambda_j[n] = 1$, so that $\lambda^+[n] = 1$. Note also that in this case $\mathbf{x}[n]$ will be covered by $\lambda_j[i]$ for every position $i < n$ that marks the rightmost occurrence of ℓ_j in \mathbf{x}

— except for j such that $\ell_j = \mathbf{x}[i] = \mathbf{x}[n]$. Thus the letter ℓ_j satisfies exactly one of two conditions:

- either ℓ_j does not previously occur in \mathbf{x} ; or
- under rotation R_j , at any position $i < n$ such that $\mathbf{x}[i] = \ell_j$, $\lambda_j[i]$ does not cover $\mathbf{x}[n]$.

Hence Case 3 provides the basis for assigning $\mathbf{x}[n]$. To complete the proof of Theorem 7, we need one more result:

Lemma 11 *If $|P[i]| = \sigma$ for $i < n$, then $|P[i']| = \sigma$ for every $i' \in i..n$.*

Proof. We know $\mathbf{x}[i] = \ell_j$, a minimum letter under rotation R_j , and from Lemma 10 we know that $\lambda^+[i] = 1$. This is possible only if the same minimum letter occurs also at positions $i+1, i+2, \dots, n$, as required. \square

Therefore Case 3 identifies strings \mathbf{x} with suffix ℓ^{n-i+1} for some i and some ℓ .

Theorem 7 justifies Algorithm ConstructString, shown in Figure 5, that in $\mathcal{O}(\sigma n)$ time constructs the unique string \mathbf{x} on a given ordered alphabet Σ , based on σ rotations of a given Lyndon array λ .

4 Final Remarks

In this paper we have started to analyze the relationship between a string and its Lyndon arrays corresponding to cyclic orderings of the underlying alphabet. Many problem areas remain:

1. As indicated by Problem 3, we know of no efficient algorithm to compute a string on a minimum alphabet consistent with a given valid Lyndon array L^* .
2. Similarly, as we have seen, it appears to be difficult to reconstruct a string exactly from Lyndon arrays. In Section 3 we have presented an algorithm to reconstruct a string \mathbf{x} on an alphabet of size σ given Lyndon arrays of \mathbf{x} based on σ rotations of the alphabet. It appears that indeed in the worst case σ such rotations are required, as shown by the example in Figure 6.

If in Figure 6 only rotations I–III are considered, then the selection $\mathbf{x}[2] = c$ rather than d would be made; if rotations II–IV were used, no determination could be made for $\mathbf{x}[1]$; and if rotations I, III, IV were used, column $i = 3$ would become an instance of Case 2(b), and

```

procedure ConstructString ( $\lambda, \sigma, n; \mathbf{x}$ )
 $i \leftarrow 1$ ;  $cover \leftarrow \sigma^n$ 
while  $i < n$  do
     $maxlen \leftarrow 0$  — Length of maximum  $\lambda[1.. \sigma, i]$ 
    for  $j \leftarrow 1$  to  $\sigma$  do
         $\triangleright$  Treat  $\lambda$  as a two-dimensional array; find final
         $\triangleright$  letter  $jmax$  & frequency  $freqmax$  for  $maxlen$ .
        if  $\lambda[j, i] > maxlen$  then
             $maxlen \leftarrow \lambda[j, i]$ ;  $jmax \leftarrow j$ ;  $freqmax \leftarrow 1$ 
        elseif  $\lambda[j, i] = maxlen$  then
             $jmax \leftarrow j$ ;  $freqmax \leftarrow freqmax + 1$ 
         $\triangleright$  Recompute maximum range covered by letter  $j$ .
         $cover[j] \leftarrow \max(cover[j], j + \lambda[j, i])$ 
        if  $freqmax < \sigma$  then
             $\triangleright$  Cases 1 and 2 — Lemmas 8 & 9.
             $\mathbf{x}[i] \leftarrow jmax$ ;  $i \leftarrow i + 1$ 
        else
             $\triangleright$  Case 3 — Lemma 10.
             $j' \leftarrow 1$ 
             $\triangleright$  Find the letter that yields no cover of position  $i$ .
            while  $j' \leq \sigma$  and  $i < cover[j']$  do  $j' \leftarrow j' + 1$ 
             $\mathbf{x}[i] \leftarrow j'$ ;  $i \leftarrow i + 1$ 
         $\triangleright$  In accordance with Lemma 11, extend to position  $n$ .
        while  $i \leq n$  do  $\mathbf{x}[i] \leftarrow j'$ ;  $i \leftarrow i + 1$ 

```

Fig. 5. Constructing a string from rotated Lyndon arrays.

$\mathbf{x}[2] = b$ could not be selected. Thus it appears that, on an alphabet of size $\sigma > 2$, the Lyndon array does not strongly determine the underlying string.

$i =$	1	2	3	4	
$\mathbf{x} =$	a	d	b	c	
$\lambda_1 =$	4	1	2	1	$(a < b < c < d)$ I
$\lambda_2 =$	1	1	2	1	$(b < c < d < a)$ II
$\lambda_3 =$	1	2	1	1	$(c < d < a < b)$ III
$\lambda_4 =$	1	3	2	1	$(d < a < b < c)$ IV

Fig. 6. $\mathbf{x} = adbc$ and its Lyndon arrays: consideration of fewer than four rotations of the alphabet may not allow \mathbf{x} to be reconstructed.

3. Thus it would be of interest to determine minimal criteria for the reconstruction of \mathbf{x} — the least number of rotations or the size of the smallest alphabet. Is there any hope of reconstructing a string based on fewer than σ permutations of the alphabet?
4. The study of UMFFs has led to the idea of V -words, analogous structures to Lyndon words, derived from a global non-lexicographic ordering of strings called V -order [11, 16, 14, 15]. Also, linear-time algorithms for computing Lyndon border and Lyndon suffix arrays have recently been proposed [2]. There is scope to extend the results of this paper to the UMFF based on V -order [1], then further to UMFFs in their full generality.

References

1. A. Alatabbi, J. W. Daykin, J. Kärkkäinen, M. S. Rahman, and W. F. Smyth. V -Order: new combinatorial properties & a simple comparison algorithm. *Discrete Appl. Math.*, 215:41–46, 2016.
2. A. Alatabbi, J. W. Daykin, and M. S. Rahman. Linear algorithms for computing the Lyndon border array and the Lyndon suffix array. 2015. *arXiv:1506.06983*.
3. A. Alatabbi, M. S. Rahman, and W. F. Smyth. Inferring an indeterminate string from a prefix graph. *Journal of Discrete Algorithms*, 32:6–13, 2015.
4. U. Baier. Linear-time suffix sorting — a new approach for suffix array construction. M.Sc. Thesis, University of Ulm, 2015.
5. U. Baier. Linear-time suffix sorting — a new approach for suffix array construction. In Roberto Grossi and Moshe Lewenstein, editors, *27th Annual Symposium on Combinatorial Pattern Matching (CPM 2016)*, volume 54 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 23:1–23:12, Dagstuhl, Germany, 2016. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
6. H. Bannai, T. I. S. Inenaga, Y. Nakashima, M. Takeda, and K. Tsuruta. The “runs” theorem. 2014. *arXiv:1406.0263v6*.
7. 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.
8. M. Christodoulakis, P. J. Ryan, W. F. Smyth, and S. Wang. Indeterminate strings, prefix arrays and undirected graphs. *Theoretical Comput. Sci.*, 600:34–48, 2015.
9. J. Clément, M. Crochemore, and G. Rindone. Reverse engineering prefix tables. In *Proc. 26th Symp. on Theoretical Aspects of Computer Science*, pages 289–300, 2009.
10. M. Crochemore, C. S. Iliopoulos, S. P. Pissis, and G. Tischler. Cover array string reconstruction. In *Symp. on Combinatorial Pattern Matching*, volume 6129 of *Lecture Notes in Comput. Sci.*, pages 251–259. Springer, 2010.
11. T.-N. Danh and D. E. Daykin. The structure of V -order for integer vectors. *Congressus Numerantium*, 113:43–53, 1996.
12. D. E. Daykin, J. W. Daykin, C. S. Iliopoulos, and W. F. Smyth. Generic algorithms for factoring strings. In *Proc. Memorial Symp. for Rudolf Ahlswede, H. Aydinian, F. Cicalese & C. Depepe, (eds.)*, volume 7777 of *Lecture Notes in Comput. Sci.*, pages 402–418. Springer-Verlag, 2013.

13. D. E. Daykin, J. W. Daykin, and W. F. Smyth. Combinatorics of unique maximal factorization families (UMFFs). *Fund. Inform.* 97–3, *Special Issue on Stringology, R. Janicki, S. J. Puglisi & M. S. Rahman (eds.)*, pages 295–309, 2009.
14. D. E. Daykin, J. W. Daykin, and W. F. Smyth. String comparison and Lyndon-like factorization using V -order in linear time. In *22nd Annual Symp. on Combinatorial Pattern Matching*, volume 6661 of *Lecture Notes in Computer Science*, pages 65–76. Springer-Verlag, 2011.
15. D. E. Daykin, J. W. Daykin, and W. F. Smyth. A linear partitioning algorithm for hybrid Lyndons using V -order. *Theoret. Comput. Sci.*, 483:149–161, 2013.
16. D.E. Daykin and J.W. Daykin. Lyndon-like and V -order factorizations of strings. *J. Discrete Algorithms*, 1(3–4):357–365, 2003.
17. D.E. Daykin and J.W. Daykin. Properties and construction of unique maximal factorization families for strings. *Internat. J. Found. Comput. Sci.*, 19(4):1073–1084, 2008.
18. J.-P. Duval. Factorizing words over an ordered alphabet. *J. Algorithms*, 4:363–381, 1983.
19. J.-P. Duval, T. Lecroq, and A. Lefebvre. Border array on bounded alphabet. *Journal of Automata, Languages and Combinatorics*, 10(1):51–60, 2005.
20. J.-P. Duval, T. Lecroq, and A. Lefebvre. Efficient validation and construction of Knuth-Morris-Pratt arrays. In *Proc. Conference in Honour of Donald E. Knuth*, 2007.
21. F. Franek, S. Gao, W. Lu, P. J. Ryan, W. F. Smyth, Y. Sun, and L. Yang. Verifying a border array in linear time. *J. Combinatorial Math. & Combinatorial Computing*, 42:223–236, 2002.
22. F. Franek, A. S. M. S. Islam, M. S. Rahman, and W. F. Smyth. Algorithms to compute the Lyndon array. In Jan Holub and Jan Ždárek, editors, *Proceedings of the Prague Stringology Conference 2016*, pages 172–184, Czech Technical University in Prague, Czech Republic, 2016.
23. P. Gawrychowski, A. Jež, and L. Jež. Validating the Knuth-Morris-Pratt failure function, fast and on-line. In *Proc. 5th Annual Computer Science Symposium in Russia*, volume 6072 of *Lecture Notes in Comput. Sci.*, pages 132–143. Springer-Verlag, 2010.
24. P. Gawrychowski, A. Jež, and L. Jež. Validating the Knuth-Morris-Pratt failure function, fast and on-line. *Theory of Computing Systems*, 54:337–372, 2014.
25. J. Sawada and F. 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.
26. B. Smyth. *Computing Patterns in Strings*. Pearson/Addison-Wesley, 2003.