


Article

# Computing Maximal Lyndon Substrings of a String

Frantisek Franek<sup>1,†</sup>, Michael Liut<sup>2,\*,†</sup> 

<sup>1</sup> Department of Computing and Software, McMaster University, Hamilton, Ontario, Canada; franek@mcmaster.ca

<sup>2</sup> Department of Mathematical and Computational Sciences, University of Toronto Mississauga, Mississauga, Ontario, Canada michael.liut@utoronto.ca

\* Correspondence: michael.liut@utoronto.ca

† These authors contributed equally to this work.

**Abstract:** There are two reasons to have an efficient algorithm for identifying all maximal Lyndon substrings of a string: firstly, Bannai et al. introduced in 2015 a linear algorithm to compute all runs of a string that relies on knowing all maximal Lyndon substrings of the input string, and secondly, Franek et al. showed in 2017 a linear equivalence of sorting suffixes and sorting maximal Lyndon substrings of a string, inspired by a novel suffix sorting algorithm of Baier. In 2016, Franek et al. presented a brief overview of algorithms for computing the Lyndon array that encodes the knowledge of maximal Lyndon substrings of the input string. Among the presented were two well-known algorithms for computing the Lyndon array: a quadratic in-place algorithm based on iterated Duval's algorithm for Lyndon factorization, and a linear algorithmic scheme based on linear suffix sorting, computing inverse suffix array, and applying to it the *Next Smaller Value* algorithm. Duval's algorithm works for strings over any ordered alphabet, while for linear suffix sorting, a constant or an integer alphabet is required. The authors at that time were not aware of Baier's algorithm. In 2017, our research group proposed a novel algorithm for the Lyndon array. Though the proposed algorithm is linear in the average case and has  $O(n \log(n))$  worst-case complexity, it is interesting as it emulates the fast Fourier algorithm's recursive approach and introduces  $\tau$ -reduction that might be of independent interest. In 2018, we presented a linear algorithm to compute the Lyndon array of a string inspired by Phase I of Baier's algorithm for suffix sorting. This paper presents theoretical analysis of these two algorithms and provides empirical comparisons of both their C++ implementations with respect to iterated Duval's algorithm.

**Keywords:** Combinatorics on Words; String Algorithms; Regularities in Strings; Suffix Sorting; Lyndon Substrings; Lyndon Arrays; Maximal Lyndon Substrings; Tau-Reduction Algorithm; Baier's Sort Algorithm; Iterative Duval Algorithm;

## 1. Introduction

There are at least two reasons for having an efficient algorithm for identifying all maximal Lyndon substrings of a string: firstly, Bannai et al. introduced in 2015 (arXiv, [1]) and published in 2017, [2], a linear algorithm to compute all runs in a string that depends on knowing all maximal Lyndon substrings of the input string, and secondly, in 2017, Franek et al. in [3] showed a linear equivalence of sorting suffixes and sorting maximal Lyndon substrings, inspired by Phase II of a suffix sorting algorithm introduced by Baier in 2015 (Master's thesis, [4]), and published in 2016, [5].

The most significant feature of the runs algorithm presented in [2] is that it relies on knowing the maximal Lyndon substrings of the input string for some order of the alphabet and for the inverse of that order, while all other linear algorithms for runs rely on Lempel-Ziv factorization of the input string. Thus, computing runs became yet another application of Lyndon words. It also raised the issue

34 which approach may be more efficient: to compute the Lempel-Ziv factorization or to compute all  
 35 maximal Lyndon substrings. There are several efficient linear algorithms for Lempel-Ziv factorization  
 36 (e.g. see [6,7] and the references therein).

37 Interestingly, Kosolobov [8] shows that for a general alphabet, in the decision tree model, the runs  
 38 problem is easier than the Lempel-Ziv decomposition. His result supports the conjecture that there  
 39 must be a linear Random Access Memory model algorithm finding all runs.

40 Baier introduced in [4], and published in [5], a new algorithm for suffix sorting. Though  
 41 Lyndon strings are never mentioned in [4,5], it was noticed by Cristoph Diegelmann in a personal  
 42 communication, [9], that Phase I of the Baier's suffix sort identifies and sorts all maximal Lyndon  
 43 substrings.

44 The maximal Lyndon substrings of a string  $x = x[1..n]$  can be best encoded in the so-called  
 45 *Lyndon array*, introduced in [10]: an integer array  $\mathcal{L}[1..n]$  so that for any  $i \in 1..n$ ,  $\mathcal{L}[i] =$  the length of  
 46 the maximal Lyndon substring starting at the position  $i$ .

```

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

```

Figure 1. Algorithm IDLA

47 In an overview [10], Franek et al. discussed an algorithm based on an iterative application  
 48 of Duval's Lyndon factorization algorithm, [11], which we refer here as IDLA, and an algorithmic  
 49 scheme based on Hohlweg and Reutenauer's work [12], which we refer to as SSLA. The authors  
 50 were not aware of Baier's algorithm at that time. Two additional algorithms were presented there, a  
 51 quadratic recursive application of Duval's algorithm, and an algorithm NSV\* with possibly  $\mathcal{O}(n \log(n))$   
 52 worst-case complexity based on ranges that can be compared in constant time for constant alphabets.  
 53 The correctness of NSV\* and its complexity were discussed there just informally.

54 The algorithm IDLA, see Figure 1, is simple and in-place, so no space is required except the storage  
 55 for the string and the storage for the Lyndon array. It is completely independent of the alphabet of the  
 56 string and does not require the alphabet to be sorted, all it requires is that the alphabet be ordered, i.e.,  
 57 only pairwise comparisons of the alphabet symbols are needed. Its weakness is its quadratic worst-case  
 58 complexity that becomes a problem for longer strings with long maximal Lyndon substrings as one of  
 59 our experiments showed.

60 In our empirical work, we used IDLA as a control for comparison and as a verifier of the results.  
 61 Note that the reason the procedure MaxLyn, see Fig. 1, really computes the maximal Lyndon prefix  
 62 is not obvious and is based on the properties of periods of prefixes, see [11], or Observation 6 and  
 63 Lemma 11 in [10]. The following lemma characterizes maximal Lyndon substrings in terms of the  
 64 relationships of the suffixes, and follows from the work of Hohlweg and Reutenauer [12].

65 **Lemma 1.** Consider a string  $x[1..n]$  over an alphabet ordered by  $\prec$ . The substring  $x[i..j]$  is Lyndon if  $x[i..n] \prec$   
 66  $x[k..n]$  for any  $i < k \leq j$ , and is maximal Lyndon if it is Lyndon and either  $j = n$  or  $x[j+1..n] \prec x[i..n]$ .

67 Thus, Lyndon array is an *NSV* (Next Smaller Value) array of the inverse suffix array. Consequently,  
 68 the Lyndon array can be computed by sorting the suffixes – i.e. computing the suffix array, then  
 69 computing the inverse suffix array, and then applying *NSV* to it, see [10]. Computing the inverse suffix  
 70 array and applying *NSV* are “naturally” linear, and computing the suffix array can be implemented to  
 71 be linear, see [10,13] and the references therein. The execution and space characteristics are dominated  
 72 by those of the first step, i.e., computation of the suffix array. We refer to this scheme as *SSLA*.

73 In 2018, a linear algorithm to compute the Lyndon array from a given Burrows-Wheeler transform  
 74 was presented, [14]. Since the Burrows-Wheeler transform is computed in linear time from the suffix  
 75 array, it is yet another scheme of how to obtain the Lyndon array via suffix sorting: compute the suffix  
 76 array, from the suffix array compute the Burrows-Wheeler transform, and then compute the Lyndon  
 77 array during the inversion of the Burrows-Wheeler transform. We refer to this scheme as *BWLA*.

78 The introduction of Baier’s suffix sort in 2015 and the consequent realization of the connection to  
 79 maximal Lyndon substrings brought up the realization that there was an elementary<sup>1</sup> algorithm to  
 80 compute the Lyndon array, and that, despite its original clumsiness, could be eventually refined to  
 81 outperform any *SSLA* or *BWLA* implementation: any implementation of suffix-sorting-based scheme  
 82 requires a full suffix sort and then some additional processing, while Baier’s approach is “just” a partial  
 83 suffix sort, see [15].

84 In this work, we present two additional algorithms for Lyndon array not discussed in [10]. The  
 85 C++ source code of the three implementations, *IDLA*, *TRLA*, and *BSLA* is available, see [16]. Note that  
 86 the procedure *IDLA* is in `lynarr.hpp` file.

87 The first algorithm presented here is *TRLA*. *TRLA* is a  $\tau$ -reduction based Lyndon array algorithm  
 88 that follows Farach’s approach used in his remarkable linear algorithm for suffix tree construction [17],  
 89 and reproduced very successfully in all linear algorithms for suffix sorting (e.g. see [13,18] and the  
 90 references therein). Farach’s approach follows Cooley-Tukey algorithm for the fast Fourier transform  
 91 relying on recursion to lower the quadratic complexity to  $O(n \log(n))$  complexity, see [19]. *TRLA* was  
 92 first introduced by the authors in 2019, see [20,21], and presented as a part of Liut’s Ph.D. thesis [22].

93 The second algorithm, *BSLA*, is a Baier’s sort-based Lyndon array algorithm. *BSLA* is based on  
 94 the idea of Phase I of Baier’s suffix sort, though our implementation necessarily differs from Baier’s.  
 95 *BSLA* was first introduced at the Prague Stringology Conference 2018, [15], and also presented as a  
 96 part of Liut’s Ph.D. thesis [22] in 2019; here we present a complete and refined theoretical analysis of  
 97 the algorithm and a more efficient implementation than that initially introduced.

98 The paper is structured as follows: in Section 2, the basic notions and terminology are presented,  
 99 in Section 3, the *TRLA* algorithm is presented and analyzed. In Section 4, the *BSLA* algorithm is  
 100 presented and analyzed. In Section 7, the empirical measurements of the performance of *IDLA*, *TRLA*,  
 101 and *BSLA* are presented on data sets with random strings of various lengths and over various alphabets.  
 102 The results are presented both in tabular and graphical forms. In Section 6, the conclusion of the  
 103 research and the future work are presented.

## 104 2. Basic Notation and Terminology

105 For two integers  $i \leq j$ , the *range*  $i..j = \{k \text{ integer} : i \leq k \leq j\}$ . An *alphabet* is a finite or infinite sets  
 106 of *symbols*, or equivalently called *letters*. We assume that a sentinel symbol  $\$$  is not in the alphabet and  
 107 is always assumed to be lexicographically the smallest. A *string* over an alphabet  $\mathcal{A}$  is a finite sequence  
 108 of symbols from  $\mathcal{A}$ . A  *$\$$ -terminated string* over  $\mathcal{A}$  is a string over  $\mathcal{A}$  terminated by  $\$$ , where  $\$ \notin \mathcal{A}$ . We  
 109 use the array notation indexing from 1 for strings, thus  $x[1..n]$  indicates a string of length  $n$ , the first

<sup>1</sup> not relying on a pre-processed global data structure such as a suffix array or a Burrows-Wheeler transform

110 symbol is the symbol with index 1, i.e.  $x[1]$ , the second symbol is the symbol with index 2, i.e.  $x[2]$ ,  
 111 etc. Thus,  $x[1..n] = x[1]x[2]...x[n]$ . For a  $\$$ -terminated string  $x$  of length  $n$ ,  $x[n+1] = \$$ . The **alphabet**  
 112 **of string**  $x$ , denoted as  $\mathcal{A}_x$ , is the set of all distinct alphabet symbols occurring in  $x$ . By a **constant**  
 113 **alphabet** we mean a fixed finite alphabet. A string  $x$  is over an **integer alphabet** if  $\mathcal{A}_x \subseteq \{0, 1, \dots, |x|\}$ .  
 114 Thus, the class of **strings over integer alphabets** =  $\{x \mid x \text{ is a string over } \{0, 1, \dots, |x|\}\}$ . A string  $x$  over  
 115 an integer alphabet is **tight** if  $\mathcal{A}_x = \{0, 1, \dots, k\}$  for some  $k \leq |x|$ . For instance,  $x = 010$  is tight as  
 116  $\mathcal{A}_x = \{0, 1\}$ , while  $y = 020$  is not as  $\mathcal{A}_y = \{0, 2\}$ ; i.e. 1 is missing from  $\mathcal{A}_y$ .

117 We use a **bold font** to denote strings, thus  $x$  denotes a string, while  $x$  denotes some  
 118 other mathematical entity such as an integer. The **empty string** is denoted by  $\varepsilon$  and has  
 119 length 0. The **length** or **size** of string  $x = x[1..n]$  is  $n$ . The length of a string  $x$  is denoted  
 120 by  $|x|$ . For two strings  $x = x[1..n]$  and  $y = y[1..m]$ , the **concatenation**  $xy$  is a string  $u$

$$121 \text{ where } u[i] = \begin{cases} x[i] \text{ for } i \leq n, \\ y[i-n] \text{ for } n < i \leq n+m. \end{cases}$$

122  
 123  
 124 If  $x = uvw$ , then  $u$  is a **prefix**,  $v$  a **substring**, and  $w$  a **suffix** of  $x$ . If  $u$  (respectively,  $v$ ,  $w$ ) is empty,  
 125 then it is called a **trivial prefix** (respectively, **trivial substring**, **trivial suffix**), if  $|u| < |x|$  (respectively,  
 126  $|v| < |x|$ ,  $|w| < |x|$ ) then it is called a **proper prefix** (respectively, **proper substring**, **proper suffix**). If  
 127  $x = uv$ , then  $vu$  is called a **rotation** or a **conjugate** of  $x$ ; if either  $u = \varepsilon$  or  $v = \varepsilon$ , then the rotation is  
 128 called **trivial**. A non-empty string  $x$  is **primitive** if there is no string  $y$  and no integer  $k \geq 2$  so that  
 129  $x = y^k = \underbrace{yy \cdots y}_{k \text{ times}}$ .

130 A non-empty string  $x$  has a non-trivial border  $u$  if  $u$  is both a non-trivial proper prefix and a  
 131 non-trivial proper suffix of  $x$ . Thus, both  $\varepsilon$  and  $x$  are trivial borders of  $x$ . A string without a non-trivial  
 132 border is call **unbordered**.

133 Let  $\prec$  be a total order of an alphabet  $\mathcal{A}$ . The order is extended to all finite strings over the alphabet  
 134  $\mathcal{A}$ : for  $x = x[1..n]$  and  $y = y[1..m]$ ,  $x \prec y$  if either  $x$  is a proper prefix of  $y$ , or there is a  $j \leq \min\{n, m\}$   
 135 so that  $x[1] = y[1], \dots, x[j-1] = y[j-1]$  and  $x[j] \prec y[j]$ . This total order induced by the order of the  
 136 alphabet is called a **lexicographic** order of all non-empty strings over  $\mathcal{A}$ . We denote by  $x \preceq y$  if either  
 137  $x \prec y$  or  $x = y$ . A string  $x$  over  $\mathcal{A}$  is **Lyndon** for a given order  $\prec$  of  $\mathcal{A}$  if  $x$  is strictly lexicographically  
 138 smaller than any non-trivial rotation of  $x$ . In particular:

$$139 \quad x \text{ is Lyndon} \Rightarrow x \text{ is unbordered} \Rightarrow x \text{ is primitive}$$

140 Note that the reverse implications do not hold:  $aba$  is primitive but neither unbordered, nor Lyndon,  
 141 while  $acaab$  is unbordered, but not Lyndon. A substring  $x[i..j]$  of  $x[1..n]$ ,  $1 \leq i \leq j \leq n$  is a **maximal**  
 142 **Lyndon substring** of  $x$  if it is Lyndon and either  $j = n$  or for any  $k > j$ ,  $x[i..k]$  is not Lyndon. The  
 143 **Lyndon array** of a string  $x = x[1..n]$  is an integer array  $\mathcal{L}[1..n]$  so that  $\mathcal{L}[i] = j$  where  $j \leq n-i$  is a  
 144 maximal integer such that  $x[i..i+j-1]$  is Lyndon. Alternatively, we can define it as an integer array  
 145  $\mathcal{L}'[1..n]$  so that  $\mathcal{L}'[i] = j$  where  $j$  is the last position of the maximal Lyndon substring starting at the  
 146 position  $i$ . The relationship between those two definitions is straightforward:  $\mathcal{L}'[i] = \mathcal{L}[i] + i - 1$ , or  
 147  $\mathcal{L}[i] = \mathcal{L}'[i] - i + 1$ .

### 148 3. $\mathcal{T}$ -Reduction Algorithm – TRLA

149 The first idea of the algorithm was proposed in Paracha's 2017 Ph.D. thesis, [23]. It follows  
 150 Farach's approach [17]:

- 151 (1) reduce the input string  $x$  to  $y$ ,
- 152 (2) by recursion compute the Lyndon array of  $y$ , and
- 153 (3) from the Lyndon array of  $y$  compute the Lyndon array of  $x$ .

155 The input strings are  $\$$ -terminated strings over integer alphabets. The reduction computed in  
 156 (1) is important. All linear algorithms for suffix array computations use the proximity property of  
 157 suffixes: comparing  $x[i..n]$  and  $x[j..n]$  can be done by comparing  $x[i]$  and  $x[j]$ , and if they are the  
 158 same, comparing  $x[i+1..n]$  with  $x[j+1..n]$ . For instance, in the first linear algorithm for suffix array by  
 159 Kärkkäinen and Sanders [24], obtaining the sorted suffixes for positions  $i \equiv 0 \pmod{3}$  and  $i \equiv 1 \pmod{3}$   
 160 via the recursive call is sufficient to determine the order of suffixes for  $i \equiv 2 \pmod{3}$  positions, and  
 161 then to merge both lists together. However, there is no such proximity property for maximal Lyndon  
 162 substrings, so the reduction itself must have a property that helps determine some of the values of the  
 163 Lyndon array of  $x$  from the Lyndon array of  $y$  and compute the rest. We present such a reduction that  
 164 we call  $\tau$ -reduction, and it may be of some general interest as it preserves order of some suffixes and  
 165 hence, by Lemma 1, some maximal Lyndon substrings.

166 The algorithm computes  $y$  as a  $\tau$ -reduction of  $x$  in step (1) in linear time and in step (3) it expands  
 167 the Lyndon array of the reduced string computed by step (2) to an incomplete Lyndon array of the  
 168 original string also in linear time. However, it computes the missing values of the incomplete Lyndon  
 169 array in at most  $O(n \log(n))$  steps resulting in the overall worst-case complexity of  $O(n \log(n))$ . When  
 170 the input string is such that the missing values of the incomplete Lyndon array of  $x$  can be computed in  
 171 linear time, the overall execution of algorithm is linear as well, and thus, the average case complexity  
 172 is linear in the length of the input string. Since for  $\tau$ -reduction, the size of  $\tau(x)$  is at most  $\frac{2}{3}|x|$ , we  
 173 eventually obtain through the recursion of step (2) applied to  $\tau(x)$  a partially filled Lyndon array of  
 174 the input string; the array is about  $\frac{1}{2}$  to  $\frac{2}{3}$  full and for every position  $i$  with an unknown value, the  
 175 values at positions  $i-1$  and  $i+1$  are known. In particular, the values at position 1 and position  $n$  are  
 176 both known. So, a lot of information is provided by the recursive step. For instance, for 00011001, via  
 177 the recursive call we would identify the maximal Lyndon substrings that are underlined in 00011001  
 178 and would need to compute the missing maximal Lyndon substrings that are underlined in 00011001.  
 179 We describe the  $\tau$ -reduction in several steps: first the  $\tau$ -pairing, then choosing the  $\tau$ -alphabet, and  
 180 finally the computation of the  $\tau$ -reduction of  $x$ .

### 181 3.1. $\tau$ -Pairing

182 Consider a  $\$$ -terminated string  $x = x[1..n]$  whose alphabet  $\mathcal{A}_x$  is ordered by  $\prec$  where  
 183  $x[n+1] = \$$  and  $\$ \prec a$  for any  $a \in \mathcal{A}_x$ . A  $\tau$ -pair consists of a pair of adjacent positions from  
 184 the range  $1..n+1$ . The  $\tau$ -pairs are computed by induction:

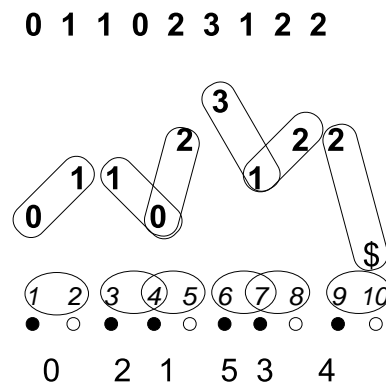
- 185 • the initial  $\tau$ -pair is  $(1, 2)$ ;
- 186 • if  $(i-1, i)$  is the last  $\tau$ -pair computed, then:
  - 187 **if**  $i = n-1$  **then**
  - 188 the next  $\tau$ -pair is set to  $(n, n+1)$
  - 189 **stop**
  - 190 **elseif**  $i \geq n$  **then**
  - 191 **stop**
  - 192 **elseif**  $x[i-1] \succ x[i]$  **and**  $x[i] \preceq x[i+1]$  **then**
  - 193 the next  $\tau$ -pair is set to  $(i, i+1)$
  - 194 **else**
  - 195 the next  $\tau$ -pair is set to  $(i+1, i+2)$

196 Every position of the input string that occurs in some  $\tau$ -pair as the first element is labeled **black**, all  
 197 others are labeled **white**. Note that most of the  $\tau$ -pairs do not overlap; if two  $\tau$ -pairs overlap, they  
 198 overlap in a position  $i$  such that  $1 < i < n$  and  $x[i-1] \succ x[i]$  and  $x[i] \preceq x[i+1]$ . Moreover, a  $\tau$ -pair can  
 199 be involved in at most one overlap; for an illustration see Figure 2, for the formal proof see Lemma 2.

200 **Lemma 2.** Let  $(i_1, i_1+1) \dots (i_k, i_k+1)$  be the  $\tau$ -pairs of a strings  $x = x[1..n]$ . Then for any  $j, \ell \in 1..k$ :

- 201 (1) if  $|(i_j, i_j+1) \cap (i_\ell, i_\ell+1)| = 1$ , then for any  $m \neq j, \ell$ ,  $|(i_j, i_j+1) \cap (i_m, i_m+1)| = 0$ ,
- 202 (2)  $|(i_j, i_j+1) \cap (i_\ell, i_\ell+1)| \leq 1$ .





**Figure 2.** Illustration of  $\tau$ -reduction of a string **011023122**

The rounded rectangles indicate symbol  $\tau$ -pairs, the ovals indicate the  $\tau$ -pairs below are the colour labels of positions, at the bottom is the  $\tau$ -reduction

203

204

205 **Proof.** By induction. Trivially true for  $|x| = 1$  as  $(1,2)$  is the only  $\tau$ -pair. Assume it true for  $|x| \leq n-1$ .

- 206 • Case  $(i_k, i_k+1) = (n, n+1)$   
 207 Then  $(i_{k-1}, i_{k-1}+1) = (n-2, n-1)$ , and so  $(i_1, i_1+1) \dots (i_{k-1}, i_{k-1}+1)$  are  $\tau$ -pairs of  $x[1..n-1]$ , and  
 208 thus they satisfy (1) and (2) by the induction hypothesis. However,  $(n, n+1) \cap (i_\ell, i_{\ell+1}) = \emptyset$  for  
 209  $1 \leq \ell < k$ , so (1) and (2) hold for  $(i_1, i_1+1) \dots (i_k, i_k+1)$ .
- 210 • Case  $(i_k, i_k+1) = (n-1, n)$  and  $(i_{k-1}, i_{k-1}+1) = (n-2, n-1)$ .  
 211 Therefore,  $(i_1, i_1+1) \dots (i_{k-1}, i_{k-1}+1)$  are  $\tau$ -pairs of  $x[1..n-1]$ , and thus they satisfy (1) and (2)  
 212 by the induction hypothesis. However,  $(i_k, i_k+1) \cap (i_\ell, i_{\ell+1}) = \emptyset$  for  $1 \leq \ell < k-1$ , and  
 213  $(i_k, i_k+1) \cap (i_{k-1}, i_{k-1}+1) = \{i_{k-1}\} = n-1$ , so  $|(i_k, i_k+1) \cap (i_{k-1}, i_{k-1}+1)| \leq 1$  and so (1) and (2)  
 214 hold for  $(i_1, i_1+1) \dots (i_k, i_k+1)$ .
- 215 • Case  $(i_k, i_k+1) = (n-1, n)$  and  $(i_{k-1}, i_{k-1}+1) = (n-3, n-2)$ .  
 216 Then  $(i_1, i_1+1) \dots (i_{k-1}, i_{k-1}+1)$  are  $\tau$ -pairs of  $x[1..n-2]$ , so satisfy (1) and (2) by the induction  
 217 hypothesis. However,  $(i_k, i_k+1) \cap (i_\ell, i_{\ell+1}) = \emptyset$  for  $1 \leq \ell < k$ , so (1) and (2) hold for  
 218  $(i_1, i_1+1) \dots (i_k, i_k+1)$ .

219

□

### 220 3.2. $\tau$ -Reduction

221 For each  $\tau$ -pair  $(i, i+1)$ , we consider the pair of alphabet symbols  $(x[i], x[i+1])$ . We call them  
 222 **symbol  $\tau$ -pairs**. They are in a total order  $\triangleleft$  induced by  $\prec : (x[i_j], x[i_j+1]) \triangleleft (x[i_\ell], x[i_\ell+1])$  if either  
 223  $x[i_j] \prec x[i_\ell]$ , or  $x[i_j] = x[i_\ell]$  and  $x[i_j+1] \prec x[i_\ell+1]$ . They are sorted using radix sort with keys of size 2,  
 224 and assigned letters from a chosen  $\tau$ -alphabet that is a subset of  $\{0, 1, \dots, |\tau(x)|\}$  so that the assignment  
 225 preserves the order. Since the input string is over an integer alphabet, the radix sort is linear.

226 In the example, Figure 2, the  $\tau$ -pairs are  $(1,2)(3,4)(4,5)(6,7)(7,8)(9,10)$  and so the  
 227 symbol  $\tau$ -pairs are  $(0,1)(1,0)(0,2)(3,1)(1,2)(2,\$)$ . The sorted symbol  $\tau$ -pairs are  $(0,1)(0,2)(1,0)$   
 228  $(1,2)(2,\$)(3,1)$ . Thus we chose as our  $\tau$ -alphabet  $\{0, 1, 2, 3, 4, 5\}$  and so the symbol  $\tau$ -pairs are assigned  
 229 these letters:  $(0,1) \rightarrow 0$ ,  $(0,2) \rightarrow 1$ ,  $(1,0) \rightarrow 2$ ,  $(1,2) \rightarrow 3$ ,  $(2,\$) \rightarrow 4$  and  $(3,1) \rightarrow 5$ . Note that the  
 230 assignments respect the order  $\triangleleft$  of the symbols  $\tau$ -pairs, and the natural order  $<$  of  $\{0, 1, 2, 3, 4, 5\}$ .

231 The  $\tau$ -letters are substituted for the symbol  $\tau$ -pairs and the resulting string is terminated with  
 232  $\$$ . This string is called the  **$\tau$ -reduction** of  $x$  and denoted  $\tau(x)$ , and it is a  $\$$ -terminated string over an

233 integer alphabet. For our running example from Figure 2,  $\tau(x) = 021534$ . The next lemma justifies  
234 calling the above transformation a reduction.

235 **Lemma 3.** For any string  $x$ ,  $\frac{1}{2}|x| \leq |\tau(x)| \leq \frac{2}{3}|x|$ .

236 **Proof.** There are two extreme cases; the first is when all the  $\tau$ -pairs do not overlap at all, then  
237  $|\tau(x)| = \frac{1}{2}|x|$ ; and the second is when all the  $\tau$ -pairs overlap, then  $|\tau(x)| = \frac{2}{3}|x|$ . Any other case must  
238 be in between.  $\square$

239 Let  $\mathcal{B}(x)$  denote the set of all black positions of  $x$ . For any  $i \in 1..|\tau(x)|$ ,  $b(i) = j$  where  $j$  is a black  
240 position in  $x$  of the  $\tau$ -pair corresponding to the new symbol in  $\tau(x)$  at position  $i$ , while  $t(j)$  assigns  
241 each black position of  $x$  the position in  $\tau(x)$  where the corresponding new symbol is, i.e.  $b(t(j)) = j$   
242 and  $t(b(i)) = i$ . Thus,

$$1..|\tau(x)| \xrightleftharpoons[t]{b} \mathcal{B}(x)$$

243 In addition, we define  $p$  as the mapping of the  $\tau$ -pairs to the  $\tau$ -alphabet.

244

245 In our running example from Figure 2,  $t(1) = 1$ ,  $t(3) = 2$ ,  $t(4) = 3$ ,  $t(6) = 4$ ,  $t(7) = 5$ , and  
246  $t(9) = 6$ , while  $b(1) = 1$ ,  $b(2) = 3$ ,  $b(3) = 4$ ,  $b(4) = 6$ ,  $b(5) = 7$ , and  $b(6) = 9$ . For the letter mapping,  
247 we get  $p(1,2) = 0$ ,  $p(3,4) = 2$ ,  $p(4,5) = 1$ ,  $p(6,7) = 5$ ,  $p(7,8) = 3$ , and  $p(9,10) = 4$ .

### 248 3.3. Properties Preserved by $\tau$ -Reduction

249 The most important property of  $\tau$ -reduction is a preservation of maximal Lyndon substrings of  
250  $x$ , which start at black positions. This means there is a closed formula that gives, for every maximal  
251 Lyndon substring of  $\tau(x)$ , a corresponding maximal Lyndon substring of  $x$ . Moreover, the formula  
252 for any black position can be computed in constant time. It is simpler to present the following results  
253 using  $\mathcal{L}'$ , the alternative form of Lyndon array, the one where the end positions of maximal Lyndon  
254 substrings are stored rather than their lengths. More formally:

255 **Theorem 1.** Let  $x = x[1..n]$ , let  $\mathcal{L}'_{\tau(x)}[1..m]$  be the Lyndon array of  $\tau(x)$ , and let  $\mathcal{L}'_x[1..n]$  be the Lyndon  
256 array of  $x$ .

257

258 Then for any black  $i \in 1..n$ ,  $\mathcal{L}'_x[i] = \begin{cases} b(\mathcal{L}'_{\tau(x)}[t(i)]) & \text{if } x[b(\mathcal{L}'_{\tau(x)}[t(i)])+1] \preceq x[i] \\ b(\mathcal{L}'_{\tau(x)}[t(i)])+1 & \text{otherwise.} \end{cases}$

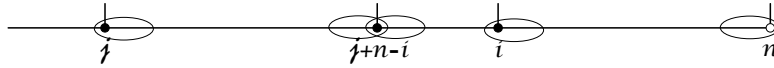
259 The proof of the theorem requires a series of lemmas that are presented below. First we show that  
260  $\tau$ -reduction preserves relationships of certain suffixes of  $x$ .

261 **Lemma 4.** Let  $x = x[1..n]$  and let  $\tau(x) = \tau(x)[1..m]$ . Let  $1 \leq i, j \leq n$ . If  $i$  and  $j$  are both black positions,  
262 then  $x[i..n] \prec x[j..n]$  implies  $\tau(x)[t(i)..m] \prec \tau(x)[t(j)..m]$ .

263 **Proof.** Since  $i$  and  $j$  are both black positions, both  $t(i)$  and  $t(j)$  are defined. Let us assume that  
264  $x[i..n] \prec x[j..n]$ .

- 265 • Case:  $x[i..n]$  is a proper prefix of  $x[j..n]$ .  
266 Then  $j < i$  and so  $x[j..j+n-i] = x[i..n]$  and thus  $x[i..n]$  is a border of  $x[j..n]$ .
- 267 • Case:  $j+n-i$  is black.  
268 Since  $n$  may be black or white, we need to discuss both cases.

- 269 • Case:  $n$  is white.



270

271

272

273

274

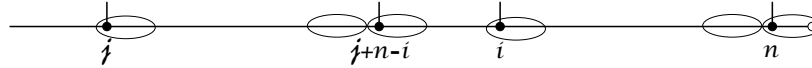
275

Then the last  $\tau$ -pair overlapping  $j..j+n-i$  must be  $(j+n-i-1, j+n-i)$  followed by a  $\tau$ -pair  $(j+n-i, j+n-i+1)$ , and the last  $\tau$ -pair overlapping  $i..n$  must be the last  $\tau$ -pair  $(n-1, n)$ . Thus,  $\tau(x)[t(i)..m] = \tau(x)[t(i)..t(n-1)] = \tau(x)[t(j)..t(j+n-i-1)]$ , and  $\tau(x)[t(j)..m] = \tau(x)[t(j)..t(n-1)]$ , and so  $\tau(x)[t(i)..m]$  is a proper prefix of  $\tau(x)[t(j)..m]$ .

276

277

- Case:  $n$  is black.



278

279

280

281

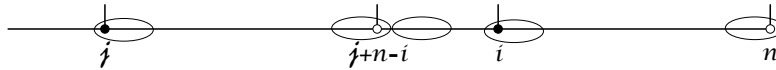
282

Then the last two  $\tau$ -pairs overlapping  $j..j+n-i$  must be  $(j+n-i-1, j+n-i)$  and  $(j+n-i, j+n-i+1)$ , and the last two  $\tau$ -pairs overlapping  $i..n$  must be  $(n-1, n)$  and  $(n, n+1)$ . Thus,  $\tau(x)[t(i)..m] = \tau(x)[t(i)..t(n)] \prec \tau(x)[t(j)..t(j+n-i)]$ , which is a prefix of  $\tau(x)[t(j)..m]$ .

283

284

- Case:  $j+n-i$  is white.



285

286

287

288

289

Then  $n$  must also be white as the  $\tau$ -pair overlapping  $j..j+n-i$  must be  $(j+n-i-1, j+n-i)$  followed by  $(j+n-i+1, j+n-i+2)$ , and the last  $\tau$ -pair overlapping  $i..n$  must be the very last  $\tau$ -pair  $(n-1, n)$ . Then  $\tau(x)[t(i)..m] = \tau(x)[t(i)..t(n-1)] = \tau(x)[t(j)..t(j+n-i-1)]$  which is a proper prefix of  $\tau(x)[t(j)..m]$ .

290

291

- Case:  $x[i] \prec x[j]$ .

292

293

294

Then  $\tau(x)[t(i)] = p(i, i+1)$  and  $\tau(x)[t(j)] = p(j, j+1)$ . Since  $x[i] \prec x[j]$ , we have  $(x[i], x[i+1]) \prec (x[j], x[j+1])$  and so  $p(i, i+1) \prec p(j, j+1)$ , giving  $\tau(x)[t(i)] \prec \tau(x)[t(j)]$  and so  $\tau(x)[t(i)..m] \prec \tau(x)[t(j)..m]$ .

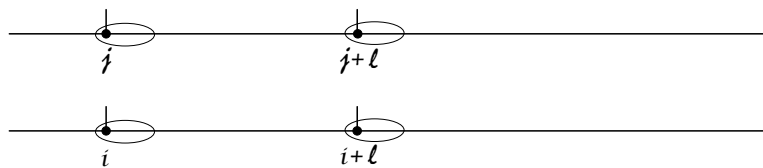
295

296

- Case: for some  $\ell$ ,  $x[i..i+\ell-1] = x[j..j+\ell-1]$  while  $x[i+\ell] \prec x[j+\ell]$ .

297

- Case both  $i+\ell$  and  $j+\ell$  are black.



298

299

Consider  $i+\ell-1$  and  $j+\ell-1$ . Either they are both black, or both are white.

300

301

- ( $\alpha$ ) Case: both  $i+\ell-1$  and  $j+\ell-1$  are black,

302

303

304

305

then  $\tau(x)[t(i)..t(i+\ell-1)] = \tau(x)[t(j)..t(j+\ell-1)]$  and so  $\tau(x)[t(i)..t(i+\ell)] = \tau(x)[t(i)..t(i+\ell-1)]\tau(x)[i+\ell] = \tau(x)[t(j)..t(j+\ell-1)]\tau(x)[i+\ell] \prec \tau(x)[t(j)..t(j+\ell-1)]\tau(x)[j+\ell] = \tau(x)[t(j)..t(j+\ell)]$  and so  $\tau(x)[t(i)..m] \prec \tau(x)[t(j)..m]$ .

306

307

- ( $\beta$ ) Case: both  $i+\ell-1$  and  $j+\ell-1$  are white,

308

309

310

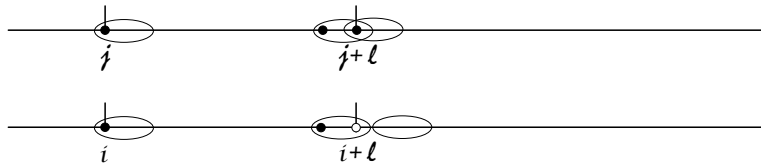
311

then both  $i+\ell-2$  and  $j+\ell-2$  are black. So,  $\tau(x)[t(i)..t(i+\ell-2)] = \tau(x)[t(j)..t(j+\ell-2)]$  and so  $\tau(x)[t(i)..t(i+\ell)] = \tau(x)[t(i)..t(i+\ell-2)]\tau(x)[i+\ell] = \tau(x)[t(j)..t(j+\ell-2)]\tau(x)[i+\ell] \prec \tau(x)[t(j)..t(j+\ell-2)]\tau(x)[j+\ell] = \tau(x)[t(j)..t(j+\ell)]$  and so  $\tau(x)[t(i)..m] \prec \tau(x)[t(j)..m]$ .

312

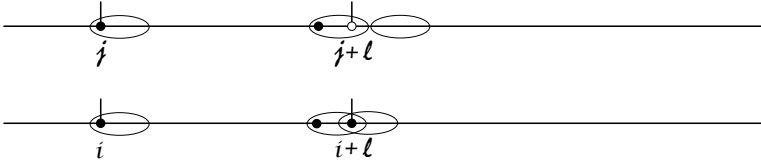


- 313 • Case:  $j+l$  is black and  $i+l$  is white.



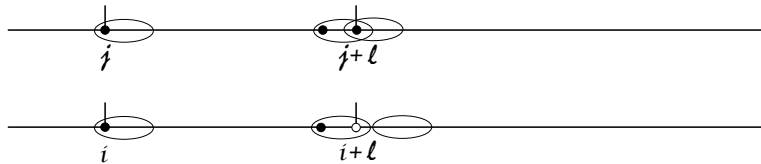
314  
315 Then both  $i+l-1$  and  $j+l-1$  are black, so proceed as in (α).

- 316 • Case:  $j+l$  is white and  $i+l$  is black.



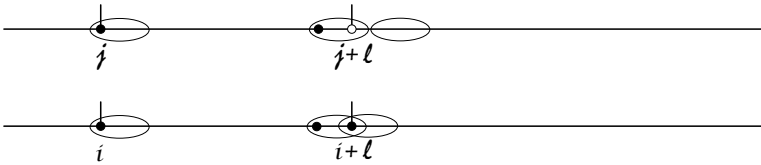
318  
319 Then both  $i+l-1$  and  $j+l-1$  are black, so proceed as in (α).

- 320 • Case: both  $j+l$  and  $i+l$  are white.



322  
323 Then both  $i+l-1$  and  $j+l-1$  are black, so proceed as in (α).

- 324 • Case:  $j+l$  is white and  $i+l$  is black.



326  
327 Then both  $i+l-1$  and  $j+l-1$  are black, so proceed as in (α).

328 □

329 Lemma 5 shows that  $\tau$ -reduction preserves the Lyndon property of certain Lyndon substrings.

330 **Lemma 5.** Let  $x = x[1..n]$  and let  $\tau(x) = \tau(x)[1..m]$ . Let  $1 \leq i < j \leq n$ . Let  $x[i..j]$  be a Lyndon substring of  
331  $x$ , and let  $i$  be a black position.

332  
333 Then  $\begin{cases} \tau(x)[t(i)..t(j)] \text{ is Lyndon} & \text{if } j \text{ is black} \\ \tau(x)[t(i)..t(j-1)] \text{ is Lyndon} & \text{if } j \text{ is white.} \end{cases}$

334 **Proof.** Let us first assume that  $j$  is black.

335 Let  $i_1 = t(i)$ ,  $j_1 = t(j)$  and consider  $k_1$  so that  $i_1 < k_1 \leq j_1$ . Let  $k = b(k_1)$ . Then  $i < k \leq j$   
336 and so  $x[i..n] \prec x[k..n]$  by Lemma 1. Hence,  $\tau(x)[t(i)..m] \prec \tau(x)[t(k)..m]$  by Lemma 4.  
337 Therefore,  $\tau(x)[t(i)..t(j)]$  is Lyndon by Lemma 1.

338 Now, let us assume that  $j$  is white.

339 Then  $j-1$  is black and  $x[i..j-1]$  is Lyndon, so as in the previous case,  $\tau(x)[t(i)..t(j-1)]$  is Lyndon.

341 □

342 Now we can show that  $\tau$ -reduction preserves some maximal Lyndon substrings.

343 **Lemma 6.** Let  $x = x[1..n]$  and let  $\tau(x) = \tau(x)[1..m]$ . Let  $1 \leq i < j \leq n$ . Let  $x[i..j]$  be a maximal Lyndon  
344 substring, and let  $i$  be a black position.

345

346 Then  $\begin{cases} \tau(x)[t(i)..t(j)] \text{ is a maximal Lyndon substring} & \text{if } j \text{ is black} \\ \tau(x)[t(i)..t(j-1)] \text{ is a maximal Lyndon substring} & \text{if } j \text{ is white.} \end{cases}$

347 **Proof.** Since  $x[i..j]$  is maximal Lyndon,  $x[j+1..n] \prec x[i..n]$  by Lemma 1, giving  $x[j+1] \preceq x[i]$ . Since  
348  $x[i..j]$  is Lyndon,  $x[i] \prec x[j]$ . Thus,  $x[j+1] \preceq x[i] \prec x[j]$ .

349

350 We will proceed by discussing two possible cases, one where  $j$  is black, and the other where  $j$  is white.

351 • Case:  $j$  is black.

352 Since  $j$  is black,  $(j, j+1)$  is a  $\tau$ -pair and  $t(j)$  is defined, and by Lemma 5,  $\tau(x)[t(i)..t(j)]$  is Lyndon,  
353 and hence by Lemma 1,  $\tau(x)[t(i)..m] \prec \tau(x)[k..m]$  for any  $t(i) < k \leq t(j)$ . Thus, we must show  
354 the maximality, i.e.  $\tau(x)[t(j)+1..m] \prec \tau(x)[t(i)..m]$ .

355 • Case:  $x[j+1] \preceq x[j+2]$ .

356 Then  $x[j] \succ x[j+1] \preceq x[j+2]$  and so  $j+1$  is black. It follows that  $t(j)+1 = t(j+1)$ . By  
357 Lemma 4,  $\tau(x)[t(j+1)..m] \prec \tau(x)[t(i)..m]$  because  $x[j+1..n] \prec x[i..n]$ , thus  $\tau(x)[t(j)+1..m] \prec$   
358  $\tau(x)[t(i)..m]$ .

359

360 • Case:  $x[j+1] \succ x[j]$ .

361 Then  $x[j] \succ x[i] \succeq x[j+1] \succ x[j+2]$ , then  $\tau(x)[t(i)] = p(i, i+1)$ , and  $\tau(x)[t(j)] = p(j, j+1)$ ,  
362 and  $\tau(x)[t(j)+1] = p(j+2, j+3)$ . It follows that  $\tau(x)[t(j)] = p(j, j+1) \succ \tau(x)[t(i)] = p(i, i+1)$   
363  $\succ \tau(x)[t(j)+1] = p(j+2, i+3)$ . Thus,  $\tau(x)[t(j)+1] \prec \tau(x)[t(i)]$ , and so  $\tau(x)[t(j)+1..m] \prec$   
364  $\tau(x)[t(i)..m]$ .

365

366 • Case:  $j$  is white.

367 By Lemma 5,  $\tau(x)[t(i)..t(j-1)]$  is Lyndon. Since  $j$  is white, it follows that  $j-1$  and  $j+1$  are  
368 black and  $t(j-1)+1 = t(j+1)$ . Since  $x[i..n] \succ x[j+1..n]$ , by Lemma 5 we get  $\tau(x)[t(i)..m] \succ$   
369  $\tau(x)[t(j+1)..m] = \tau(x)[t(j-1)+1..m]$ .

370

□

371 Now we are ready to tackle the proof of Theorem 1.

372

373 **Proof of Theorem 1.**

374

375 Let  $\mathcal{L}'_x[i] = j$  where  $i$  is black. Then  $t(i)$  is defined and  $x[i..j]$  is a maximal Lyndon substring of  $x$ .

376 • Case:  $j$  is black.

377 Then by Lemma 6,  $\tau(x)[t(i)..t(j)]$  is a maximal Lyndon substring of  $\tau(x)$ , hence  $\mathcal{L}'_{\tau(x)}[t(i)] = t(j)$ .  
378 Therefore,  $b(\mathcal{L}'_{\tau(x)}[t(i)]) = b(t(j)) = j = \mathcal{L}'_x[i]$ . Since  $x[i..j]$  is maximal,  $x[j+1] \preceq x[i]$ , i.e.  
379  $x[b(\mathcal{L}'_{\tau(x)}[t(i)])+1] = x[j+1] \preceq x[i]$ .

380

381 • Case:  $j$  is white.

382 Then  $j-1$  is black and the  $\tau(x)[t(j-1)] = p(j-1, j)$ . By Lemma 6,  $\tau(x)[t(i)..t(j-1)]$  is a maximal  
383 Lyndon substring of  $\tau(x)$ , hence  $\mathcal{L}'_{\tau(x)}[t(i)] = t(j-1)$ , so  $b(\mathcal{L}'_{\tau(x)}[t(i)]) = b(t(j-1)) = j-1$ ,  
384 giving  $b(\mathcal{L}'_{\tau(x)}[t(i)]+1) = j$ . Since  $x[i..j]$  is maximal,  $x[i] \prec x[j]$ , i.e.  $x[b(\mathcal{L}'_{\tau(x)}[t(i)]+1)] =$   
385  $x[j] \succ x[i]$ .

386

□

387 3.4. Computing  $\mathcal{L}'_x$  from  $\mathcal{L}'_{\tau(x)}$ 

388 Theorem 1 indicates how to compute the partial  $\mathcal{L}'_x$  from  $\mathcal{L}'_{\tau(x)}$ . The procedure is given in  
 389 Figure 3.

```

for  $i \leftarrow 1$  to  $n$ 
  if  $i = 1$  or  $(x[i-1] \succ x[i] \text{ and } x[i] \preceq x[i+1])$  then
    if  $x[b(\mathcal{L}'_{\tau(x)}[t(i)]+1) \preceq x[i]$  then
       $\mathcal{L}'_x[i] \leftarrow b(\mathcal{L}'_{\tau(x)}[t(i)])$ 
    else
       $\mathcal{L}'_x[i] \leftarrow b(\mathcal{L}'_{\tau(x)}[t(i)]+1$ 
    else
       $\mathcal{L}'_x[i] \leftarrow nil$ 

```

Figure 3. Computing the partial Lyndon array of the input string

390 To compute the missing values, the partial array is processed from right to left. When a missing  
 391 value at position  $i$  is encountered (note that it is recognized by  $\mathcal{L}'_x[i] = nil$ ), the Lyndon array  
 392  $\mathcal{L}'_x[i+1..n]$  is completely filled and also  $\mathcal{L}'_x[i-1]$  is known. Note that  $\mathcal{L}'_x[i+1]$  is the ending position  
 393 of the maximal Lyndon substring starting at the position  $i+1$ . If  $x[i] \succ x[i+1]$ , then the maximal  
 394 Lyndon substring from position  $i+1$  cannot be extended to the left, and hence the maximal Lyndon  
 395 substring at the position  $i$  has a length of 1 and so ends in  $i$ . Otherwise, let us call such a point *hard*,  
 396  $x[i..\mathcal{L}'_x[i+1]]$  is Lyndon and we have to test if the maximal Lyndon substring that immediatly follows  
 397 is lexicographically bigger. If it is, the two substrings can be concatenated to a bigger Lyndon substring  
 398 starting at  $i$ , and so on. But of course, this is all happening inside the maximal Lyndon substring  
 399 starting at  $i-1$  and ending at  $\mathcal{L}'_x[i-1]$  due to Monge property<sup>2</sup> of the maximal Lyndon substrings.

400  
 401 The **while** loop, seen in Figure 4's procedure, is the likely cause of the  $\mathcal{O}(n \log(n))$  complexity.  
 402 At first glance, it may seem that the complexity might be  $\mathcal{O}(n^2)$ , however, the doubling of the length  
 403 of the string when a hard point is introduced actually trims it down to an  $\mathcal{O}(n \log(n))$  worst-case  
 404 complexity. See Subsection 3.5 for more details and Section 7 for the measurements and graphs.

```

 $\mathcal{L}'_x[n] \leftarrow n$ 
for  $i \leftarrow n-1$  downto 2
  if  $\mathcal{L}'[i] = nil$  then
    if  $x[i] \succ x[i+1]$  then
       $\mathcal{L}'[i] \leftarrow i$ 
    else
      if  $\mathcal{L}'[i-1] = i-1$  then
         $stop \leftarrow n$ 
      else
         $stop \leftarrow \mathcal{L}'[i-1]$ 
         $\mathcal{L}'[i] \leftarrow \mathcal{L}'[i+1]$ 
        while  $\mathcal{L}'[i] < stop$  do
          if  $x[i..\mathcal{L}'[i]] \prec x[\mathcal{L}'[i]+1..\mathcal{L}'[\mathcal{L}'[i]+1]]$  then
             $\mathcal{L}'[i] \leftarrow \mathcal{L}'[\mathcal{L}'[i]+1]$ 
          else
            break

```

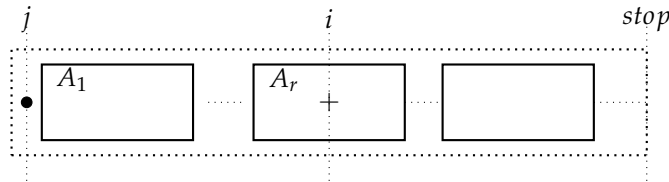
Figure 4. Computing missing values of the Lyndon array of the input string

<sup>2</sup> two maximal Lyndon substring are either disjoint or one completely includes the other

405 Consider our running example from Figure 2. Since  $\tau(x) = 021534$ , we have  
 406  $\mathcal{L}'_{\tau(x)}[1..6] = 6, 2, 6, 4, 6, 6$  giving  $\mathcal{L}'_x[1..9] = 9, \bullet, 3, 9, \bullet, 6, 9, \bullet, 9$ . Computing  $\mathcal{L}'_x[8]$  is easy as  $x[8] = x[9]$   
 407 and so  $\mathcal{L}'_x[8] = 8$ .  $\mathcal{L}'_x[5]$  is more complicate: we can extend the maximal Lyndon substrng from  $\mathcal{L}'_x[6]$   
 408 to the left to 23, but no more, so  $\mathcal{L}'_x[5] = 6$ . Computing  $\mathcal{L}'_x[2]$  is again easy as  $x[2] = x[3]$  and so  
 409  $\mathcal{L}'_x[2] = 2$ . Thus  $\mathcal{L}'_x[1..9] = 9, 2, 3, 9, 6, 6, 9, 8, 9$ .

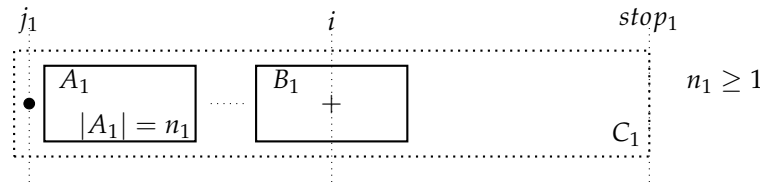
410 3.5. The Complexity of TRLA

411 To determine the complexity of the algorithm, we attach to each position  $i$  a counter  $red[i]$   
 412 initialized to 0. When computing a missing value  $\mathcal{L}'_x[j]$  with a configuration shown below where  
 413  $stop = \begin{cases} \mathcal{L}'_x[j-1] & \text{if } \mathcal{L}'_x[j-1] > j-1 \\ n & \text{otherwise} \end{cases}$ , we have to check if  $A_1$  can be extended by  $j$  to a Lyndon  
 414 substrng (i.e. if  $x[j] \preceq x[j+1]$ ). If so, we have to compare  $jA_1$  with  $A_2$  and if  $jA_1A_2$  is Lyndon (i.e.  
 415  $jA_1 \prec A_2$ ), then we must check if  $jA_1A_2A_3$  is Lyndon (i.e.  $jA_1A_2 \prec A_3$ ) so on and so forth. In general,  
 416 we are checking if  $jA_1..A_r$  is Lyndon (i.e.  $jA_1A_2..A_{r-1} \prec A_r$ ). When comparing the Lyndon substrng  
 417  $jA_1..A_{r-1}$  with  $A_r$ , at every position  $i$  of  $A_r$ , we increment the counter  $red[i]$ . When done, the value of  
 418  $red[i]$  represents how many how many comparisons were done at the position  $i$ .



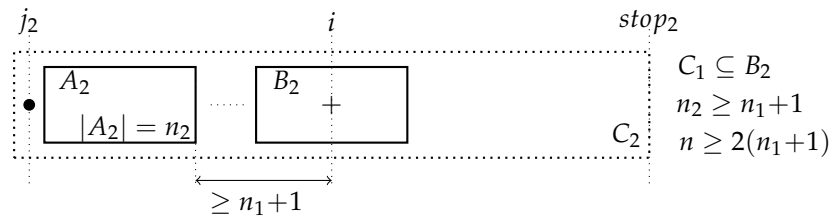
419  
 420 Consider a position  $i$  that was used  $k$  times for  $k \geq 4$ , i.e.  $red[i] = k$ . The next diagram indicates  
 421 the configuration when the counter  $red[i]$  was incremented for the 1st time in the comparison of  $j_1A_1..$   
 422 and  $B_1$  during the computation of the missing value  $\mathcal{L}'_x[j_1]$  where

423  $stop_1 = \begin{cases} \mathcal{L}'_x[j_1-1] & \text{if } \mathcal{L}'_x[j_1-1] > j_1-1 \\ n & \text{otherwise} \end{cases}$



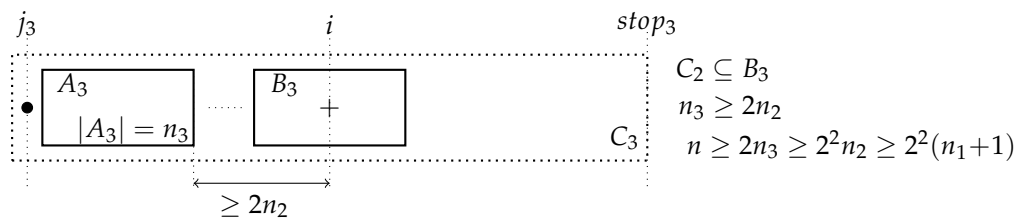
424  
 425 The next diagram indicates the configuration when the counter  $red[i]$  was incremented for the 2nd  
 426 time in the comparison of  $j_2A_2..$  and  $B_2$  during the computation of the missing value  $\mathcal{L}'_x[j_2]$  where

427  $stop_2 = \begin{cases} \mathcal{L}'_x[j_2-1] & \text{if } \mathcal{L}'_x[j_2-1] > j_2-1 \\ n & \text{otherwise} \end{cases}$



428  
 429 The next diagram indicates the configuration when the counter  $red[i]$  was incremented for the 3rd  
 430 time in the comparison of  $j_3A_3..$  and  $B_3$  during the computation of the missing value  $\mathcal{L}'_x[j_3]$  where

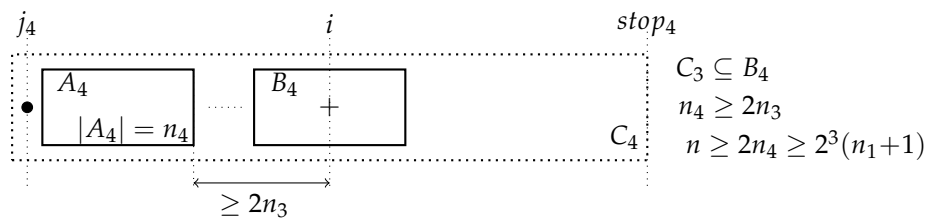
431  $stop_3 = \begin{cases} \mathcal{L}'_x[j_3-1] & \text{if } \mathcal{L}'_x[j_3-1] > j_3-1 \\ n & \text{otherwise} \end{cases}$



432

433 The next diagram indicates the configuration when the counter  $red[i]$  was incremented for the 4th  
 434 time in the comparison of  $j_4 A_4 \dots$  and  $B_4$  during the computation of the missing value  $\mathcal{L}'_x[j_4]$  where

$$435 \quad stop_4 = \begin{cases} \mathcal{L}'_x[j_4 - 1] & \text{if } \mathcal{L}'_x[j_4 - 1] > j_4 - 1 \\ n & \text{otherwise} \end{cases}$$



436

437 Thus, if  $red[i] = k$ , then  $n \geq 2^{k-1}(n_1 + 1) \geq 2^k$  as  $n_1 + 1 \geq 2$ . Thus,  $n \geq 2^k$  and so  $k \leq \log(n)$ . Thus,  
 438 either  $k < 4$  or  $k \leq \log(n)$ . Therefore, the overall complexity is  $\mathcal{O}(n \log(n))$ .

439 To show that the average case complexity is linear, we first recall that the overall complexity of  
 440 *TRLA* is determined by the procedure filling the missing values. We showed above that there are at  
 441 most  $\log(n)$  missing values that cannot be determined in constant time. We overestimate the number  
 442 of strings of length  $n$  over an alphabet of size  $C$ ,  $2 \leq C \leq n$ , that will force non-linear computation,  
 443 by assuming that every possible  $\log(n)$  subset of indices with any possible letters assignment forces  
 444 worst performance. Thus, there are  $C^n - \binom{n}{\log(n)} C^{\log(n)}$  strings that are processed in linear time, say  
 445 with a constant  $K_1$ , and there are  $\binom{n}{\log(n)} C^{\log(n)}$  strings that are processed in the worst time, with a  
 446 constant  $K_2$ . Let  $K = \max(K_1, K_2)$ . Then the average time is bounded by

447

$$\begin{aligned} & \frac{(C^n - \binom{n}{\log(n)} C^{\log(n)}) K n + \binom{n}{\log(n)} C^{\log(n)} K n \log(n)}{C^n} = \\ & K n + K n \frac{\binom{n}{\log(n)} C^{\log(n)}}{C^n} (\log(n) - 1) \leq \\ & K n + K n \frac{\binom{n}{\log(n)} C^{\log(n)}}{C^n} \log(n) \leq \\ & K n + K n \frac{n^{\log(n)} C^{\log(n)}}{\log(n)! C^n} \log(n) \leq \\ & K n + K n \frac{n^{2\log(n)}}{2^n} \leq \\ & K n + K n = 2K n \end{aligned}$$

448 for  $n \geq 2^7$ . The last step follows from the fact that  $n^{2\log(n)} \leq 2^n$  for any  $n \geq 2^7$ .

449 The combinatorics of the processing is too complicated to ascertain whether the worst-time  
 450 complexity is linear or not. We tried to generate strings that might give the worst performance. We  
 451 used three different formulas to generate the strings, nesting the white indices that might require  
 452 non-constant computation: the data set `extreme_trla` of binary strings is created using the recursive  
 453 formula  $u_{k+1} = 00u_k0u_k$ , using the first 100 shortest binary Lyndon words as the start  $u_0$ . The moment

454 the size  $u_k$  exceeds the required length of the string, the recursion stops and the string is trimmed to  
 455 the required length. For `extreme_tr1a1` data set, we used the same approach with the formula  $u_{k+1} =$   
 456  $000u_k00u_k$  and for `extreme_tr1a2` data set, we used the formula  $u_{k+1} = 0000u_k00u_k$ .

457 The space complexity of our C++ implementation is bounded by  $9n$  integers. This upper bound  
 458 is derived from the fact that a Tau object (see `Tau.hpp`, [16]) requires  $3n$  integers of space for a  
 459 string of length  $n$ . So the first call to `TRLA` requires  $3n$ , the next recursive call requires at most  
 460  $3\frac{2}{3}n$ , the next recursive call requires at most  $3(\frac{2}{3})^2n$ , ... Thus,  $3n + 3\frac{2}{3}n + 3(\frac{2}{3})^2n + 3(\frac{2}{3})^3n + \dots =$   
 461  $3n(1 + \frac{2}{3} + (\frac{2}{3})^2 + (\frac{2}{3})^3 + (\frac{2}{3})^4 + \dots) = 3n\frac{1}{1-\frac{2}{3}} = 9n$ . However, it should be possible to bring it down  
 462 to  $6n$  integers.

#### 463 4. The Algorithm *BSLA*

464 The input strings for *BSLA* are tight strings over integer alphabets. Note that this requirement  
 465 does not detract from the applicability of the algorithm as any string over an integer alphabet can  
 466 easily be transformed in  $\mathcal{O}(|x|)$  time to a tight string so that the original string and the transformed  
 467 string have the same Lyndon array. Thus, computing the Lyndon array for the transformed string  
 468 gives also the Lyndon array for the original string.

469 The algorithm is based on a refinement of a list of groups of indices of the input string  $x$ . The  
 470 refinement is driven by a group that is already complete and the refinement process makes the  
 471 immediately preceding group also complete. In turn, this newly completed group is used as the driver  
 472 of the next round of the refinement. In this fashion, the refinement proceeds from right to left until all  
 473 the groups in the list are complete. The initial list of groups consists of the groups of indices with the  
 474 same alphabet symbol.  
 475

476 Each group is assigned a specific substring of the input string referred to as the *context* of the  
 477 group. Every index in a group has the property that there is an occurrence of the group's context at that  
 478 position. Throughout the process the list of the groups is maintained in an increasing lexicographic  
 479 order by their contexts. Moreover, at every stage, the contexts of all the groups are Lyndon substrings  
 480 of  $x$  with an additional property that the contexts of complete groups are maximal Lyndon substrings.  
 481 Hence, when the refinement is complete, the contexts of all the groups in the list represent all maximal  
 482 Lyndon substrings of  $x$ .  
 483

484 In order to verify the process and prove it correct, and to be able to describe the refinement in  
 485 technical details, we must formally define several notions and properties in the next section.  
 486

##### 487 4.1. Notation and Basic Notions of *BSLA*

488 For the sake of simplicity, we fix a string  $x = x[1..n]$  for the whole Subsection 4.1; all the definitions  
 489 and the observations apply and refer to this  $x$ .

490 A *group*  $G$  is a non-empty set of indices of  $x$ . The group  $G$  is assigned a *context*, i.e. a substring  
 491  $\text{con}(G)$  of  $x$  with the property that for any  $i \in G$ ,  $x[i..i+|\text{con}(G)|-1] = \text{con}(G)$ . If  $i \in G$ , then  $C(i)$   
 492 denotes the occurrence of the context of  $G$  at the position  $i$ , i.e.  $C(i) = x[i..i+|\text{con}(G)|-1]$ . We say that  
 493 a group  $G'$  is *smaller than* or *precedes* a group  $G''$  if  $\text{con}(G') \prec \text{con}(G'')$ .

494 **Definition 1.** An ordered list of groups  $\langle G_k, G_{k-1}, \dots, G_2, G_1 \rangle$  is a *group configuration* if

- 495 (C<sub>1</sub>)  $G_k \cup G_{k-1} \cup \dots \cup G_2 \cup G_1 = 1..n$ ;
- 496 (C<sub>2</sub>)  $G_j \cap G_\ell = \emptyset$  for any  $1 \leq \ell < j \leq k$ ;
- 497 (C<sub>3</sub>)  $\text{con}(G_k) \prec \text{con}(G_{k-1}) \prec \dots \prec \text{con}(G_2) \prec \text{con}(G_1)$ ;
- 498 (C<sub>4</sub>) For any  $j \in 1..k$ ,  $\text{con}(G_j)$  is a Lyndon substring of  $x$ .



500 Note that  $(C_1)$  and  $(C_2)$  guarantee that  $\langle G_k, G_{k-1}, \dots, G_2, G_1 \rangle$  is a disjoint partitioning of  $1..n$ . For  
 501  $i \in n$ ,  $gr(i)$  denotes the unique group to which  $i$  belongs, i.e. if  $i \in G_t$ , then  $gr(i) = G_t$ . Note that  
 502 using this notation,  $C(i) = x[i..i+|con(gr(i))|-1]$ .

503 The mapping  $prev$  is defined by  $prev(i) = \max\{j < i : con(gr(j)) \prec con(gr(i))\}$  if such  $j$  exists,  
 504 otherwise  $prev(i) = nil$ .

505 For a group  $G$  from a group configuration, we define an equivalence  $\sim$  on  $G$  as follows:  $i \sim j$  iff  
 506  $gr(prev(i)) = gr(prev(j))$  or  $prev(i) = prev(j) = nil$ . The symbol  $[i]_{\sim}$  denotes the class of equivalence  
 507  $\sim$  that contains  $i$ , i.e.  $[i]_{\sim} = \{j \in G \mid j \sim i\}$ . If  $prev(i) = nil$ , then the class  $[i]_{\sim}$  is called trivial. An  
 508 interesting observation states that if  $G$  is viewed as an ordered set of indices, then a non-trivial  $[i]_{\sim}$  is  
 509 an interval:  
 510

511 **Observation 7.** Let  $G$  be a group from a group configuration for  $x$ . Consider an  $i \in G$  such that  $prev(i) \neq nil$ .  
 512 Let  $j_1 = \min[i]_{\sim}$  and  $j_2 = \max[i]_{\sim}$ . Then  $[i]_{\sim} = \{j \in G \mid j_1 \leq j \leq j_2\}$ .

513 **Proof.** Since  $prev(j_1)$  is a candidate to be  $prev(j)$ ,  $prev(j) \neq nil$  and  $prev(j_1) \leq prev(j) \leq prev(j_2) =$   
 514  $prev(j_1)$ , so  $prev(j) = prev(j_1) = prev(j_2)$ .  
 515 □

516 On each non-trivial class of  $\sim$ , we define a relation  $\approx$  as follows:  $i \approx j$  iff  $|j-i| = |con(G)|$ ; in  
 517 simple terms it means that the occurrence  $C(i)$  of  $con(G)$  is immediately followed by the occurrence  $C(j)$  of  
 518  $con(G)$ . The transitive closure of  $\approx$  is a relation of equivalence, which we also denote by  $\approx$ . The  
 519 symbol  $[i]_{\approx}$  denotes the class of equivalence  $\approx$  containing  $i$ , i.e.  $[i]_{\approx} = \{j \in [i]_{\sim} \mid j \approx i\}$ .  
 520

521 For each  $j$  from a non-trivial  $[i]_{\sim}$ , we define the *valence* by  $val(j) = |[i]_{\approx}|$ . In simple terms,  $val(i)$   
 522 is the number of elements from  $[i]_{\sim}$  that are  $\approx i$ . Thus,  $1 \leq val(i) \leq |G|$ .  
 523

524 Interestingly, if  $G$  is viewed as an ordered set of indices, then  $[i]_{\approx}$  is a subinterval of the interval  
 525  $[i]_{\sim}$ :

526 **Observation 8.** Let  $G$  be a group from a group configuration for  $x$ . Consider an  $i \in G$  such that  $prev(i) \neq nil$ .  
 527 Let  $j_1 = \min[i]_{\approx}$  and  $j_2 = \max[i]_{\approx}$ . Then  $[i]_{\approx} = \{j \in [i]_{\sim} \mid j_1 \leq j \leq j_2\}$ .

528 **Proof.** Argue by contradiction. Assume that there is an  $j \in [i]_{\sim}$  so that  $j_1 < j < j_2$  so that  $j \notin [i]_{\approx}$ .  
 529 Take the minimal such  $j$ . Consider  $j' = j - |con(G)|$ . Then  $j' \in [i]_{\sim}$  and since  $j' < j$ ,  $j' \in [i]_{\approx}$  due to the  
 530 minimality of  $j$ . So  $i \approx j' \approx j$  and so  $j \approx i$ , a contradiction.  
 531 □

532 **Definition 2.** A group  $G$  is *complete* if for any  $i \in G$ , the occurrence  $C(i)$  of  $con(G)$  is a maximal Lyndon  
 533 substring of  $x$ .

534 A group configuration  $\langle G_k, G_{k-1}, \dots, G_2, G_1 \rangle$  is *t-complete*,  $1 \leq t \leq k$ , if

- 535 (C<sub>5</sub>) the groups  $G_t, \dots, G_1$  are complete;  
 536 (C<sub>6</sub>) the mapping  $prev$  is **proper** on  $G_t$ :  
 537 for any  $i \in G_t$ , if  $prev(i) \neq nil$  and  $v = val(i)$ , then there are  $i_1, \dots, i_v \in G_t$ ,  
 538  $i \in \{i_1, \dots, i_v\}$ ,  $prev(i) = prev(i_1) = \dots = prev(i_v)$ , and so that  $C(prev(i))C(i_1)\dots C(i_v)$   
 539 is a prefix of  $x[j..n]$ ;  
 540 (C<sub>7</sub>) the family  $\{C(i) \mid i \in 1..n\}$  is **proper**:  
 541 (a) if  $C(j)$  is proper substring of  $C(i)$ , i.e.  $C(j) \subsetneq C(i)$ , then  $con(G_t) \prec con(gr(j))$ ,  
 542 (b) if  $C(i)$  is followed immediately by  $C(j)$ , i.e. when  $i+|con(gr(i))| = j$ , and  
 543  $C(i) \prec C(j)$ , then  $con(gr(j)) \preceq con(G_t)$ ;

544 (C<sub>8</sub>) the family  $\{C(i) \mid i \in 1..n\}$  has the **Monge** property, i.e. if  $C(i) \cap C(j) \neq \emptyset$ , then  $C(i) \subseteq C(j)$   
 545 or  $C(j) \subseteq C(i)$ .

546 The condition (C<sub>6</sub>) is all-important for carrying out the refinement process (see (R<sub>3</sub>) below). The  
 547 conditions (C<sub>7</sub>) and (C<sub>8</sub>) are necessary for asserting that the condition (C<sub>6</sub>) is preserved during the  
 548 refinement process.

#### 549 4.2. The Refinement

550 For the sake of simplicity, we fix a string  $x = x[1..n]$  for the whole Subsection 4.2; all the definitions,  
 551 lemmas and theorems apply and refer to this  $x$ .

552 **Lemma 9.** Let  $A_x = \{a_1, \dots, a_k\}$  and  $a_1 \prec a_2 \prec \dots \prec a_k$ . For  $1 \leq \ell \leq k$ , define  $G_\ell = \{i \in 1..n : x[i] =$   
 553  $a_{k+1-\ell}\}$  with context  $a_{k+1-\ell}$ . Then  $\langle G_k, \dots, G_1 \rangle$  is a 1-complete group configuration.

554 **Proof.** (C<sub>1</sub>), (C<sub>2</sub>), (C<sub>3</sub>), and (C<sub>4</sub>) are straightforward to verify. To verify (C<sub>5</sub>), we need to show that  $G_1$   
 555 is complete. Any occurrence of  $a_k$  in  $x$  is a maximal Lyndon substring, so  $G_1$  is complete.

556 To verify (C<sub>6</sub>), consider  $j = \text{prev}(i)$  and  $\text{val}(i) = v$  for  $i \in G_1$ . Consider any  $j < \ell < i$ . If  $x[\ell] \neq a_k$ ,  
 557 then  $\text{prev}(i) < \ell$  which contradicts the definition of  $\text{prev}$ . Hence  $x[\ell] = a_k$  and so  $x[j+1] = \dots = x[i] =$   
 558  $\dots x[j+v+1] = a_k$  while  $x[j] = a_\ell$  for some  $\ell < k$ . It follows that  $x[j..n]$  has  $a_\ell(a_k)^v$  as a prefix.

559 The condition (C<sub>7</sub>(a)) is trivially satisfied as no  $C(i)$  can have a proper substring. If  $C(i)$  is immediately  
 560 followed by  $C(j)$  and  $C(i) \prec C(j)$ , then  $C(i) = x[i], j = i+1, C(j) = x[i+1]$  and  $x[i] \prec x[i+1]$ . Then  
 561  $\text{con}(C(j)) = x[i+1] \preceq a_k = \text{con}(G_1)$ , so (C<sub>7</sub>(b)) is also satisfied.

562 To verify (C<sub>8</sub>), consider  $C(i) \cap C(j) \neq \emptyset$ . Then  $C(i) = x[i] = x[j] = C(j)$ .

563 □

564 Let  $\langle G_k, \dots, G_t, \dots, G_1 \rangle$  by a  $t$ -complete group configuration. The refinement is driven by the group  $G_t$   
 565 and it might only partition the groups that precede it, i.e. the groups  $G_k, \dots, G_{t+1}$ , while the groups  
 566  $G_t, \dots, G_1$  remain unchanged.

567 (R<sub>1</sub>) Partition  $G_t$  into classes of the equivalence  $\sim$ .

568  $G_t = [i_1]_{\sim} \cup [i_2]_{\sim} \cup \dots \cup [i_p]_{\sim} \cup X$  where  $X = \{i \in G_t : \text{prev}(i) = \text{nil}\}$  may be possibly  
 569 empty and  $i_1 < i_2 < \dots < i_p$ .

570 (R<sub>2</sub>) Partition every class  $[i_\ell]_{\sim}, 1 \leq \ell \leq p$ , into classes of the equivalence  $\approx$ .

571  $[i_\ell]_{\sim} = [j_{\ell,1}]_{\approx} \cup [j_{\ell,2}]_{\approx} \cup \dots \cup [j_{\ell,m_\ell}]_{\approx}$  where  $\text{val}(j_{\ell,1}) < \text{val}(j_{\ell,2}) < \dots < \text{val}(j_{\ell,m_\ell})$ .

572 (R<sub>3</sub>) So we have a list of classes in this order:  $[j_{1,1}]_{\approx}, [j_{1,2}]_{\approx}, \dots, [j_{1,m_1}]_{\approx}, [j_{2,1}]_{\approx}, [j_{2,2}]_{\approx}, \dots, [j_{2,m_2}]_{\approx},$   
 573  $\dots, [j_{p,1}]_{\approx}, [j_{p,2}]_{\approx}, \dots, [j_{p,m_p}]_{\approx}$ . This list is processed from left to right. Note that for each  
 574  $i \in [j_{\ell,k}]_{\approx}, \text{prev}(i) \in \text{gr}(j_{\ell,k})$  and  $\text{val}(i) = \text{val}(j_{\ell,k})$ .

575 For each  $j_{\ell,k}$ , move all elements  $\{\text{prev}(i) : i \in [j_{\ell,k}]_{\approx}\}$  from the group  $\text{gr}(\text{prev}(j_{\ell,k}))$  into  
 576 a new group  $H$  and place  $H$  in the list of groups right after the group  $\text{gr}(\text{prev}(j_{\ell,k}))$  and  
 577 set its context to  $\text{con}(\text{gr}(\text{prev}(j_{\ell,k})))\text{con}(\text{gr}(j_{\ell,k}))^{\text{val}(j_{\ell,k})}$ . (Note, that this "doubling of the  
 578 contexts" is possible due to (C<sub>6</sub>)). Then update  $\text{prev}$ :

579 all values of  $\text{prev}$  are correct except possibly the values of  $\text{prev}$  for indices from  
 580  $H$ . It may be the case that for  $i \in H$ , there is  $i' \in \text{gr}(j_{\ell,k})$  so that  $\text{prev}(i) < i'$ , so  
 581  $\text{prev}(i)$  must be reset to maximal such  $i'$ . (Note that before the removal of  $H$  from  
 582  $\text{gr}(j_{\ell,k})$ , the index  $i'$  was not eligible to be considered for  $\text{prev}(i)$  as  $i$  and  $i'$  were both  
 583 from the same group.)

584 Theorem 2 shows that having a  $t$ -complete group configuration  $\langle G_k, \dots, G_{t+1}, G_t, \dots, G_1 \rangle$  and refining  
 585 it by  $G_t$ , then the resulting system of groups is a  $(t+1)$ -complete group configuration. This allows to  
 586 carry on the refinement in an iterative fashion.

587 **Theorem 2.** Let  $Conf = \langle G_k, \dots, G_{t+1}, G_t, \dots, G_1 \rangle$  be a  $t$ -complete group configuration,  $1 \leq t$ . After performing  
 588 the refinement of  $Conf$  by group  $G_t$ , the resulting system of groups denoted as  $Conf'$  is a  $(t+1)$ -complete group  
 589 configuration.

590 **Proof.** We carry the proof in a series of claims. The symbols  $gr()$ ,  $con()$ ,  $C()$ ,  $prev()$ , and  $val()$   
 591 denote the functions for  $Conf$ , while  $gr'()$ ,  $con'()$ ,  $C'()$ ,  $prev'()$ , and  $val'()$  denote the functions for  $Conf'$ .

592  
 593 **Claim 1.**  $Conf'$  is a group configuration, i.e.  $(C_1)$ ,  $(C_2)$ ,  $(C_3)$  and  $(C_4)$  for  $Conf'$  hold.

594 *Proof of Claim 1.*

595  $(C_1)$  and  $(C_2)$  follow from the fact that the process is a refinement, i.e. a group is either preserved as is,  
 596 or is partitioned into two or more groups. The doubling of the contexts in step  $(R_3)$  guarantees that the  
 597 increasing order of the contexts is preserved, i.e.  $(C_3)$  holds. For any  $j \in G_t$  so that  $j = prev(i) \neq nil$ ,  
 598  $con(gr(prev(j)))$  is Lyndon and  $con(gr(j))$  is also Lyndon, and  $con(gr(prev(j))) \prec con(gr(j))$ , so  
 599  $con(gr(prev(j)))con(gr(j))^{val(j)}$  is Lyndon as well and thus  $(C_4)$  holds. That concludes the proof of  
 600 Claim 1.

601 **Claim 2.**  $\{C'(i) \mid i \in 1..n\}$  is proper and has the Monge property, i.e.  $(C_7)$  and  $(C_8)$  for  $Conf'$  hold.

602 *Proof of Claim 2.*

603 Consider  $C'(i)$  for some  $i \in 1..n$ . There are two possibilities:

- 604 •  $C'(i) = C(i)$ , or
- 605 •  $C'(i) = C(i)C(i_1)...C(i_v)$ , for some  $i_1, i_2, \dots, i_v \in G_t$ , so that for any  $1 \leq \ell \leq v$ ,  $i = prev(i_\ell)$ ,  
 606 and  $C(i_\ell) = con(G_t)$ ,  $v = val(i_\ell)$ , and for any  $1 \leq \ell < k$ , and  $i_{\ell+1} = i_\ell + |con(G_t)|$ . Note that  
 607  $con(gr(i)) \prec con(G_t)$ .

608 Consider  $C'(i)$  and  $C'(j)$  for some  $1 \leq i < j \leq n$ .

- 610 • Case  $C'(i) = C(i)$  and  $C'(j) = C(j)$ .
  - 611 • Show that  $(C_7(a))$  holds.  
 612 If  $C'(j) \subsetneq C'(i)$ , then  $C(j) \subsetneq C(i)$ , and so by  $(C_7(a))$  for  $Conf$ ,  
 613  $con(G_t) \prec con(gr(j))$ , and thus  $con'(H_{t+1}) \prec con(G_t) \prec con(gr(j)) = con'(gr'(j))$ .  
 614 Therefore,  $(C_7(a))$  for  $Conf'$  holds.
  - 615 • Show that  $(C_8)$  holds. If  $C'(i) \cap C'(j) \neq \emptyset$ , then  $C(i) \cap C(j) \neq \emptyset$ , so  $C(j) \subseteq C(i)$ , and so  
 616  $C'(j) \subseteq C'(i)$ , so  $(C_8)$  for  $Conf'$  holds.
- 617  
 618 • Case  $C'(i) = C(i)$  and  $C'(j) = C(j)C(j_1)..C(j_w)$ ,  
 619 where  $w = val(j_1)$ ,  $C(j_1) = \dots = C(j_w) = con(G_t)$ , and  $j_1 \approx \dots \approx j_w$ .
  - 620 • Show that  $(C_7(a))$  holds.  
 621 If  $C'(j) \subsetneq C'(i)$ , then  $C(j)C(j_1)..C(j_w) \subsetneq C(i)$ , hence  $C(j) \subsetneq C(i)$ , and so by  $(C_7(a))$  for  $Conf$ ,  
 622  $con(G_t) \prec con(gr(j))$ . By  $t$ -completeness of  $Conf$ ,  $C(j)$  is a maximal Lyndon substring, a  
 623 contradiction with  $C(j)C(j_1)..C(j_w)$  being Lyndon. This is an impossible case.
  - 624 • Show that  $(C_8)$  holds.  
 625 If  $C'(i) \cap C'(j) \neq \emptyset$ , then  $C(j) \subseteq C(i)$  by  $(C_8)$  for  $Conf$ . By  $(C_7(a))$  for  $Conf$ ,  $C(j)$  cannot be  
 626 a suffix of  $C(i)$  as  $con(gr(j)) \prec con(G_t)$ . Hence  $C(i) \cap C(j_1) \neq \emptyset$ , and so  $C(j)C(j_1) \subseteq C(i)$   
 627 and since  $C(j_1)$  cannot be a suffix of  $C(i)$  as  $gr(j_1) = G_t$ , it follows that  $C(i) \cap C(j_2) \neq \emptyset$ , ...,  
 628 ultimately giving  $C(j)C(j_1)..C(j_w) \subseteq C(i)$ . So  $(C_8)$  for  $Conf'$  holds.

629

- 630 • Case  $C'(i) = C(i)C(i_1)..C(i_v)$  and  $C'(j) = C(j)$ ,  
 631 where  $v = \text{val}(i_1)$ ,  $C(i_1) = \dots = C(i_v) = \text{con}(G_t)$ , and  $i_1 \approx \dots \approx i_v$ .
- 632 • Show that  $(C_7(a))$  holds.  
 633 If  $C'(j) \subsetneq C'(i)$ , then either  $C(j) \subsetneq C(i)$ , which implies by  $(C_7(a))$  for  $\text{Conf}$  that  $\text{con}(G_t) \prec$   
 634  $\text{con}(\text{gr}(j))$ , giving  $\text{con}'(H_{t+1}) \prec \text{con}'(G_t) = \text{con}(G_t) \prec \text{con}(\text{gr}(j)) = \text{con}'(\text{gr}'(j))$ , or  $C(j) \subseteq$   
 635  $C(i_\ell)$  for some  $1 \leq \ell \leq v$ . If  $C(j) = C(i_\ell)$ , then  $\text{gr}(j) = \text{gr}(i_\ell) = G_t$ , giving  $\text{con}'(H_{t+1}) \prec$   
 636  $\text{con}(G_t) = \text{con}(\text{gr}(j))$ . So  $(C_7(a))$  for  $\text{Conf}'$  holds.
- 637 • Show that  $(C_8)$  holds.  
 638 Let  $C'(i) \cap C'(j) \neq \emptyset$ . Consider  $\mathcal{D} = \{i_\ell \mid 1 \leq \ell \leq v \ \& \ C(j) \cap C(i_\ell) \neq \emptyset\}$ .  
 639 Assume that  $\mathcal{D} \neq \emptyset$ :  
 640 By  $(C_8)$  for  $\text{Conf}$ , either  $C(j) \subseteq \bigcup_{i_\ell \in \mathcal{D}} C(i_\ell) \subseteq C'(i)$  and we are done, or  $\bigcup_{i_\ell \in \mathcal{D}} C(i_\ell) \subseteq$   
 641  $C(j)$ . Let  $i_k$  be the smallest element of  $\mathcal{D}$ . Since  $C(i_k)$  cannot be prefix of  $C(j)$ , it means  
 642 that  $i_k = i_1$ . Since  $C(i_1)$  cannot be a prefix of  $C(j)$ , it means that  $C(i) \cap C(j) \neq \emptyset$ , and so  
 643  $C(j) \subseteq C(i)$ , which contradicts the fact that  $C(j) \subseteq \bigcup_{i_\ell \in \mathcal{D}} C(i_\ell) \subseteq C'(i)$ .
- 644 Assume that  $\mathcal{D} = \emptyset$ :  
 645 Then  $C(i) \cap C(j) \neq \emptyset$ , and so by  $(C_8)$  for  $\text{Conf}$ ,  $C(j) \subseteq C(i) \subseteq C'(i)$  as  $i < j$ .
- 646 • Case  $C'(i) = C(i)C(i_1)..C(i_v)$  and  $C'(j) = C(j)C(j_1)..C(j_w)$ ,  
 647 where  $v = \text{val}(i_1)$ ,  $C(i_1) = \dots = C(i_v) = \text{con}(G_t)$ , and  $i_1 \approx \dots \approx i_v$ , and where  $v = \text{val}(j_1)$ ,  
 648  $C(j_1) = \dots = C(j_w) = \text{con}(G_t)$ , and  $j_1 \approx \dots \approx j_w$ .
- 649 • Show that  $(C_7(a))$  holds.  
 650 Let  $C'(j) \subsetneq C'(i)$ . Then either  $C(j) \subseteq C(i)$  and so  $\text{con}(G_t) \prec \text{con}(\text{gr}(j))$ , implying that  $C(j)$   
 651 is maximal contradicting  $C(j)C(j_1)..C(j_w)$  being Lyndon. Thus  $C(j) \subsetneq C(i_\ell)$  for some  $1 \leq$   
 652  $\ell \leq v$ . But then  $\text{con}(G_t) \prec \text{con}(\text{gr}(j))$ , implying that  $C(j)$  is maximal, again a contradiction.  
 653 This is an impossible case.
- 654 • Show that  $(C_8)$  holds.  
 655 Let  $C'(i) \cap C'(j) \neq \emptyset$ . Let us first assume that  $C(i) \cap C(j) \neq \emptyset$ . Then  $C(j) \subseteq C(i)$ . Since  
 656  $C(j)$  cannot be a suffix of  $C(i)$ , it follows that  $C(i) \cap C(j_1) \neq \emptyset$ . Therefore,  $C(j)C(j_1) \subseteq C(i)$ .  
 657 Repeating this argument leads to  $C(j)C(j_1)..C(j_w) \subseteq C(i)$  and we are done.  
 658 So assume that  $C(i) \cap C(j) = \emptyset$ . Let  $1 \leq \ell \leq v$  be the smallest such that  $C(i_\ell) \cap C(j) \neq \emptyset$ .  
 659 Such  $\ell$  must exist. Then  $i_\ell \leq j$ . If  $i_\ell = j$ , then either  $C(i_\ell)$  is a prefix of  $C(j)$  or vice versa,  
 660 both impossibilities, hence  $i_\ell < j$ . Repeating the same arguments as for  $i$ , we get that  
 661  $C(j)C(j_1)..C(j_w) \subseteq C(i_\ell)$  and so we are done.

662 It remains to show that  $(C_7(b))$  for  $\text{Conf}'$  holds.

663

664 Consider  $C'(i)$  immediately followed by  $C'(j)$  with  $C'(i) \prec C'(j)$ .

- 665 • Assume that  $\text{gr}'(j) \in \{G_{t-1}, \dots, G_1\}$ .  
 666 Then  $\text{con}(G_t) = \text{con}'(G_t)$ ,  $\text{gr}(j) = \text{gr}'(j)$  and  $\text{con}(\text{gr}(j)) = \text{con}'(\text{gr}'(j))$ . If  $C'(i) = C(i)$ , then  
 667  $C(i) \prec C(j)$  and  $C(i)$  is immediately followed by  $C(j)$ , so by  $(C_7(b))$  for  $\text{Conf}$ , we have a  
 668 contradiction. Thus  $C'(i) = C(i)C(i_1)..C(i_v)$  for  $v = \text{val}(i)$  and  $\text{con}(\text{gr}(i_v)) = \text{con}(G_t) \prec$   
 669  $\text{con}(\text{gr}(j))$  and  $C(i_v)$  is immediately followed by  $C(j)$ , a contradiction by  $(C_7(b))$  for  $\text{Conf}$ .
- 670 • Assume that  $\text{gr}'(j) = G_t$ .  
 671 Then the group  $\text{gr}(i)$  were partitioned when refining by  $G_t$ , and so  $C'(i) = \text{con}'(\text{gr}'(i)) =$   
 672  $\text{con}(\text{gr}(i))C(j)^v$  for  $v = \text{val}(j)$ . Since  $C'(i)$  is immediately followed by  $C'(j) = \text{con}(G_t)$ , we have  
 673 again a contradiction as it implies that  $\text{val}(j) = v+1$ .

674 That concludes the proof of Claim 2.

675 **Claim 3.** *The function  $prev'$  is proper on  $H_{t+1}$ , i.e.  $(C_6)$  for  $Conf'$  holds.*

676 *Proof of Claim 3.*

677 Let  $j = prev'(i)$  and  $i \in H_{t+1}$  with  $val'(i) = v$ . Then  $|[i]_{\approx}| = v$  and so  $[i]_{\approx} = \{i_1, \dots, i_v\}$ ,  
678 where  $i_1 < i_2 < \dots < i_v$ . Hence  $i_1, \dots, i_v \in H_{t+1}$  and  $C'(i_1) = \dots = C'(i_v) = con'(H_{t+1})$  and  
679  $j = prev'(i) = prev'(i_1) = \dots = prev'(i_v)$  and so  $j < i_1$ . It remains to show that  $C'(j)C'(i_1)\dots C'(i_v)$  is a  
680 prefix of  $x[j..n]$ . It suffices to show that  $C'(j)$  is immediately followed by  $C'(i_1)$ .

681  
682 If  $C'(j) \cap C'(i_1) \neq \emptyset$ , then by the Monge property  $(C_8)$ ,  $C'(i_1) \subseteq C'(j)$  as  $j < i_1$ , and so by  $(C_7(a))$ ,  
683  $con'(H_{t+1}) \prec con'(gr'(i_1)) = con'(H_{t+1})$ , a contradiction.

684  
685 Thus,  $C'(j) \cap C'(i_1) = \emptyset$ . Set  $j_1 = j + |con'(gr'(j))|$ . It follows that  $j_1 \leq i_1$ . Assume that  $j_1 < i_1$ .  
686 Since  $j = prev'(i_1)$  and  $j < i_1$ ,  $con'(gr'(j_1)) \succeq con'(gr'(i_1)) = con'(H_{t+1})$ . Since  $j_1 \notin H_{t+1}$ ,  
687  $con'(gr'(j_1)) \succ con'(H_{t+1})$ . Consider  $C'(j_1)$ . If  $C'(j_1) \cap C'(i_1) \neq \emptyset$ , then by  $(C_8)$ ,  $C'(i_1) \subseteq C'(j_1)$ , and  
688 so by  $(C_7(a))$ ,  $con'(H_{t+1}) \prec con'(gr'(i_1)) = con'(H_{t+1})$ , a contradiction. Thus  $C'(j_1) \cap C'(i_1) = \emptyset$ .  
689 Since  $C'(j_1)$  immediately follows  $C'(j)$ , by  $(C_7(b))$ ,  $con'(gr'(j_1)) \preceq con'(H_{t+1})$ , a contradiction.  
690 Therefore  $j_1 = i_1$ , and so  $prev'$  is proper on  $H_{t+1}$ .

691

692 That concludes the proof of Claim 3.

693 **Claim 4.**  *$H_{t+1}$  is a complete group, i.e.  $(C_5)$  for  $Conf'$  holds.*

694 *Proof of Claim 4.*

695 Assume that there is  $i \in H_{t+1}$  so that  $C'(i)$  is not maximal, i.e. for some  $k \geq i + |con'(H_{t+1})|$ ,  $x[i..k]$  is a  
696 maximal Lyndon substring of  $x$ .

697

698 Either  $k = n$  and so  $con'(gr'(k)) = x[k]$  and so  $C'(k)$  is a suffix of  $x[i..k]$ , or  $k < n$  and then  
699  $x[k+1] \prec x[k]$  since  $x[k+1] \preceq x[k]$  implies that  $x[i..k+1]$  is Lyndon, a contradiction with the maximality  
700 of  $x[i..k]$ . Consider  $C'(k)$ , then  $C'(k) \subseteq x[i..k]$  and so  $C'(k) = x[k]$ .

701

702 Therefore, there is  $j_1$  so that  $i + |con'(H_{t+1})| \leq j_1 \leq k$  and  $C'(j_1)$  is a suffix of  $x[i..k]$ . Take the smallest  
703  $j_1$  such. If  $j_1 = i + |con'(H_{t+1})|$ , then  $C'(i) \prec C'(j_1)$  as  $x[i..k] = C'(i)C'(j_1)$  is Lyndon. By  $(C_7(b))$ ,  
704  $C'(j_1) \preceq con'(H_{t+1})$ , so we have  $con'(H_{t+1}) = C'(i) \prec C'(j_1) \preceq con'(H_{t+1})$ , a contradiction.

705

706 Therefore,  $j_1 > i + |con'(H_{t+1})|$ . Consider  $x[j_1 - 1]$ . If  $x[j_1 - 1] \preceq x[j_1]$ ,  $x[j_1 - 1..k]$  is Lyndon, and since  
707  $x[j_1..k] = C'(j_1)$ ,  $x[j_1 - 1..k]$  would be a context of  $gr'(j_1 - 1)$ , and this contradicts the fact the  $j_1$  was  
708 chosen to be the smallest such. Therefore,  $x[j_1 - 1] \succ x[j_1]$  and so  $con'(gr'(j_1 - 1)) = x[j_1 - 1]$ . Thus,  
709 there is  $j_2$ ,  $i + |con'(H_{t+1})| \leq j_2 < j_1 \leq k$  and  $C'(j_2)$  is a suffix of  $x[i..j_1 - 1]$ . Take the smallest such  
710  $j_2$ . If  $C'(j_2) \prec C'(j_1)$ , then by  $(C_7(b))$ ,  $C'(j_1) \preceq con'(H_{t+1})$ , a contradiction. Hence  $C'(j_2) \succeq C'(j_1)$ .  
711 If  $j_2 = i + |con'(H_{t+1})|$ , then  $x[i..k] = C'(i)C'(j_2)C'(j_1)$  and so by  $(C_7(b))$ ,  $C'(j_2) \preceq con'(H_{t+1})$ , a  
712 contradiction. Hence  $i + |con'(H_{t+1})| < j_2$ .

713

714 The same argument done for  $j_2$  can be now done for  $j_3$ . We end up with  $i + |con'(H_{t+1})| \leq j_3 < j_2 < j_1 \leq k$   
715 and with  $C'(j_3) \succeq C'(j_2) \succeq C'(j_1) \succ con'(H_{t+1})$ . If  $i + |con'(H_{t+1})| = j_3$ , then we have a contradiction,  
716 so  $i + |con'(H_{t+1})| < j_3$ . These argument can be repeated only finitely many times, and we obtain  
717  $i + |con'(H_{t+1})| = j_\ell < j_{\ell-1} < \dots < j_2 < j_1 \leq k$  so that  $x[i..k] = C'(i)C'(j_\ell)C'(j_{\ell-1})\dots C'(j_2)C'(j_1)$ , which is  
718 a contradiction.

719

720 Therefore, our initial assumption that  $C'(i)$  is not maximal always leads to a contradiction.

721

722 That concludes the proof of Claim 4.

723 The four claims show that all the conditions  $(C_1) \dots (C_8)$  are satisfied for  $Conf'$ , and that proves  
724 Theorem 2.

725

□

726 As the last step, we show that when the process of refinement is completed, all maximal Lyndon  
727 substrings of  $x$  are identified and sorted via the contexts of the groups of the final configuration.

### 728 Theorem 3.

729 Let  $Conf_1 = \langle G_{k_1}^1, G_{k_1-1}^1, \dots, G_2^1, G_1^1 \rangle$  with  $gr_1()$ ,  $con_1()$ ,  $C_1()$ ,  $prev_1()$ , and  $val_1()$  be the initial 1-complete  
730 group configuration from Lemma 9.

731 Let  $Conf_2 = \langle G_{k_2}^2, G_{k_2-1}^2, \dots, G_2^2, G_1^2 \rangle$  with  $gr_2()$ ,  $con_2()$ ,  $C_2()$ ,  $prev_2()$ , and  $val_2()$  be the 2-complete group  
732 configuration obtained from  $Conf_1$  through the refinement by the group  $G_1^1$ .

733 Let  $Conf_3 = \langle G_{k_3}^3, G_{k_3-1}^3, \dots, G_2^3, G_1^3 \rangle$  with  $gr_3()$ ,  $con_3()$ ,  $C_3()$ ,  $prev_3()$ , and  $val_3()$  be the 3-complete group  
734 configuration obtained from  $Conf_2$  through the refinement by the group  $G_2^2$ .

735 ...

736 Let  $Conf_r = \langle G_{k_r}^r, G_{k_r-1}^r, \dots, G_2^r, G_1^r \rangle$  with  $gr_r()$ ,  $con_r()$ ,  $C_r()$ ,  $prev_r()$ , and  $val_r()$  be the  $r$ -complete group  
737 configuration obtained from  $Conf_{r-1}$  through the refinement by the group  $G_{r-1}^{r-1}$ . Let  $Conf_r$  be the final  
738 configuration after the refinement runs out.

739 Then  $x[i..k]$  is a maximal Lyndon substring of  $x$  iff  $x[i..k] = C_r(i) = con_r(gr_r(i))$ .

740 **Proof.** That all the groups of  $Conf_r$  are complete follows from Theorem 2 and hence every  $C_r(i)$  is a  
741 maximal Lyndon string. Let  $x[i..k]$  be a maximal Lyndon substring of  $x$ . Consider  $C_r(i)$ , since it is  
742 maximal, it must be equal to  $x[i..k]$ .

743

□

#### 744 4.3. Motivation for the Refinement

745 The process of refinement is in fact a process of gradual revealing of the Lyndon substrings which  
746 we call the *water draining method*.

- 747 (a) lower the water level by one
- 748 (b) extend the existing Lyndon substrings  
749 *the revealed letters are used to extend the existing Lyndon substrings where possible, or became Lyndon*  
750 *substrings of length 1 otherwise;*
- 751 (c) consolidate the new Lyndon substrings  
752 *processed from the right, if several Lyndon substrings are adjacent and can be joined to a longer Lyndon*  
753 *substring, they are joined.*

754

755 The diagram in Figure 5 and the description that follows it illustrate the method for a string  
756 011023122. The input string is visualized as a curve and the height at each point is the value of the  
757 letter at that position.



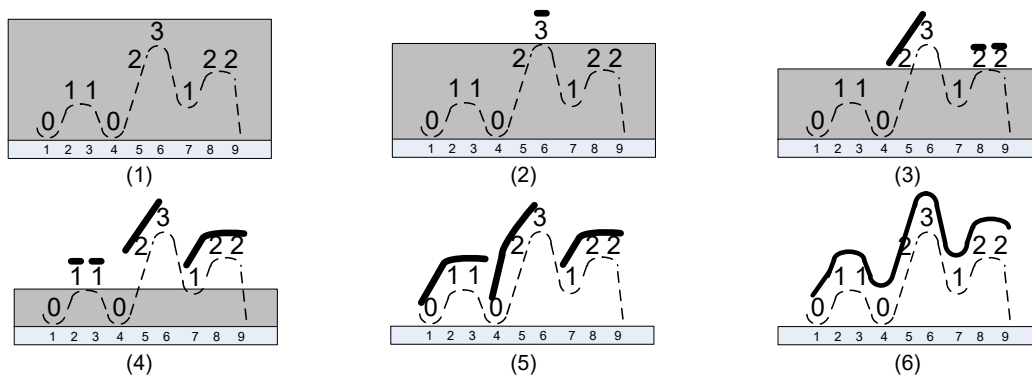


Figure 5. The water draining method for 011023122

In Figure 5, we illustrate the process:

- 758 (1) We start with the string 011023122 and a full tank of water.
- 759 (2) We drain one level, only 3 is revealed, nothing to extend, nothing to consolidate.
- 760 (3) We drain one more level and three 2's are revealed, the first 2 extends 3 to 23 and the remaining
- 761 two 2's form Lyndon substrings 2 of length 1, nothing to consolidate.
- 762 (4) We drain one more level and three 1's are revealed, the first two 1's form Lyndon substrings 1
- 763 of length 1, the third 1 extends 22 to 122, nothing to consolidate.
- 764 (5) We drain one more level and two 0's are revealed, the first 0 extends 11 to 011, the second 0
- 765 extends 23 to 023. In the consolidation phase, 023 is joined with 122 to form a Lyndon substring
- 766 023122, and then 011 is joined with 023122 to form a Lyndon substring 011023122.
- 767

768 So, during the process, the following maximal Lyndon substrings were identified: 3 at position 6,

769 23 at position 5, 2 at positions 8 and 9, 1 at positions 2 and 3, 122 at position 7, 023 at position 4, and

770 finally 011023122 at position 1. Note that all positions are accounted for, we really have all maximal

771 Lyndon substrings of the string 011023122.

772

0	1	1	0	2	3	1	2	2
1	2	3	4	5	6	7	8	9

$$\begin{aligned}
 &G_0 = \{1,4\} \quad G_1 = \{2,3,7\} \quad G_2 = \{5,8,9\} \quad G_3 = \{6\} \\
 &G_0 = \{1,4\} \quad G_1 = \{2,3,7\} \quad G_2 = \{8,9\} \quad G_{23} = \{5\} \quad G_3 = \{6\} \\
 &G_0 = \{1\} \quad G_{023} = \{4\} \quad G_1 = \{2,3,7\} \quad G_2 = \{8,9\} \quad G_{23} = \{5\} \quad G_3 = \{6\} \\
 &G_0 = \{1\} \quad G_{023} = \{4\} \quad G_1 = \{2,3\} \quad G_{122} = \{7\} \quad G_2 = \{8,9\} \quad G_{23} = \{5\} \quad G_3 = \{6\} \\
 &G_0 = \{1\} \quad G_{023122} = \{4\} \quad G_1 = \{2,3\} \quad G_{122} = \{7\} \quad G_2 = \{8,9\} \quad G_{23} = \{5\} \quad G_3 = \{6\} \\
 &G_{011} = \{1\} \quad G_{023122} = \{4\} \quad G_1 = \{2,3\} \quad G_{122} = \{7\} \quad G_2 = \{8,9\} \quad G_{23} = \{5\} \quad G_3 = \{6\} \\
 &G_{011023122} = \{1\} \quad G_{023122} = \{4\} \quad G_1 = \{2,3\} \quad G_{122} = \{7\} \quad G_2 = \{8,9\} \quad G_{23} = \{5\} \quad G_3 = \{6\}
 \end{aligned}$$

Figure 6. Group Refinement for 011023122

773 In Figure 6 we present an illustrative example for the string 011023122, where the arrows  
 774 represent the *prev* mapping shown only on the group used for the refinement. The groups are  
 775 indicated by the bold font.

776

#### 777 4.4. The Complexity of BSLA

778 The computation of the initial configuration can be done in linear time as the input string is tight.  
 779 The tightness of the input string is used to guarantee the linearity of computing the initial value of  
 780 *prev*. Since all groups are non-empty, there can never be more groups than  $n$ . Theorem 2 is at the  
 781 heart of the algorithm. The refinement by the last completed group is linear in the size of the group,  
 782 including the update of *prev*. Therefore, the overall all worst case complexity of *BSLA* is linear in the  
 783 length of the input string.

## 784 5. Measurements

785 Initially, computations were performed on the Department of Computing and Software's *moore*  
 786 server; Memory: 32GB (DDR4 @ 2400 MHz), CPU: 8 x Intel Xeon E5-2687W v4 @ 3.00GHz, OS:  
 787 Linux version 2.6.18-419.el5 (gcc version 4.1.2 and Red Hat version 4.1.2-55). To verify correctness,  
 788 new randomized data was produced and computed independently on the University of Toronto  
 789 Mississauga's *octolab* cluster; Memory: 8 x 32GB (DDR4 @ 3200 MHz), CPU: 8 x AMD Ryzen  
 790 Threadripper 1920X (12-Core) @ 4.00GHz, OS: Ubuntu 16.04.6 LTS (gcc version 5.4.0). The results of  
 791 both were extremely similar, those reported herein are those generated using the *moore* server. All the  
 792 programs were compiled without any additional level of optimization<sup>3</sup>. The CPU time was measured  
 793 in clock ticks with 1,000,000 clock ticks per second. Since the execution time was negligible for short  
 794 strings, the processing of the same string was repeated several times (the repeat factor varied from  $10^6$ ,  
 795 for strings of length 10, to 1, for strings of length  $5 \times 10^6$ ), resulting in a higher precision. Thus, for  
 796 graphing, the logarithmic scale was used for both, the  $x$ -axis representing the length of the strings, and  
 797 the  $y$ -axis representing the time. We used 4 categories of randomly generated data sets:

- 798 (1) `bin`  
 799 random tight binary strings over the alphabet  $\{0, 1\}$
- 800 (2) `dna`  
 801 random tight 4-ary (kind of random DNA) over the alphabet  $\{0, 1, 2, 3\}$
- 802 (3) `eng`  
 803 random tight 26-ary strings (kind of random English) over the alphabet  $\{0, 1, \dots, 25\}$
- 804 (4) `int`  
 805 random tight strings over integer alphabet  $\{0, 1, \dots, n\}$ .

806 Each data set contains 100 randomly generated strings of the same length. For each category, there  
 807 were data sets for length 10, 50,  $10^2$ ,  $5 \times 10^2$ , ...,  $10^5$ ,  $5 \times 10^5$ ,  $10^6$ , and  $5 \times 10^6$ . The minimum, average,  
 808 and maximum time for each data set was computed. Since the variance for each dataset was minimal,  
 809 the results for minimum times and the results for maximum times completely mimicked the results for  
 810 the average times, so we only present here the averages.

811 Tables 1, 2, 3, 4 and the graphs in Figures 7, 8, 9, 10 from Section 7 clearly indicate that the  
 812 performance of the three algorithms is linear and virtually indistinguishable. We expected *IDLA* and  
 813 *TRLA* to exhibit linear behaviour on random strings as such strings tend to have. However, we did not  
 814 expect the results to be so close.

---

<sup>3</sup> i.e. neither `-O1`, nor `-O2`, nor `-O3` flag were specified for the compilation

815 The data set `extreme_idla` contains individual strings  $0123\dots n$  of the required lengths designed  
 816 to force *IDLA* into its worst quadratic performance. Table 5 and the graph in Figure 11 from Section 7  
 817 show this clearly.

818 In Subsection 3.5, we described how the three data sets `extreme_trla`, `extreme_trla1`, and  
 819 `extreme_trla2` are generated and why. The results of experimenting with these data sets do not  
 820 suggest that the worst-case complexity for *TRLA* is  $\mathcal{O}(n \log(n))$ . Yet again, the performances of  
 821 the three algorithms are linear and virtually indistinguishable, see Tables 6, 7 and the graphs in  
 822 Figures 12, 13 from Section 7.

823

## 824 6. Conclusion and Future Work

825 We presented two novel algorithms for computing maximal Lyndon substrings. The first one,  
 826 *TRLA*, has a simple implementation with a complicated theory behind it. Its average time complexity  
 827 is linear in the length of the input string, and its worst-case complexity is no worse than  $\mathcal{O}(n \log(n))$ .  
 828 The  $\tau$ -reduction used in the algorithm is an interesting reduction preserving maximal Lyndon  
 829 substrings, a fact used significantly in the design of the algorithm. Interestingly, it seem to slightly  
 830 outperform *BSLA*, at least on the data sets used for our experimentations. *BSLA*, the second algorithm,  
 831 is linear and elementary in the sense that it does not require a pre-processed global data structure.  
 832 Being linear and elementary, *BSLA* is more interesting and it is possible that its performance could be  
 833 more streamlined. However, both the theory and implementation of *BSLA* are rather complex.

834

835 On random strings, none of the two algorithms were significantly better than the simple *IDLA*,  
 836 whose implementation is just a few lines. However, its quadratic worst-case complexity is an obstacle  
 837 as our experiments indicated.

838

839 Additional effort needs to go into proving *TRLA*'s worst-case complexity. The experiments  
 840 performed have not indicated that it is not linear even in the worst case. Both algorithms need to be  
 841 compared to some efficient implementation of *SSLA* and *BWLA*.

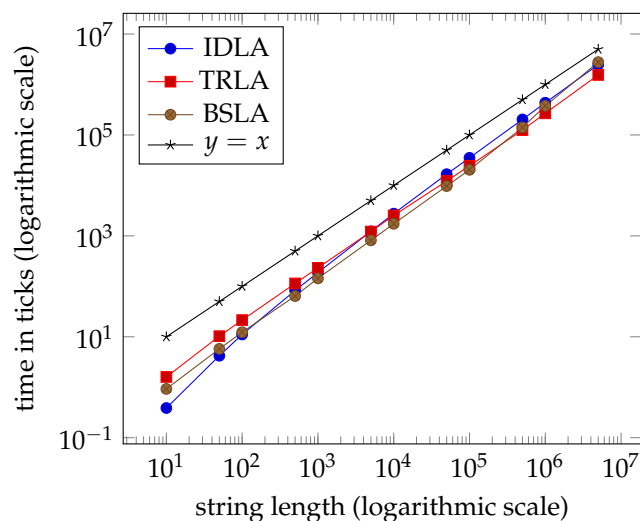
842

## 843 7. Results

844 This section contains the measurements of average times for the data sets discussed in the previous  
 845 section. For better understanding of the data, we present them in both, tabular and graphic form.

string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$3.876 \times 10^{-1}$	1.600	$9.269 \times 10^{-1}$
50	4.213	$1.028 \times 10$	5.762
100	$1.111 \times 10$	$2.143 \times 10$	$1.233 \times 10$
500	$8.270 \times 10$	$1.141 \times 10^2$	$6.414 \times 10$
1000	$1.872 \times 10^2$	$2.324 \times 10^2$	$1.438 \times 10^2$
5000	$1.233 \times 10^3$	$1.216 \times 10^3$	$8.155 \times 10^2$
10000	$2.758 \times 10^3$	$2.537 \times 10^3$	$1.735 \times 10^3$
50000	$1.669 \times 10^4$	$1.222 \times 10^4$	$9.702 \times 10^3$
100000	$3.523 \times 10^4$	$2.447 \times 10^4$	$2.039 \times 10^4$
500000	$2.031 \times 10^5$	$1.249 \times 10^5$	$1.409 \times 10^5$
1000000	$4.345 \times 10^5$	$2.687 \times 10^5$	$3.764 \times 10^5$
5000000	$2.479 \times 10^6$	$1.547 \times 10^6$	$2.775 \times 10^6$

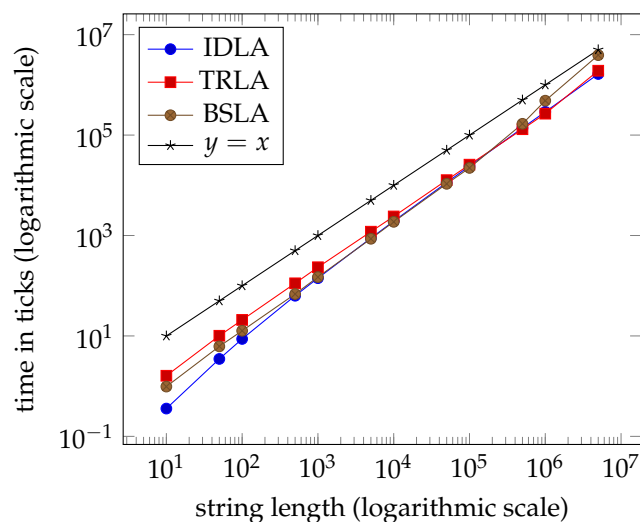
846 **Table 1.** Average times for data set bin  
 (10<sup>6</sup> clock ticks per second)



**Figure 7.** Average times for data set bin  
( $10^6$  clock ticks per second)  
for reference the  $y = x$  curve is shown

string length	time in ticks	time in ticks	time in ticks
	IDLA	TRLA	BSLA
10	$3.576 \times 10^{-1}$	1.611	$9.803 \times 10^{-1}$
50	3.484	$1.008 \times 10$	6.196
100	8.711	$2.084 \times 10$	$1.268 \times 10$
500	$6.313 \times 10$	$1.117 \times 10^2$	$6.833 \times 10$
1000	$1.413 \times 10^2$	$2.332 \times 10^2$	$1.488 \times 10^2$
5000	$8.861 \times 10^2$	$1.191 \times 10^3$	$8.605 \times 10^2$
10000	$1.931 \times 10^3$	$2.396 \times 10^3$	$1.872 \times 10^3$
50000	$1.133 \times 10^4$	$1.266 \times 10^4$	$1.070 \times 10^4$
100000	$2.418 \times 10^4$	$2.575 \times 10^4$	$2.218 \times 10^4$
500000	$1.385 \times 10^5$	$1.309 \times 10^5$	$1.677 \times 10^5$
1000000	$2.904 \times 10^5$	$2.670 \times 10^5$	$4.840 \times 10^5$
5000000	$1.645 \times 10^6$	$1.913 \times 10^6$	$3.907 \times 10^6$

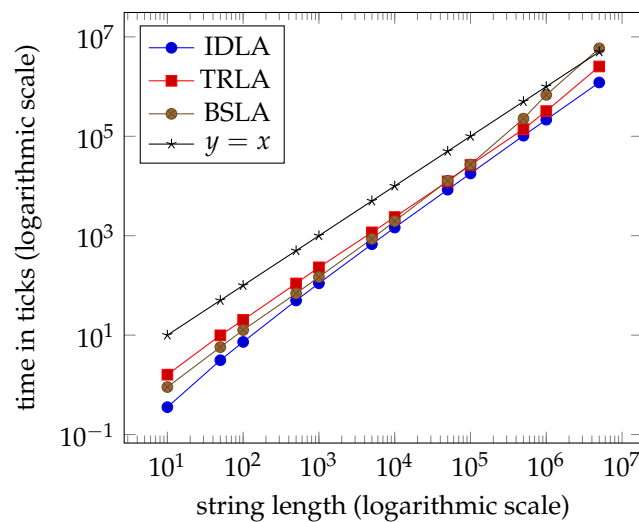
**Table 2.** Average times for data set dna  
( $10^6$  clock ticks per second)



**Figure 8.** Average times for data set dna  
( $10^6$  clock ticks per second)  
for reference the  $y = x$  curve is shown

string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$3.556 \times 10^{-1}$	1.602	$9.021 \times 10^{-1}$
50	3.128	9.909	5.729
100	7.293	$2.022 \times 10$	$1.266 \times 10$
500	$4.976 \times 10$	$1.095 \times 10^2$	$6.917 \times 10$
1000	$1.104 \times 10^2$	$2.313 \times 10^2$	$1.478 \times 10^2$
5000	$6.745 \times 10^2$	$1.167 \times 10^3$	$8.647 \times 10^2$
10000	$1.459 \times 10^3$	$2.371 \times 10^3$	$1.968 \times 10^3$
50000	$8.433 \times 10^3$	$1.240 \times 10^4$	$1.274 \times 10^4$
100000	$1.786 \times 10^4$	$2.649 \times 10^4$	$2.703 \times 10^4$
500000	$1.028 \times 10^5$	$1.392 \times 10^5$	$2.267 \times 10^5$
1000000	$2.157 \times 10^5$	$3.210 \times 10^5$	$6.819 \times 10^5$
5000000	$1.208 \times 10^6$	$2.538 \times 10^6$	$5.873 \times 10^6$

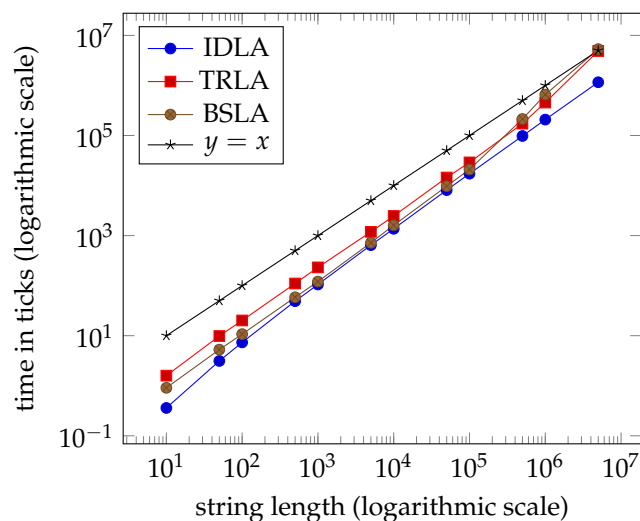
**Table 3.** Average times for data set `eng`  
( $10^6$  clock ticks per second)



**Figure 9.** Average times for data set `eng`  
( $10^6$  clock ticks per second)  
for reference the  $y = x$  curve is shown

string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$3.601 \times 10^{-1}$	1.591	$9.091 \times 10^{-1}$
50	3.128	9.843	5.269
100	7.341	$2.014 \times 10$	$1.074 \times 10$
500	$4.893 \times 10$	$1.103 \times 10^2$	$5.863 \times 10$
1000	$1.055 \times 10^2$	$2.315 \times 10^2$	$1.206 \times 10^2$
5000	$6.450 \times 10^2$	$1.188 \times 10^3$	$7.191 \times 10^2$
10000	$1.368 \times 10^3$	$2.463 \times 10^3$	$1.609 \times 10^3$
50000	$8.106 \times 10^3$	$1.446 \times 10^4$	$9.722 \times 10^3$
100000	$1.727 \times 10^4$	$2.899 \times 10^4$	$2.097 \times 10^4$
500000	$9.774 \times 10^4$	$1.725 \times 10^5$	$2.136 \times 10^5$
1000000	$2.078 \times 10^5$	$4.550 \times 10^5$	$6.558 \times 10^5$
5000000	$1.159 \times 10^6$	$4.876 \times 10^6$	$5.316 \times 10^6$

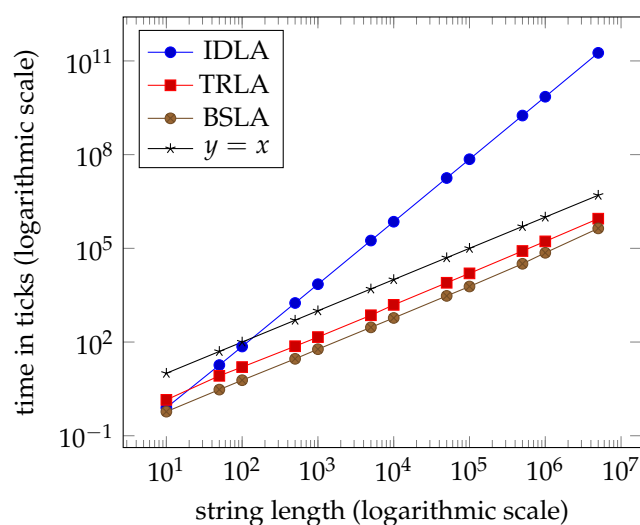
**Table 4.** Average times for data set `int`  
( $10^6$  clock ticks per second)



**Figure 10.** Average times for data set `int`  
( $10^6$  clock ticks per second)  
for reference the  $y = x$  curve is shown

string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$7.900 \times 10^{-1}$	1.410	$5.900 \times 10^{-1}$
50	$1.830 \times 10$	8.300	3.000
100	$7.190 \times 10$	$1.580 \times 10$	6.000
500	$1.778 \times 10^3$	$7.350 \times 10$	$2.900 \times 10$
1000	$7.105 \times 10^3$	$1.440 \times 10^2$	$5.900 \times 10$
5000	$1.776 \times 10^5$	$7.200 \times 10^2$	$2.950 \times 10^2$
10000	$7.111 \times 10^5$	$1.550 \times 10^3$	$5.900 \times 10^2$
50000	$1.784 \times 10^7$	$7.900 \times 10^3$	$3.000 \times 10^3$
100000	$7.130 \times 10^7$	$1.580 \times 10^4$	$6.000 \times 10^3$
500000	$1.783 \times 10^9$	$8.200 \times 10^4$	$3.200 \times 10^4$
1000000	$7.137 \times 10^9$	$1.670 \times 10^5$	$7.200 \times 10^4$
5000000	$1.813 \times 10^{11}$	$8.900 \times 10^5$	$4.350 \times 10^5$

**Table 5.** Average times for data set `extreme_idla`  
( $10^6$  clock ticks per second)

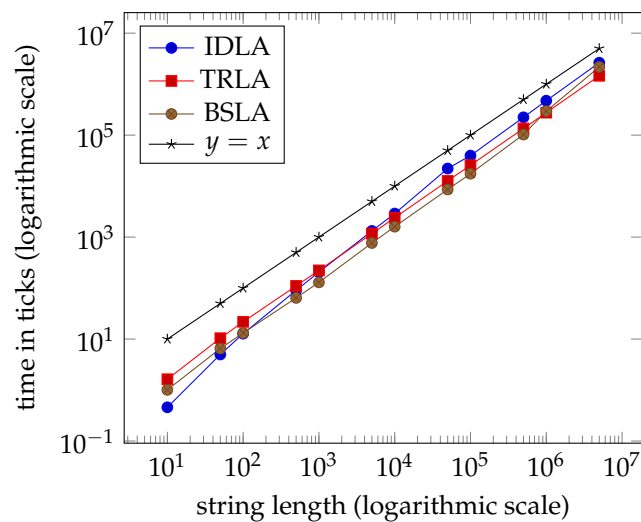


**Figure 11.** Average times for data set `extreme_idla`  
( $10^6$  clock ticks per second)  
for reference the  $y = x$  curve is shown



string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$4.588 \times 10^{-1}$	1.628	1.010
50	4.987	$1.039 \times 10$	6.527
100	$1.275 \times 10$	$2.179 \times 10$	$1.315 \times 10$
500	$9.033 \times 10$	$1.101 \times 10^2$	$6.418 \times 10$
1000	$2.060 \times 10^2$	$2.222 \times 10^2$	$1.295 \times 10^2$
5000	$1.319 \times 10^3$	$1.171 \times 10^3$	$7.662 \times 10^2$
10000	$2.896 \times 10^3$	$2.394 \times 10^3$	$1.600 \times 10^3$
50000	$2.209 \times 10^4$	$1.263 \times 10^4$	$8.598 \times 10^3$
100000	$3.965 \times 10^4$	$2.567 \times 10^4$	$1.752 \times 10^4$
500000	$2.233 \times 10^5$	$1.349 \times 10^5$	$1.027 \times 10^5$
1000000	$4.734 \times 10^5$	$2.759 \times 10^5$	$2.950 \times 10^5$
5000000	$2.632 \times 10^6$	$1.458 \times 10^6$	$2.174 \times 10^6$

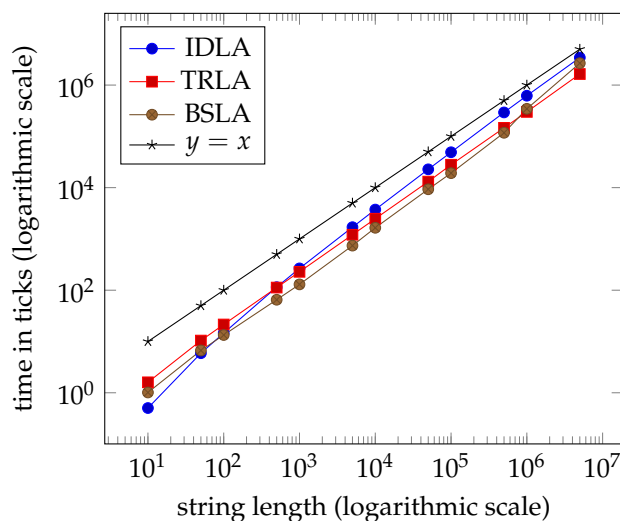
**Table 6.** Average times for data set `extreme_trla`  
( $10^6$  clock ticks per second)



**Figure 12.** Average times for data set `extreme_trla`  
( $10^6$  clock ticks per second)  
for reference the  $y = x$  curve is shown

string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$5.040 \times 10^{-1}$	1.600	1.013
50	5.910	$1.042 \times 10$	6.689
100	$1.460 \times 10$	$2.145 \times 10$	$1.338 \times 10$
500	$1.146 \times 10^2$	$1.126 \times 10^2$	$6.514 \times 10$
1000	$2.662 \times 10^2$	$2.284 \times 10^2$	$1.298 \times 10^2$
5000	$1.694 \times 10^3$	$1.205 \times 10^3$	$7.408 \times 10^2$
10000	$3.734 \times 10^3$	$2.477 \times 10^3$	$1.648 \times 10^3$
50000	$2.276 \times 10^4$	$1.310 \times 10^4$	$9.294 \times 10^3$
100000	$4.901 \times 10^4$	$2.796 \times 10^4$	$1.917 \times 10^4$
500000	$2.928 \times 10^5$	$1.465 \times 10^5$	$1.174 \times 10^5$
1000000	$6.199 \times 10^5$	$3.000 \times 10^5$	$3.443 \times 10^5$
5000000	$3.432 \times 10^6$	$1.642 \times 10^6$	$2.671 \times 10^6$

**Table 7.** Average times for data set `extreme_trla1`  
( $10^6$  clock ticks per second)



**Figure 13.** Average times for data set `extreme_tr1a1`

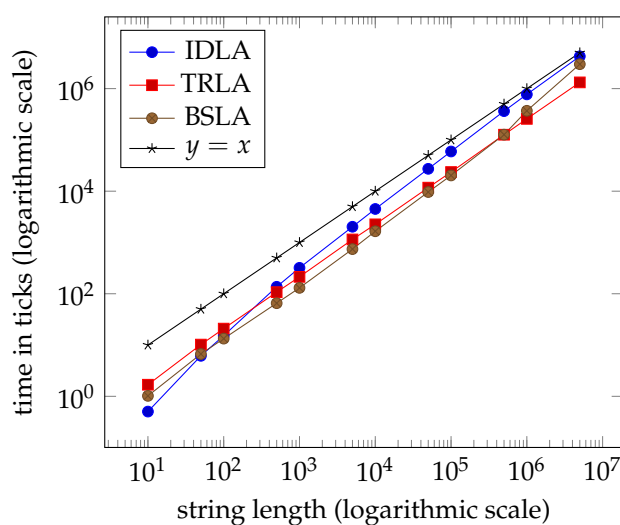
( $10^6$  clock ticks per second)

for reference the  $y = x$  curve is shown

string length	time in ticks IDLA	time in ticks TRLA	time in ticks BSLA
10	$5.041 \times 10^{-1}$	1.683	1.020
50	6.160	$1.020 \times 10$	6.684
100	$1.526 \times 10$	$2.090 \times 10$	$1.334 \times 10$
500	$1.367 \times 10^2$	$1.074 \times 10^2$	$6.543 \times 10$
1000	$3.202 \times 10^2$	$2.135 \times 10^2$	$1.302 \times 10^2$
5000	$2.024 \times 10^3$	$1.145 \times 10^3$	$7.421 \times 10^2$
10000	$4.500 \times 10^3$	$2.257 \times 10^3$	$1.654 \times 10^3$
50000	$2.728 \times 10^4$	$1.172 \times 10^4$	$9.623 \times 10^3$
100000	$5.941 \times 10^4$	$2.362 \times 10^4$	$2.027 \times 10^4$
500000	$3.639 \times 10^5$	$1.262 \times 10^5$	$1.266 \times 10^5$
1000000	$7.719 \times 10^5$	$2.571 \times 10^5$	$3.695 \times 10^5$
5000000	$4.263 \times 10^6$	$1.323 \times 10^6$	$2.994 \times 10^6$

**Table 8.** Average times for data set `extreme_tr1a2`

( $10^6$  clock ticks per second)



**Figure 14.** Average times for data set `extreme_tr1a2`

( $10^6$  clock ticks per second)

for reference the  $y = x$  curve is shown

846 **Author Contributions:** Conceptualization, Frantisek Franek and Michael Liut; Data curation, Frantisek Franek  
847 and Michael Liut; Formal analysis, Frantisek Franek and Michael Liut; Funding acquisition, Frantisek Franek;  
848 Investigation, Frantisek Franek and Michael Liut; Methodology, Frantisek Franek and Michael Liut; Project  
849 administration, Frantisek Franek; Resources, Frantisek Franek and Michael Liut; Software, Frantisek Franek  
850 and Michael Liut; Supervision, Frantisek Franek; Validation, Frantisek Franek and Michael Liut; Visualization,  
851 Frantisek Franek and Michael Liut; Writing – original draft, Frantisek Franek and Michael Liut; Writing – review  
852 & editing, Frantisek Franek and Michael Liut.

853 **Funding:** This research was funded by the National Sciences and Research Council of Canada (NSERC) grant  
854 RGPIN/5504-2018.

855 **Conflicts of Interest:** The authors declare no conflict of interest.

## 856 Abbreviations

857 The following abbreviations are used in this manuscript:

858 BSLA Baier’s Sort Lyndon Array  
859 IDLA Iterative Duval Lyndon Array  
TRLA Tau Reduction Lyndon Array

## 860 References

- 861 1. Bannai, H.; I, T.; Inenaga, S.; Nakashima, Y.; Takeda, M.; Tsuruta, K. The “Runs” Theorem. Available at  
862 <https://arxiv.org/abs/1406.0263>.
- 863 2. Bannai, H.; I, T.; Inenaga, S.; Nakashima, Y.; Takeda, M.; Tsuruta, K. The “Runs” Theorem. *SIAM J.*  
864 *COMPUT.* **2017**, *46*, 1501–1514.
- 865 3. Franek, F.; Paracha, A.; Smyth, W. The linear equivalence of the suffix array and the partially sorted  
866 Lyndon array. Proc. Prague Stringology Conference, 2017, pp. 77–84.
- 867 4. Baier, U. Linear-time suffix sorting — a new approach for suffix array construction. M.Sc. Thesis,  
868 University of Ulm, Ulm, Germany.
- 869 5. Baier, U. Linear-time suffix sorting — a new approach for suffix array construction. 27th Annual  
870 Symposium on Combinatorial Pattern Matching (CPM 2016); Grossi, R.; Lewenstein, M., Eds.; Schloss  
871 Dagstuhl–Leibniz-Zentrum fuer Informatik: Dagstuhl, Germany, 2016; Vol. 54, *Leibniz International*  
872 *Proceedings in Informatics (LIPIcs)*, pp. 1–12.
- 873 6. Chen, G.; Puglisi, S.; Smyth, W. Lempel-Ziv factorization using less time & space. *Mathematics in Computer*  
874 *Science* **2013**, *1*, 605–623.
- 875 7. Crochemore, M.; Ilie, L.; Smyth, W. A simple algorithm for computing the Lempel-Ziv factorization. Proc.  
876 18th Data Compression Conference, 2008, pp. 482–488.
- 877 8. Kosolobov, D. Lempel-Ziv factorization may be harder than computing all runs. Proceedings of 32  
878 International Symposium on Theoretical Aspects of Computer Science – STACS 2015, 2015, pp. 582–593.
- 879 9. Digelmann, C. Personal communication.
- 880 10. Franek, F.; Sohidull Islam, A.; Sohel Rahman, M.; Smyth, W. Algorithms to compute the Lyndon array.  
881 Proceedings of Prague Stringology Conference 2016, 2016, pp. 172–184.
- 882 11. Duval, J.P. Factorizing words over an ordered alphabet. *J. Algorithms* **1983**, *4*, 363–381.
- 883 12. Hohlweg, C.; Reutenauer, C. Lyndon words, permutations and trees. *Theoretical Computer Science* **2003**,  
884 *307*, 173–178. doi:[http://dx.doi.org/10.1016/S0304-3975\(03\)00099-9](http://dx.doi.org/10.1016/S0304-3975(03)00099-9).
- 885 13. Nong, G.; Zhang, S.; Chan, W.H. Linear suffix array construction by almost pure induced-sorting. 2009  
886 Data Compression Conference, 2009, pp. 193–202.
- 887 14. Louza, F.; Smyth, W.; Manzini, G.; Telles, G. Lyndon array construction during Burrows–Wheeler inversion.  
888 *Journal of Discrete Algorithms* **2018**, *50*, 2–9.
- 889 15. Franek, F.; Liut, M.; Smyth, W. On Baier’s sort of maximal Lyndon substrings. Proceedings of Prague  
890 Stringology Conference 2018, 2018, pp. 63–78.
- 891 16. C++ code for IDLA, TRLA and BSLA algorithms. Available at:  
892 <https://github.com/MichaelLiut/Computing-LyndonArray>.
- 893 17. Farach, M. Optimal suffix tree construction with large alphabets. Proc. 38th IEEE Symp. Foundations of  
894 Computer Science. IEEE, 1997, pp. 137–143.

- 895 18. Nong, G. Practical linear-time  $O(1)$ -workspace suffix sorting for constant alphabets. *ACM Trans. Inf. Syst.*  
896 **2013**, *31*, 1–15.
- 897 19. Cooley, J.; Tukey, J. An algorithm for the machine calculation of complex Fourier series. *Mathematics of*  
898 *Computation* **1965**, *19*, 297–301.
- 899 20. Franek, F.; Liut, M. Computing maximal Lyndon substrings of a string, AdvOL Report 2019/2, McMaster  
900 University. Available at:  
901 [http://optlab.mcmaster.ca//component/option,com\\_docman/task\\_cat\\_view/gid,77/Itemid,92](http://optlab.mcmaster.ca//component/option,com_docman/task_cat_view/gid,77/Itemid,92).
- 902 21. Franek, F.; Liut, M. Algorithms to compute the Lyndon array revisited. *Proc. of Prague Stringology*  
903 *Conference 2019*, 2019, pp. 16–28.
- 904 22. Liut, M. Computing Lyndon Arrays. Ph.D. Thesis, McMaster University, Hamilton, Ontario, Canada.
- 905 23. Paracha, A. Lyndon factors and periodicities in strings. Ph.D. Thesis, McMaster University, Hamilton,  
906 Ontario, Canada.
- 907 24. Kärkkäinen, J.; Sanders, P. Simple linear work suffix array construction. *Proceedings of the 30th*  
908 *international conference on Automata, languages and programming*; Springer-Verlag: Berlin, Heidelberg,  
909 2003; ICALP'03, pp. 943–955.