

1 Finding an Optimal Alphabet Ordering for Lyndon 2 Factorization is Hard

3 Daniel Gibney

4 Department of Computer Science, University of Central Florida, USA

5 <https://www.cs.ucf.edu/~dgibney/>

6 daniel.gibney@ucf.edu

7 Sharma V. Thankachan

8 Department of Computer Science, University of Central Florida, USA

9 <http://www.cs.ucf.edu/~sharma/>

10 sharma.thankachan@ucf.edu

11 — Abstract —

12 This work establishes several strong hardness results on the problem of finding an ordering on a
13 string's alphabet that either minimizes or maximizes the number of factors in that string's Lyndon
14 factorization. In doing so, we demonstrate that these ordering problems are sufficiently complex
15 to model a wide variety of ordering constraint satisfaction problems (OCSPs). Based on this, we
16 prove that (i) the decision versions of both the minimization and maximization problems are NP-
17 complete, (ii) for both the minimization and maximization problems there does not exist a constant
18 approximation algorithm running in polynomial time under the Unique Game Conjecture and (iii)
19 there does not exist an algorithm to solve the minimization problem in time $\text{poly}(|T|) \cdot 2^{o(\sigma \log \sigma)}$ for
20 a string T over an alphabet of size σ under the Exponential Time Hypothesis (essentially the brute
21 force approach of trying every alphabet order is hard to improve significantly).

22 **2012 ACM Subject Classification**

23 **Keywords and phrases** Lyndon Factorization, String Algorithms, Burrows-Wheeler Transform

24 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23



© John Q. Open and Joan R. Access;
licensed under Creative Commons License CC-BY
42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:15



Leibniz International Proceedings in Informatics
LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

25 **1** Introduction

26 A Lyndon word is a string that is lexicographically strictly smallest among all of its cyclic shifts.
 27 Letting \circ denote concatenation, the Lyndon factorization of a string T is the partitioning
 28 of T into Lyndon words T_1, T_2, \dots, T_f that are lexicographically non-increasing and $T =$
 29 $T_1 \circ T_2 \circ \dots \circ T_f$. For example, the Lyndon factorization of 0, 1, 0, 0, 2, 1, 1, 0, 0, 1, 0, 1, 1, 2 is
 30 $(0, 1)(0, 0, 2, 1, 1)(0, 0, 1, 0, 1, 1, 2)$, assuming the usual ordering, $0 < 1 < 2$.

31 Lyndon words and Lyndon factorization are well-studied, and play an important role
 32 in string algorithms [1, 2, 10, 24, 28, 30], algebra and combinatorics [7, 17, 25], and data
 33 compression [12, 18, 20, 34, 35]. As an example, it was shown in [29] that local suffixes inside
 34 each Lyndon factor can be sorted independently and then merged to construct a string's
 35 suffix array. As another example, Lyndon factorization is used in both the construction
 36 of a string's bijective Burrows-Wheeler transform (BBWT) [13] and in performing pattern
 37 matching on indexes built from the string's BBWT [3], where the number of steps used
 38 to locate occurrences of a pattern P depends on the number of Lyndon factors within a
 39 particular suffix of P . Because of such applications, it would be beneficial to be able to
 40 control the number of factors in the Lyndon factorization of a string. Unfortunately, the
 41 Lyndon factorization of a string is fixed and unique under a fixed ordering of its alphabet
 42 [26]. However, it can vary under different alphabet orderings. For instance, if we change the
 43 alphabet ordering to $2 < 0 < 1$ in our example above, we obtain the Lyndon factorization
 44 $(0, 1), (0), (0), (2, 1, 1, 0, 0, 1, 0, 1, 1), (2)$. This leads to the following problems:

45 ► **Problem 1** (Lyndon Factor Minimization - Decision Version). *Given an integer A and text*
 46 *T over alphabet Σ , does there exist an ordering on Σ such that the number of Lyndon factors*
 47 *of T is at most A ?*

48 ► **Problem 2** (Lyndon Factor Maximization - Decision Version). *Given an integer A and text*
 49 *T over alphabet Σ , does there exist an ordering on Σ such that the number of Lyndon factors*
 50 *of T is at least A ?*

51 We will also consider the *optimization variants* of these problems. The objective cost
 52 of a solution is the number of factors in its Lyndon factorization. In particular, for the
 53 minimization problem, a λ -approximation for $\lambda > 1$, is a polynomial-time algorithm that
 54 outputs an alphabet ordering where the number of factors is at most λ times the minimum
 55 possible number of factors over all possible alphabet orderings. Similarly, for the maximization
 56 problem, a λ -approximation for $\lambda < 1$, is a polynomial-time algorithm that outputs an
 57 alphabet ordering where the number of factors is at least λ times the maximum number of
 58 possible factors over all possible alphabet orderings.

59 These problems were first considered by Clare and Daykin, who proposed a polynomial-
 60 time greedy algorithm that can be adjusted to provide either a small number of factors or
 61 a large number of factors [8]. Through experiments, the authors showed that the number
 62 of factors can be significantly affected by their algorithm. Another approach that uses
 63 evolutionary algorithms to find alphabet orderings to optimize the number of Lyndon factors
 64 was considered in [9] and in [27]. Again, it was shown that there is often a significant effect on
 65 the number of factors, which can be controlled by the use of different fitness functions within
 66 the evolutionary algorithms. These techniques, although appearing to have a significant
 67 impact on the number of factors, do not provide any approximation guarantee. Motivated
 68 by this, and by other alphabet ordering hardness results presented in [4], we present the first
 69 set of hardness results on these problems.

70 ► **Theorem 1.** *The decision version of Lyndon Factor Minimization is NP-complete.*

71 ▶ **Theorem 2.** *Under the Exponential Time Hypothesis, the optimization version of Lyndon*
 72 *Factor Minimization cannot be solved in time $\text{poly}(|T|) \cdot 2^{o(\sigma \log \sigma)}$.*

73 ▶ **Theorem 3.** *Under the Unique Games Conjecture, the optimization version of Lyndon*
 74 *Factor Minimization does not admit a λ -approximation for any constant $\lambda > 1$.*

75 ▶ **Theorem 4.** *The decision version of Lyndon Factor Maximization is NP-complete.*

76 ▶ **Theorem 5.** *Under the Unique Games Conjecture, the optimization version of Lyndon*
 77 *Factor Maximization does not admit a λ -approximation for any constant $\lambda < 1$.*

78 We will prove these theorems in Section 3.1, Section 3.2, Section 3.3, Section 4.1, and Section
 79 4.2, respectively. We leave open whether it is possible to have a result similar to Theorem 2
 80 for Lyndon Factor Maximization.

81 Our main line of attack is to model ordering constraint satisfaction problems (OCSPs), a
 82 subject of extensive research in its own right [5, 6, 15, 16, 31, 33]. In these problems, the
 83 task is to find a linear ordering on a set of variables subject to some additional constraints.
 84 Our work shows that a solver for these Lyndon factorization problems would be powerful
 85 enough to solve difficult OCSP instances. Our results make use of strings that allow us to
 86 model different constraint satisfaction problems and thus prove our hardness results.

87 2 Preliminaries

88 We denote the concatenation of the strings u and v using the ‘ \circ ’ symbol, writing their
 89 concatenation as $u \circ v$. However, we omit ‘ \circ ’ where the concatenation is clear from context
 90 and it would be cumbersome to use. Throughout this paper, we will use ‘ $<$ ’ and ‘ $>$ ’ to refer
 91 to alphabet order between symbols, the lexicographic order between strings, and the usual
 92 ordering between real numbers. Again, context will make it clear which type of order is
 93 meant. A suffix of a string T is a string v such that $T = u \circ v$ for some string u . The suffix
 94 array $SA[\cdot]$ of a string $T[1, n]$ is a length n array where $SA[i]$ is equal to the starting index
 95 of the i^{th} lexicographically smallest suffix of T . The inverse suffix array $ISA[\cdot]$ is defined as
 96 the length n array such that $i = ISA[SA[i]]$, i.e., the position in the lexicographic order of
 97 the suffix starting at index i .

98 The Lyndon factorization (defined in Section 1) of a string can be computed in linear
 99 time. This can be done using the well known Duval’s algorithm [11], or by using the inverse
 100 suffix array, ISA , which can be constructed in linear time [22]. Lemma 6 makes it clear why
 101 the latter technique works.

102 ▶ **Lemma 6** (Theorem 2.2 [29]). *The starting index, i , of a suffix in T that is lexicographically*
 103 *smaller than any suffix starting at index $j < i$ is an index where a Lyndon factor begins.*

104 We will use this observation to construct strings that model constraints that occur in an
 105 ordering constraint satisfaction problem (OCSP). The definition of an OCSP used here is
 106 less general than the one given in [14], but still sufficient for our purposes.

107 ▶ **Definition 7.** *An OCSP of arity k is specified by a set $\Lambda \subseteq S_k$ where S_k is the set of*
 108 *permutations of $\{1, 2, \dots, k\}$. An instance of such an OCSP consists of a set of variables,*
 109 *$V = \{x_1, \dots, x_n\}$, and m constraints, C_1, \dots, C_m , each of which is an ordered k -tuple of*
 110 *V . The objective is to find a global ordering σ of V that maximizes $\sum_{i=1}^m \chi_\Lambda(\sigma_{|C_i})$, where*
 111 *$\sigma_{|C_i} \in S_k$ is the ordering of the k elements of C_i induced by the global ordering σ , and*
 112 *$\chi_\Lambda(\sigma_{|C_i}) = 1$ if $\sigma_{|C_i} \in \Lambda$ and 0 otherwise. If $\chi_\Lambda(\sigma_{|C_i}) = 1$, we say that C_i is satisfied.*

23:4 Finding an Optimal Alphabet Ordering for Lyndon Factorization is Hard

113 Note that $m \leq n!/(n-k)! \leq n^k$. Additionally, we will only consider OCSF instances
 114 where each variable appears in at least two constraints. Under this last assumption, we can
 115 relate the number of variables, n , to the number of clauses, m .

116 ► **Lemma 8.** *For OCSFs with arity k constraints, n variables, and m constraints, where
 117 every variable appears in at least two clauses, $n \leq \frac{k}{2}m$.*

118 **Proof.** Since every variable appears in at least two constraints,

$$119 \quad 2n \leq \sum_{i=1}^n (\text{the number of times variable } x_i \text{ appears in total}) = km. \quad \blacktriangleleft$$

120 One of the simplest OCSFs is the *Maximum Acyclic Subgraph Problem* (MAS), where
 121 $k = 2$, making constraints of the form (x_i, x_j) , and where Λ contains only the identity
 122 permutation that orders $x_i < x_j$. The dual minimization problem of MAS is known as
 123 *Feedback Arc Set* (FAS). In this problem, the cost of a solution is the number of constraints
 124 being violated, instead of the number of constraints being satisfied. The problem is otherwise
 125 identical. The following hardness result for FAS is used when proving Theorem 3.

126 ► **Lemma 9** ([14]). *Conditioned on the Unique Games Conjecture, for every constant $C > 1$,
 127 it is NP-hard to find a C -approximation for FAS.*

128 The Unique Games Conjecture is described in [21]. We will use the term Unique-Games-hard
 129 to refer to problems that, conditioned on the Unique Games conjecture, are NP-hard.

130 We can always assume that at least half of the constraints in an instance of MAS can
 131 be satisfied. To see this, take an arbitrary ordering of the variables. Either this ordering
 132 or its reversal must satisfy at least $m/2$ constraints. This is just a specific instance of a
 133 more general result. We can always assume our optimal solution satisfies at least $|\Lambda|m/k!$
 134 constraints. Since the expected number of constraints satisfied by a random ordering on the
 135 variables is $|\Lambda|m/k!$, we know the maximum number of constraints satisfied by any ordering
 136 is bounded below by this quantity. It turns out, however, that finding a solution that does
 137 better than this expected value is computationally difficult. We give a simplified statement
 138 of the main result in [14], maintaining only the pertinent details for our problem.

139 ► **Theorem 10** ([14]). *For an OCSF with arity k , for every constant $\varepsilon > 0$, it is Unique-
 140 Games-hard to find an ordering for the variables that achieves a ratio of satisfied constraints
 141 over total constraints that is at least $|\Lambda|/k! + \varepsilon$.*

142 Our results also make use of the OCSF known as the *Betweenness Problem*. In this problem
 143 $k = 3$ and Λ is of size two. For a constraint (x_i, x_j, x_k) to be satisfied either $x_i < x_j < x_k$
 144 or $x_k < x_j < x_i$. For example, the ordering $x_4 < x_5 < x_3 < x_2 < x_1$ satisfies the constraint
 145 (x_1, x_2, x_5) , but not the constraint (x_4, x_2, x_5) . By applying Theorem 10 to the Betweenness
 146 problem, we obtain that it is Unique-Games-hard to achieve a ratio of satisfied constraints
 147 to total constraints better than $2/3! = 1/3$.

148 For hardness under the Exponential Time Hypothesis (ETH) [19], we will use a result by
 149 Kim and Gonçalves appearing in [23]. An Arity k Permutation CSP as defined in [23] is a
 150 OCSF where Λ consists of only the identity permutation, i.e., a constraint $(x_{i_1}, x_{i_2}, \dots, x_{i_j})$,
 151 is satisfied iff it is ordered $x_{i_1} < x_{i_2} < \dots < x_{i_j}$, and constraints up to arity k are allowed.
 152 This is different from our definition of OCSFs, where all constraints are of exactly arity k .
 153 The differences between these two definitions are accommodated for whenever Lemma 11 is
 154 used. In [23] the authors prove the following.

155 ► **Lemma 11** ([23]). *Assuming ETH, there is no $2^{o(n \log n)}$ -algorithm for Arity 4 Permutation
 156 CSP (and thus for Arity k Permutation CSP, $k \geq 4$).*

3 Hardness of Lyndon Factor Minimization

The first reduction is from the Betweenness problem to the Lyndon Factor Minimization Problem. It is used to demonstrate NP-completeness. An alternative proof can be done with a reduction from MAS. Our reasoning for choosing one over the other is we believe that the Betweenness problem provides a good initial illustration of the power of a hypothetical solver to these Lyndon factorization problems. It also provides a warm-up for the techniques used in Section 3.2. Moreover, we will use a reduction from MAS as a short proof to illustrate NP-completeness for the maximization problem, before introducing a more involved reduction to prove an inapproximability result.

3.1 NP-Completeness of Lyndon Factor Minimization

We are given as input an instance ϕ of the Betweenness problem consisting of n variables x_1, x_2, \dots, x_n and m constraints C_1, C_2, \dots, C_m . Let $F(T)$ denote the number of Lyndon factors of a string T under the alphabet ordering currently under consideration. We will use $F_T(T_1)$ to denote the number of Lyndon factors of T starting within the *first occurrence* of the substring T_1 of T . The subscript T is to remind us that the factors starting in T_1 are sensitive to the other symbols in T . By a *run* of a symbol, we mean a maximal unary substring containing that symbol.

► **Lemma 12.** *Let T be any string of the form $T = T_1 \circ (x_0)^\alpha \circ (x_1^\gamma x_2^\gamma \dots x_n^\gamma)^\beta$ where T_1 is over the alphabet $\{x_0, \dots, x_n\}$, α is greater than the length of any run of x_0 in T_1 , γ is greater than the length of any run of any symbol other than x_0 in T_1 , and $\beta > 1$. If x_0 is the smallest symbol in the ordering, then $F(T) \leq F_T(T_1) + 1$.*

Proof. If T_1 does not end with an x_0 , then the first x_0 in the $(x_0)^\alpha$ marks the start of a new Lyndon factor in T since $(x_0)^\alpha$ is lexicographically smaller than any preceding suffix. Then this factor includes the remaining suffix of T . In this case $F(T) = F_T(T_1) + 1$. If T_1 contains a suffix consisting of only x_0 's, then a new Lyndon factor must start at the first of these x_0 's, and again this factor contains the remaining suffix of T . In this case, $F(T) = F_T(T_1)$. ◀

► **Lemma 13.** *Let T be defined as in Lemma 12. If x_0 is not the smallest symbol in the ordering, $F(T) \geq \beta - 1$.*

Proof. In this case, the smallest symbol must be one of x_1, \dots, x_n . Suppose the smallest is x_i . Then the first symbol in the first x_i^γ marks the beginning of a Lyndon factor. This factor is of the form $x_i^\gamma x_{i+1}^\gamma \dots x_n^\gamma x_1^\gamma \dots x_{i-1}^\gamma$ and is repeated at least $\beta - 1$ times. In particular, the suffix $x_{i+1}^\gamma \dots x_n^\gamma$ is preceded by $\beta - 1$ factors of the form $x_i^\gamma x_{i+1}^\gamma \dots x_n^\gamma x_1^\gamma \dots x_{i-1}^\gamma$. ◀

Lemmas 12 and 13 will be useful in proving that x_0 must be smallest in an optimal ordering. We now introduce our constraint gadgets.

► **Lemma 14.** *Let x_0 be the smallest symbol in T . For $i, j, k > 0$, consider the first instance of a substring S of T where*

$$S = x_0^\eta x_j x_0^\eta x_i x_0^\eta x_j x_0^\eta x_i x_0^\eta x_k x_0^\eta x_i x_0^\eta x_j x_0^\eta x_j x_0^\eta x_j$$

and η is larger than the length of any run of x_0 preceding S in T , and S is immediately followed by the run $x_0^{\eta+1}$. The symbols in this first instance of S , make up three complete Lyndon factors if x_j is ordered between x_i and x_k , and four complete Lyndon factors otherwise.

23:6 Finding an Optimal Alphabet Ordering for Lyndon Factorization is Hard

197 **Proof.** Since the number of times x_0 is repeated is more than the length of any previous
 198 run, it must be the case that a new factor begins at the start of S . The six possible cases
 199 and their corresponding factorizations are:

$$\begin{aligned}
 & x_0 < x_i < x_j < x_k : (x_0^\eta x_j), (x_0^\eta x_i x_0^\eta x_j x_0^\eta x_i x_0^\eta x_k), (x_0^\eta x_i x_0^\eta x_i x_0^\eta x_j x_0^\eta x_j x_0^\eta x_j) \\
 & x_0 < x_i < x_k < x_j : (x_0^\eta x_j), (x_0^\eta x_i x_0^\eta x_j), (x_0^\eta x_i x_0^\eta x_k), (x_0^\eta x_i x_0^\eta x_i x_0^\eta x_j x_0^\eta x_j x_0^\eta x_j) \\
 200 & x_0 < x_j < x_i < x_k : (x_0^\eta x_j x_0^\eta x_i x_0^\eta x_j x_0^\eta x_i x_0^\eta x_k x_0^\eta x_i x_0^\eta x_i), (x_0^\eta x_j), (x_0^\eta x_j), (x_0^\eta x_j) \\
 & x_0 < x_k < x_i < x_j : (x_0^\eta x_j), (x_0^\eta x_i x_0^\eta x_j), (x_0^\eta x_i), (x_0^\eta x_k x_0^\eta x_i x_0^\eta x_i x_0^\eta x_j x_0^\eta x_j x_0^\eta x_j) \\
 & x_0 < x_j < x_k < x_i : (x_0^\eta x_j x_0^\eta x_i x_0^\eta x_j x_0^\eta x_i x_0^\eta x_k x_0^\eta x_i x_0^\eta x_i), (x_0^\eta x_j), (x_0^\eta x_j), (x_0^\eta x_j) \\
 & x_0 < x_k < x_j < x_i : (x_0^\eta x_j x_0^\eta x_i), (x_0^\eta x_j x_0^\eta x_i), (x_0^\eta x_k x_0^\eta x_i x_0^\eta x_i x_0^\eta x_j x_0^\eta x_j x_0^\eta x_j)
 \end{aligned}$$

201 Notice that only in the first and last orderings where the constraint is satisfied are there
 202 three factors. The other cases have four. ◀

203 For each constraint $C_t = (x_i, x_j, x_k)$ in the instance ϕ of the Betweenness problem, where
 204 $1 \leq t \leq m$, we construct the gadget from Lemma 14,

$$205 \quad S(C_t) := x_0^t x_j x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_i x_0^t x_j x_0^t x_j x_0^t x_j.$$

206 We next define $S(\phi) := S(C_1) \circ S(C_2) \circ \dots \circ S(C_m) \circ (x_0)^{m+1} \circ (x_1^2 x_2^2 \dots x_n^2)^\beta$ where
 207 $\beta = 3m + 3$.

208 ▶ **Lemma 15.** *The string $S(\phi)$ has an alphabet ordering yielding at most $3m + 1$ Lyndon*
 209 *factors iff there exists a variable ordering satisfying all constraints in ϕ .*

210 **Proof.** Assuming there exists a constraint satisfying variable ordering for ϕ , make x_0 the
 211 smallest symbol and order the remaining symbols x_1, \dots, x_n according to the variable ordering.
 212 By Lemma 14, each of the substrings $S(C_t)$ for $1 \leq t \leq m$ contributes three factors, and by
 213 the analysis in Lemma 12 the remaining suffix contributes one additional factor. This creates
 214 $3m + 1$ factors in total.

215 Conversely, assume that no variable ordering exists that satisfies the constraints. If x_0 is
 216 the smallest symbol, then at least one $S(C_t)$ gadget contributes four factors while the others
 217 contribute at least three. The remaining suffix contributes one factor making the number of
 218 factors at least $4 + 3(m - 1) + 1 = 3m + 2$. If x_0 is not the smallest symbol, then by Lemma
 219 13, the number of factors is at least $\beta - 1 = (3m + 3) - 1 = 3m + 2$. ◀

220 Since determining if there exists a variable ordering satisfying all constraints in an instance
 221 of the Betweenness problem is NP-hard [32], determining whether there exists an alphabet
 222 order where there are at most $3m + 1$ Lyndon factors is NP-hard as well. With a symbol
 223 ordering as a polynomial sized certificate, the problem is clearly in NP, proving Theorem 1.

224 3.2 ETH Hardness of Lyndon Factor Minimization

225 Here we reduce Arity 4 Permutation CSP to Lyndon Factor Minimization. Assume for the
 226 moment that x_0 is the smallest symbol, and that each substring $S(C_t)$ (yet to be defined) is
 227 followed by a run of x_0 longer than any run of x_0 that precedes it.

228 For an arity 2 constraint $C_t = (x_i, x_j)$, we construct a string using the symbols x_0 ,
 229 x_i , and x_j that has either 3 or 4 factors depending on the ordering on the variables. We
 230 will demonstrate which orderings create which factorizations. The string we construct is

231 $S(C_t) = x_0^t x_i x_0^t x_i x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_i$, which has the factorizations for different
232 orderings,

233 Ordering	Factorization	# factors
234 $x_i < x_j$:	$(x_0^t x_i x_0^t x_i x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_i)$	3
235 $x_j < x_i$:	$(x_0^t x_i)(x_0^t x_i)(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_i)$	4

237 Slightly more involved are the strings to model arity 3 constraints $C_t = (x_i, x_j, x_k)$,
238 $S(C_t) = x_0^t x_i x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_i$, where

239 Ordering	Factorization	# factors
240 $x_i < x_j < x_k$:	$(x_0^t x_i x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_i)$	3
241 $x_i < x_k < x_j$:	$(x_0^t x_i x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_i)$	4
242 $x_j < x_i < x_k$:	$(x_0^t x_i)(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_i)$	4
243 $x_k < x_i < x_j$:	$(x_0^t x_i x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_i)$	4
244 $x_j < x_k < x_i$:	$(x_0^t x_i)(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_i)$	4
245 $x_k < x_j < x_i$:	$(x_0^t x_i)(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_i)$	4

247 The most involved is the gadget for an arity 4 constraint $C_t = (x_i, x_j, x_k, x_h)$,

248 $S(C_t) = x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i$
249 which has the following factorizations depending on the ordering given to its symbols,

Ordering ('<' omitted)	Factorization	#
x_i, x_j, x_k, x_h :	$(x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k)(x_0^t x_i x_0^t x_j)(x_0^t x_i)$	3
x_i, x_j, x_h, x_k :	$(x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_k)(x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k)(x_0^t x_i x_0^t x_j)(x_0^t x_i)$	4
x_i, x_k, x_j, x_h :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j)(x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)$	4
x_i, x_h, x_j, x_k :	$(x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_k)(x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)$	4
x_i, x_k, x_h, x_j :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j)(x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)$	4
x_i, x_h, x_k, x_j :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)$	4
x_j, x_i, x_k, x_h :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)$	4
x_j, x_i, x_h, x_k :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)$	4
x_k, x_i, x_j, x_h :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_h, x_i, x_j, x_k :	$(x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_k)(x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_k, x_i, x_h, x_j :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
250 x_h, x_i, x_k, x_j :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_j, x_k, x_i, x_h :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)$	4
x_j, x_h, x_i, x_k :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)$	4
x_k, x_j, x_i, x_h :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_h, x_j, x_i, x_k :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)(x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_k, x_h, x_i, x_j :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_h, x_k, x_i, x_j :	$(x_0^t x_i x_0^t x_j)(x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)(x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_j, x_k, x_h, x_i :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)$	4
x_j, x_h, x_k, x_i :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)$	4
x_k, x_j, x_h, x_i :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_h, x_j, x_k, x_i :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i)(x_0^t x_j x_0^t x_i)(x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
x_k, x_h, x_j, x_i :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i x_0^t x_h x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4
251 x_h, x_k, x_j, x_i :	$(x_0^t x_i)(x_0^t x_j x_0^t x_i)(x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)(x_0^t x_h x_0^t x_j x_0^t x_i x_0^t x_k x_0^t x_i x_0^t x_j x_0^t x_i)$	4

252 The string construction for the overall reduction is almost identical to the one for ϕ in
 253 Section 3.1. We only need to select β to be slightly different. We let $\beta = 4m + 3$. This
 254 is enough to ensure that in an optimal solution x_0 must be the smallest symbol. If x_0 is
 255 smallest, in the worst-case, when all constraints are not satisfied, there are at most $4m + 1$
 256 Lyndon factors. If x_0 is not smallest, as shown in Lemma 13, the number of factors is at
 257 least $\beta - 1 = 4m + 2$. Then, with x_0 as the minimum, each ordering on x_1, \dots, x_n gives us
 258 $3s + 4(m - s) + 1 = 4m + 1 - s$ factors, where s is the number of satisfied constraints when
 259 using the corresponding variable ordering in ϕ . Therefore, an optimal ordering for the n
 260 variables of ϕ is obtained by an order on the $(n + 1)$ symbols which minimizes the number of
 261 Lyndon factors in the string. This combined with Lemma 11 proves Theorem 2.

262 3.3 Inapproximability of Lyndon Factor Minimization

263 We will perform an approximation preserving reduction from FAS to Lyndon Factor Minim-
 264 ization. Recall that for FAS the arity k of the constraints is 2, so that constraints are of the
 265 form (x_i, x_j) and Λ consists of the identity permutation. In other words, the constraint is
 266 only satisfied if $x_i < x_j$. The cost of the solution will be the number of violated constraints,
 267 which we wish to minimize. Our gadget for constraint $C_t = (x_i, x_j)$ will be

$$268 \quad S(C_t) = (x_0^t x_i) \circ (x_0^t x_j)^{\alpha-1}$$

269 where $\alpha > 1$ will be chosen later. The whole string for our reduction will be

$$270 \quad T = S(\phi) = S(C_1) \circ S(C_2) \circ \dots \circ S(C_m) \circ (x_0)^{m+1} \circ (x_1^2 x_2^2 \dots x_n^2)^\beta$$

271 where $\beta = \alpha m + 3$. By Lemma 13, if x_0 is not smallest, then $F(T) \geq \beta - 1$. We consider
 272 next what happens in our constraint gadgets when x_0 is smallest.

273 ► **Lemma 16.** *If x_0 is smallest and $x_i < x_j$ then $F_T(S(C_t)) = 1$.*

274 **Proof.** Since x_0^t is the longest run of x_0 seen so far, the start of $S(C_t)$ marks the smallest
 275 suffix seen so far when traversing T from left to right. Then, since $x_j > x_i$, the start of all
 276 substrings of the form $x_0^t x_j$ do not mark the start of the smallest suffix seen so far. ◀

277 ► **Lemma 17.** *If x_0 is smallest and $x_j < x_i$ then $F_T(S(C_t)) = \alpha$.*

278 **Proof.** Again, since x_0^t is the longest run of x_0 seen so far, the start of $S(C_t)$ marks the
 279 smallest suffix seen so far when traversing T from left to right. However, now the start of
 280 each substring of the form $x_0^t x_j$ marks the start of the smallest suffix seen so far (recall after
 281 the last $x_0^t x_j$ there will be a longer run of x_0 than has been seen before). Hence, there are
 282 $\alpha - 1$ additional factors created. ◀

283 ► **Lemma 18.** *Any alphabet ordering where x_0 is smallest has fewer factors than an alphabet
 284 ordering where x_0 is not the smallest.*

285 **Proof.** If x_0 is smallest, $F(T) = s + \alpha(m - s) + 1$ where s is the number of satisfied constraints
 286 and the $+1$ arises from the last factor, $(x_0)^{m+1} \circ (x_1^2 x_2^2 \dots x_n^2)^\beta$. Because $\alpha > 1$, this is
 287 upper bounded by the case when $s = 0$ so that $F(T) \leq \alpha m + 1$. On the other hand, if x_0 is
 288 not smallest $F(T) \geq \beta - 1 = \alpha m + 2$. ◀

289 Henceforth, we only need to worry about the case when x_0 is the smallest. Our aim is to
 290 show that a constant approximation algorithm for Lyndon Factor Minimization allows us to
 291 construct a constant approximation algorithm for FAS. If our hypothetical approximation
 292 algorithm for Lyndon Factor Minimization ever returned a solution where x_0 is not smallest,
 293 we add the additional step of replacing that solution with any solution where x_0 is smallest,
 294 obtaining a solution that performs even better. Then our modified algorithm maintains being
 295 an approximation algorithm for Lyndon Factor Minimization (perhaps with an even smaller
 296 approximation factor).

297 Let s_F^* denote the number of constraints satisfied in an optimal solution of ϕ for FAS and
 298 let s_L^* denote the number of constraints in ϕ satisfied by the variable ordering obtained from
 299 our optimal, factor minimizing, alphabet order for the corresponding instance of Lyndon
 300 Factor Minimization. Also, let s denote the actual number of constraints satisfied by the
 301 variable ordering obtained from our approximate factor minimizing alphabet order for the
 302 corresponding instance of Lyndon Factor Minimization. A λ -approximation for Lyndon
 303 Factor Minimization with $\lambda > 1$ gives the following set of inequalities:

$$304 \quad s_L^* + \alpha(m - s_L^*) + 1 \leq s + \alpha(m - s) + 1 \leq \lambda(s_L^* + \alpha(m - s_L^*) + 1).$$

305 Which can be equivalently written as

$$306 \quad (m - s_L^*) + \frac{s_L^* + 1}{\alpha} \leq (m - s) + \frac{s + 1}{\alpha} \leq \lambda(m - s_L^*) + \lambda \frac{s_L^* + 1}{\alpha}. \quad (1)$$

307 We will show that by taking α large enough we can ensure $s_L^* = s_F^*$.

308 **► Lemma 19.** *With $\alpha = 2(m + 1) + 1$, we have that $s_L^* = s_F^*$.*

309 **Proof.** The cost of a solution of ϕ is of the form $m - s_F^*$. The solution for ϕ we get from
 310 mapping our solution for Lyndon factorization back to ϕ must have at least as many violated
 311 constraints as the optimal solution for ϕ , i.e., $m - s_L^* \geq m - s_F^*$, and so $s_F^* \geq s_L^*$. Let us
 312 suppose for the sake of contradiction that $s_F^* \geq s_L^* + 1$. This implies $m - s_L^* - (m - s_F^*) \geq 1$.
 313 Then, using in addition that $\frac{s_F^* + 1}{\alpha} \leq \frac{m + 1}{\alpha} \leq \frac{1}{2}$, we obtain

$$314 \quad \frac{s_F^* + 1}{\alpha} - \frac{s_L^* + 1}{\alpha} \leq \frac{1}{2} < 1 \leq m - s_L^* - (m - s_F^*),$$

315 which implies that

$$316 \quad m - s_F^* + \frac{s_F^* + 1}{\alpha} < m - s_L^* + \frac{s_L^* + 1}{\alpha}.$$

317 Or, written more naturally as the cost of a Lyndon Factor Minimization Problem's solution,

$$318 \quad s_F^* + \alpha(m - s_F^*) + 1 < s_L^* + \alpha(m - s_L^*) + 1.$$

319 But then this implies that the ordering on x_1, \dots, x_n that is used to obtain the optimal
 320 solution for ϕ creates fewer Lyndon factors than our supposedly optimal solution for Lyndon
 321 Factor Minimization, a contradiction. ◀

322 Let us now upper bound $m - s$ (our approximate solution cost when the solution is
 323 mapped back to FAS) in terms of $\lambda(m - s_F^*)$. Combining the inequalities in (1) with Lemma
 324 19, and the fact that $s_F^* = s_L^* \leq m$ when $\alpha = 2(m + 1) + 1$, we get that

$$325 \quad m - s \leq m - s + \frac{s + 1}{\alpha} \leq \lambda(m - s_L^*) + \lambda \frac{s_L^* + 1}{\alpha} \leq \lambda \left(m - s_F^* + \frac{1}{2} \right).$$

23:10 Finding an Optimal Alphabet Ordering for Lyndon Factorization is Hard

326 The case where $m = s_F^*$ can easily be solved in polynomial time, so we can consider that
327 check added to our hypothetical solution as well. Hence, we assume $m - s_F^* \geq 1 > 1/2$ and,

$$328 \quad m - s_F^* \leq m - s \leq \lambda \left(m - s_F^* + \frac{1}{2} \right) < \lambda(m - s_F^* + m - s_F^*) = 2\lambda(m - s_F^*).$$

329 We have shown that a λ approximation for Lyndon Factor Minimization allows us to obtain,
330 at worst, a 2λ approximation for FAS. Moreover, the α value we need to do this is polynomial
331 in m so that the whole reduction is done in polynomial time. This polynomial time constant
332 approximation algorithm is better than what is allowed by Lemma 9 under the Unique Games
333 Conjecture. This completes the proof of Theorem 3.

334 **4 Hardness of Lyndon Factor Maximization**

335 Our approach will be similar to the one taken for minimization. First, we introduce
336 some gadgetry for the NP-completeness proof that is later expanded upon to create an
337 inapproximability result. As of now, the authors have not yet found gadgets to establish the
338 same ETH hardness for the maximization problem.

339 **4.1 NP-Completeness of Lyndon Factor Maximization**

340 We perform a reduction from the dual of FAS, the Maximum Acyclic Subgraph Problem
341 (MAS). Recall MAS is identical to FAS except for the cost of a solution now being the number
342 of constraints satisfied, which we wish to maximize. For constraint $C_t = (x_i, x_j)$, we define
343 our constraint gadget as $S(C_t) = x_0^{t+1} x_j x_0^{t+1} x_i$ (note the reversal of i and j). The entire
344 string formed by our instance ϕ of FAS is

$$345 \quad T = S(\phi) = (x_0 x_1 x_2 \dots x_n) \circ S(C_1) \circ S(C_2) \circ \dots \circ S(C_m) \circ (x_0)^m.$$

346 **► Lemma 20.** *If x_0 is not the smallest symbol in the ordering, then $F(T) \leq n + m$.*

347 **Proof.** Suppose $x_i \neq x_0$ is the smallest symbol. Then the first Lyndon factor starting with
348 x_i occurs in the prefix $(x_0 x_1 \dots x_n)$. Subsequent Lyndon factors must begin with x_i . The
349 prefix contributes at most n factors and there are at most m remaining occurrences of x_i . ◀

350 **► Lemma 21.** *In an ordering where x_0 is smallest, $F(T) = 2s + (m - s) + 1 + m$, where s
351 is the number of constraints satisfied in MAS by the ordering given to x_1, \dots, x_n .*

352 **Proof.** For a substring $S(C_t)$, if $C_t = (x_i, x_j)$ is not satisfied (i.e., $x_i > x_j$) then $F_T(S(C_t)) =$
353 1. If it is satisfied (i.e., $x_i < x_j$) then $F_T(S(C_t)) = 2$. The prefix $x_0 x_1 x_2 \dots x_n$ contributes
354 exactly one additional factor. The suffix $(x_0)^m$ contributes m factors. ◀

355 **► Lemma 22.** *Any ordering where x_0 is the smallest has more factors than an ordering
356 where x_0 is not the smallest.*

357 **Proof.** By Lemma 8, we can assume that $n \leq m$. Then by Lemma 20, we have that if
358 x_0 is not smallest, $F(T) \leq n + m \leq 2m$. By Lemma 21, if x_0 is smallest then $F(T) =$
359 $2s + (m - s) + 1 + m = s + 2m + 1 > 2m$. ◀

360 The value $F(T)$ is maximized by an alphabet order which has the largest possible number
361 of satisfied constraints, say s^* . This gives $(s^* + 2m + 1)$ Lyndon factors. Clearly, this solution
362 also provides an ordering satisfying the maximum number of constraints in our MAS instance.
363 Since MAS is NP-hard, we have shown Lyndon Factor Maximization is NP-hard as well. The

364 decision problem is in NP using the ordering on $x_1 \dots x_n$ as a polynomial sized certificate,
 365 and this remains NP-hard as it could be used to solve the optimization problem. This
 366 completes the proof of Theorem 4.

367 4.2 Inapproximability of Lyndon Factor Maximization

368 First, let us describe the OCSP from which we are reducing. Let $k > 1$ be the arity of the
 369 constraints, which we will specify later. Each constraint will be satisfied iff the variables
 370 in that constraint have one of the $(k - 1)!$ orderings where the last variable is ordered
 371 first, i.e., for constraint $(x_{i_1}, x_{i_2}, \dots, x_{i_{k-1}}, x_{i_k})$, the ordering over those variables will have
 372 $x_{i_k} < x_{i_j}$ for $j \in [1, k - 1]$. According to Theorem 10, it is Unique-Games-Hard to find an
 373 approximation which beats $|\Lambda|m/k! = (k - 1)!m/k! = m/k$ constraints being satisfied.

374 Our constraint gadget is of the form

$$375 \quad S(C_t) = (x_0^{t+1} x_{i_1}) \circ (x_0^{t+1} x_{i_2}) \circ \dots \circ (x_0^{t+1} x_{i_{k-1}}) \circ (x_0^{t+1} x_{i_k})^\alpha$$

376 and our overall string constructed from our instance ϕ of OCSP is

$$377 \quad T := S(\phi) = (x_0 x_1 x_2 \dots x_n) \circ S(C_1) \circ S(C_2) \circ \dots \circ S(C_m) \circ (x_0), \quad \text{where } \alpha = mn.$$

378 **► Lemma 23.** *If x_0 is not smallest then $F(T) \leq n + m$.*

380 **Proof.** Let $x_i \neq x_0$ be the smallest symbol instead. Then the prefix $(x_0 x_1 x_2 \dots x_n)$
 381 contributes at most n factors, and each remaining factor must begin with x_i . We will show
 382 that there is at most 1 factor starting in each constraint gadget. For a given constraint
 383 containing x_i , if $x_i \neq x_{i_k}$ this is immediate. On the other hand, if $x_i = x_{i_k}$ then only its
 384 first occurrence can form a smaller suffix of T than those preceding it. In more detail, since
 385 $x_0 > x_i = x_{i_k}$, we have $x_{i_k} (x_0^t x_{i_k})^{\alpha-1} x_0 < x_{i_k} (x_0^t x_{i_k})^{\alpha-2} x_0 < x_{i_k} (x_0^t x_{i_k})^{\alpha-3} x_0 < \dots$
 386 Note that this is the reason for the final x_0 appended to T . ◀

387 **► Lemma 24.** *If x_0 is smallest, and in constraint $C_t = (x_{i_1} \dots x_{i_k})$ the symbol x_{i_k} is smallest
 388 among $x_{i_1} \dots x_{i_k}$, then $F_T(S(C_t)) \geq \alpha$.*

389 **Proof.** Since $x_0^{t+1} x_{i_k} < x_0^{t+1} x_{i_j}$ for $j \in [1, k - 1]$, and the string following $S(C_t)$ is either
 390 x_0^{t+2} (or x_0 then the empty string), the start of **each run** of x_0 in the substring $(x_0^{t+1} x_{i_k})^\alpha$
 391 marks the start of a suffix smaller than any of those preceding it. ◀

392 **► Lemma 25.** *If x_0 is the smallest in the ordering, then $F(T) \geq \alpha s + 1$ where s is the
 393 number clauses in ϕ satisfied by the ordering given to $x_1 \dots, x_n$. This is larger than the
 394 number of factors from any ordering where x_0 is not the smallest.*

395 **Proof.** By Lemma 24, when x_0 is the smallest each of the satisfied constraint gadgets
 396 contributes at least α factors. In addition, the lone x_0 symbol at the end of T forms its own
 397 factor. For the second statement, we can always assume our approximate solution satisfies at
 398 least 1 constraint, hence $s \geq 1$ and $\alpha s + 1 \geq mn + 1 > m + n$, which by Lemma 23 is an
 399 upper bound on the number of factors when x_0 is not smallest. ◀

400 From here we only need to consider when x_0 is smallest, for the same reasoning as given
 401 in Section 3.3. Now, suppose we have a λ -approximation with $\lambda < 1$ for Lyndon Factor
 402 Maximization. Let s_L^* be the number of constraint gadgets satisfied from our optimal solution
 403 of Lyndon factor maximization, and s the number from the approximate solution. Then,

$$404 \quad \lambda(\alpha s_L^* + 1 + y_L^*) \leq \alpha s + 1 + y \leq \alpha s_L^* + 1 + y_L^*$$

23:12 Finding an Optimal Alphabet Ordering for Lyndon Factorization is Hard

405 where y_L^* represents the number of additional factors contributed beyond $\alpha s_L^* + 1$ and y
 406 represents the number of factors beyond $\alpha s + 1$ for our approximate solution. We can
 407 equivalently write the above expression as

$$408 \quad \lambda s_L^* \left(1 + \frac{1}{\alpha s_L^*} + \frac{y_L^*}{\alpha s_L^*} \right) \leq s \left(1 + \frac{1}{\alpha s} + \frac{y}{\alpha s} \right) \leq s_L^* \left(1 + \frac{1}{\alpha s_L^*} + \frac{y_L^*}{\alpha s_L^*} \right). \quad (2)$$

409 ► **Lemma 26.** *For all $s \in [1, m]$, and for the corresponding y value as described above,*

$$410 \quad 1 \leq \left(1 + \frac{1}{\alpha s} + \frac{y}{\alpha s} \right) \leq 3.$$

411 **Proof.** We first bound y from above. Any factor in a constraint gadget begins at the start of
 412 a run x_0 . In a satisfied constraint gadget, there are $k - 1$ such runs outside of the $(x_0^{t+1} x_{i_k})^\alpha$
 413 substring. Hence, each satisfied constraint gadget contributes at most $k - 1$ additional factors
 414 beyond α . A constraint gadget that is not satisfied, i.e., has $x_{i_j} < x_{i_k}$ for some $j \neq k$, has
 415 the gadget's last factor beginning at the start of the substring $(x_0^{t+1} x_{i_j})$. This implies the
 416 substring $(x_0^{t+1} x_{i_k})^\alpha$ does not split into different factors. Therefore, an unsatisfied constraint
 417 gadget again contributes at most $k - 1$ factors. Because of this, the m constraint gadgets
 418 contribute at most k additional factors in total and $y \leq m(k - 1)$. Finally, $\alpha = mn$, hence

$$419 \quad \frac{y}{\alpha s} \leq \frac{y}{\alpha} \leq \frac{m(k-1)}{\alpha} \leq \frac{mn}{\alpha} = 1 \quad \text{and} \quad \frac{1}{\alpha s} \leq \frac{1}{\alpha} = \frac{1}{nm} \leq 1. \quad \blacktriangleleft$$

420 Let s_C^* be the number of constraints satisfied in an optimal solution to ϕ . Like in Section
 421 3.3, we know that $s \leq s_C^*$ and $s_L^* \leq s_C^*$. Using Lemma 26 we can easily make them differ by
 422 at most a constant factor.

423 ► **Lemma 27.** *Using the definitions above, it holds that $s_C^* \leq 3s_L^*$.*

424 **Proof.** For the sake of contradiction, assume instead that $s_C^* > 3s_L^*$. Applying the ordering
 425 given by the optimal solution of ϕ to the symbols x_1, \dots, x_n , and letting y_C^* be defined as
 426 above but for s_C^* , we have

$$427 \quad s_C^* \left(1 + \frac{1}{\alpha s_C^*} + \frac{y_C^*}{\alpha s_C^*} \right) > s_C^* > 3s_L^* \geq s_L^* \left(1 + \frac{1}{\alpha s_L^*} + \frac{y_L^*}{\alpha s_L^*} \right)$$

428 However, this implies $\alpha s_C^* + 1 + y_C^* > \alpha s_L^* + 1 + y_L^*$. Thus, s_L^* couldn't have been the number
 429 of constraints satisfied in an optimal solution to our Lyndon Factor Maximization instance,
 430 since using whichever ordering was used for the solution to ϕ would have given us more
 431 factors, a contradiction. \blacktriangleleft

432 By Lemma 27, we have $\frac{1}{3}s_C^* \leq s_L^*$. Multiplying both sides by $\lambda/3$, we obtain $\frac{\lambda}{9}s_C^* \leq \frac{\lambda}{3}s_L^*$.
 433 By Lemma 26 and our starting inequality in (2) we also have that

$$434 \quad \lambda s_L^* \leq \lambda s_L^* \left(1 + \frac{1}{\alpha s_L^*} + \frac{y_L^*}{\alpha s_L^*} \right) \leq s \left(1 + \frac{1}{\alpha s} + \frac{y}{\alpha s} \right) \leq 3s.$$

435 From which we obtain $\frac{\lambda}{3}s_L^* \leq s$. Combining these inequalities with the fact that $s \leq s_C^*$, we
 436 get $\frac{\lambda}{9}s_C^* \leq s \leq s_C^*$. That is, a λ -approximation algorithm for Lyndon Factor Maximization
 437 provides at least a $\lambda/9$ -approximation algorithm for this set of OCSP problems.

438 To finish the proof of Theorem 5, suppose for the sake of contradiction there exists
 439 a λ -approximation algorithm for Lyndon factor maximization for some constant $\lambda < 1$.
 440 Consider the set of OCSPs problems described in beginning of Section 4.2 with arity k such
 441 that $1/k < \lambda/9$. With our reduction, we obtain a polynomial-time algorithm that can find
 442 a solution with approximation ratio better than $|\Lambda|/k! = 1/k$, proving the Unique Games
 443 Conjecture false by Theorem 10.

444 — References —

- 445 1 Hideo Bannai, Tomohiro I, Shunsuke Inenaga, Yuto Nakashima, Masayuki Takeda, and
 446 Kazuya Tsuruta. The "runs" theorem. *SIAM J. Comput.*, 46(5):1501–1514, 2017. doi:
 447 10.1137/15M1011032.
- 448 2 Hideo Bannai, Juha Kärkkäinen, Dominik Köppl, and Marcin Piatkowski. Constructing the
 449 bijective BWT. *CoRR*, abs/1911.06985, 2019. URL: <http://arxiv.org/abs/1911.06985>,
 450 arXiv:1911.06985.
- 451 3 Hideo Bannai, Juha Kärkkäinen, Dominik Köppl, and Marcin Piatkowski. Indexing the
 452 bijective BWT. In *30th Annual Symposium on Combinatorial Pattern Matching, CPM 2019,*
 453 *June 18-20, 2019, Pisa, Italy*, pages 17:1–17:14, 2019. doi:10.4230/LIPIcs.CPM.2019.17.
- 454 4 Jason W. Bentley, Daniel Gibney, and Sharma V. Thankachan. On the complexity of bwt-runs
 455 minimization via alphabet reordering. In *28th Annual European Symposium on Algorithms,*
 456 *ESA 2020, September 7-9, 2020, Pisa, Italy (Virtual Conference)*, pages 15:1–15:13, 2020.
 457 doi:10.4230/LIPIcs.ESA.2020.15.
- 458 5 Moses Charikar, Venkatesan Guruswami, and Rajsekar Manokaran. Every permutation CSP
 459 of arity 3 is approximation resistant. In *Proceedings of the 24th Annual IEEE Conference on*
 460 *Computational Complexity, CCC 2009, Paris, France, 15-18 July 2009*, pages 62–73, 2009.
 461 doi:10.1109/CCC.2009.29.
- 462 6 Moses Charikar, Konstantin Makarychev, and Yury Makarychev. On the advantage over
 463 random for maximum acyclic subgraph. In *48th Annual IEEE Symposium on Foundations*
 464 *of Computer Science (FOCS 2007), October 20-23, 2007, Providence, RI, USA, Proceedings,*
 465 pages 625–633, 2007. doi:10.1109/FOCS.2007.47.
- 466 7 Kuo Tsai Chen, Ralph H Fox, and Roger C Lyndon. Free differential calculus, iv. the quotient
 467 groups of the lower central series. *Annals of Mathematics*, pages 81–95, 1958.
- 468 8 Amanda Clare and Jacqueline W. Daykin. Enhanced string factoring from alphabet orderings.
 469 *Inf. Process. Lett.*, 143:4–7, 2019. doi:10.1016/j.ipl.2018.10.011.
- 470 9 Amanda Clare, Jacqueline W. Daykin, Thomas Mills, and Christine Zarges. Evolutionary
 471 search techniques for the lyndon factorization of biosequences. In *Proceedings of the Genetic*
 472 *and Evolutionary Computation Conference Companion, GECCO 2019, Prague, Czech Republic,*
 473 *July 13-17, 2019*, pages 1543–1550, 2019. doi:10.1145/3319619.3326872.
- 474 10 Maxime Crochemore and Dominique Perrin. Two-way string matching. *J. ACM*, 38(3):651–675,
 475 1991. doi:10.1145/116825.116845.
- 476 11 Jean-Pierre Duval. Génération d’une section des classes de conjugaison et arbre des mots de lyn-
 477 don de longueur bornée. *Theor. Comput. Sci.*, 60:255–283, 1988. doi:10.1016/0304-3975(88)
 478 90113-2.
- 479 12 Isamu Furuya, Yuto Nakashima, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki
 480 Takeda. Lyndon factorization of grammar compressed texts revisited. In Gonzalo Navarro,
 481 David Sankoff, and Binhai Zhu, editors, *Annual Symposium on Combinatorial Pattern Match-*
 482 *ing, CPM 2018, July 2-4, 2018 - Qingdao, China*, volume 105 of *LIPIcs*, pages 24:1–24:10.
 483 Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018. doi:10.4230/LIPIcs.CPM.2018.24.
- 484 13 Joseph Yossi Gil and David Allen Scott. A bijective string sorting transform. *CoRR*,
 485 abs/1201.3077, 2012. URL: <http://arxiv.org/abs/1201.3077>, arXiv:1201.3077.
- 486 14 Venkatesan Guruswami, Johan Håstad, Rajsekar Manokaran, Prasad Raghavendra, and Moses
 487 Charikar. Beating the random ordering is hard: Every ordering CSP is approximation resistant.
 488 *SIAM J. Comput.*, 40(3):878–914, 2011. doi:10.1137/090756144.
- 489 15 Venkatesan Guruswami and Yuan Zhou. Approximating bounded occurrence ordering csp. In
 490 *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*
 491 *- 15th International Workshop, APPROX 2012, and 16th International Workshop, RANDOM*
 492 *2012, Cambridge, MA, USA, August 15-17, 2012. Proceedings*, pages 158–169, 2012. doi:
 493 10.1007/978-3-642-32512-0_14.

- 494 16 Johan Håstad. Some optimal inapproximability results. In *Proceedings of the Twenty-Ninth Annual ACM Symposium on the Theory of Computing, El Paso, Texas, USA, May 4-6, 1997*, pages 1–10, 1997. doi:10.1145/258533.258536.
- 497 17 Christophe Hohlweg and Christophe Reutenauer. Lyndon words, permutations and trees. *Theor. Comput. Sci.*, 307(1):173–178, 2003. doi:10.1016/S0304-3975(03)00099-9.
- 499 18 Tomohiro I, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Faster lyndon factorization algorithms for SLP and LZ78 compressed text. *Theor. Comput. Sci.*, 656:215–224, 2016. doi:10.1016/j.tcs.2016.03.005.
- 502 19 Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-sat. *J. Comput. Syst. Sci.*, 62(2):367–375, 2001. doi:10.1006/jcss.2000.1727.
- 504 20 Juha Kärkkäinen, Dominik Kempa, Yuto Nakashima, Simon J. Puglisi, and Arseny M. Shur. On the size of lempel-ziv and lyndon factorizations. In Heribert Vollmer and Brigitte Vallée, editors, *34th Symposium on Theoretical Aspects of Computer Science, STACS 2017, March 8-11, 2017, Hannover, Germany*, volume 66 of *LIPICs*, pages 45:1–45:13. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2017. doi:10.4230/LIPICs.STACS.2017.45.
- 509 21 Subhash Khot. On the unique games conjecture. In *46th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2005), 23-25 October 2005, Pittsburgh, PA, USA, Proceedings*, page 3, 2005. doi:10.1109/SFCS.2005.61.
- 512 22 Dong Kyue Kim, Jeong Seop Sim, Heejin Park, and Kunsoo Park. Linear-time construction of suffix arrays. In *Combinatorial Pattern Matching, 14th Annual Symposium, CPM 2003, Morelia, Michocán, Mexico, June 25-27, 2003, Proceedings*, pages 186–199, 2003. doi:10.1007/3-540-44888-8_14.
- 516 23 Eun Jung Kim and Daniel Gonçalves. On exact algorithms for the permutation CSP. *Theor. Comput. Sci.*, 511:109–116, 2013. doi:10.1016/j.tcs.2012.10.035.
- 518 24 Manfred Kufleitner. On bijective variants of the burrows-wheeler transform. In *Proceedings of the Prague Stringology Conference 2009, Prague, Czech Republic, August 31 - September 2, 2009*, pages 65–79, 2009. URL: <http://www.stringology.org/event/2009/p07.html>.
- 521 25 Pierre Lalonde and Arun Ram. Standard lyndon bases of lie algebras and enveloping algebras. *Transactions of the American Mathematical Society*, 347(5):1821–1830, 1995.
- 523 26 M. Lothaire. *Combinatorics on words*, volume 17. Cambridge university press, 1997.
- 524 27 Lily Major, Amanda Clare, Jacqueline W. Daykin, Benjamin Mora, Leonel Jose Peña Gamboa, and Christine Zarges. Evaluation of a permutation-based evolutionary framework for lyndon factorizations. In *Parallel Problem Solving from Nature - PPSN XVI - 16th International Conference, PPSN 2020, Leiden, The Netherlands, September 5-9, 2020, Proceedings, Part I*, pages 390–403, 2020. doi:10.1007/978-3-030-58112-1_27.
- 529 28 Sabrina Mantaci, Antonio Restivo, Giovanna Rosone, and Marinella Sciortino. Sorting suffixes of a text via its lyndon factorization. In Jan Holub and Jan Zdárek, editors, *Proceedings of the Prague Stringology Conference 2013, Prague, Czech Republic, September 2-4, 2013*, pages 119–127. Department of Theoretical Computer Science, Faculty of Information Technology, Czech Technical University in Prague, 2013. URL: <http://www.stringology.org/event/2013/p11.html>.
- 535 29 Sabrina Mantaci, Antonio Restivo, Giovanna Rosone, and Marinella Sciortino. Suffix array and lyndon factorization of a text. *J. Discrete Algorithms*, 28:2–8, 2014. doi:10.1016/j.jda.2014.06.001.
- 538 30 Marcin Mucha. Lyndon words and short superstrings. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, pages 958–972, 2013. doi:10.1137/1.9781611973105.69.
- 541 31 Alantha Newman. Cuts and orderings: On semidefinite relaxations for the linear ordering problem. In *Approximation, Randomization, and Combinatorial Optimization, Algorithms and Techniques, 7th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems, APPROX 2004, and 8th International Workshop on Randomization*

- 545 *and Computation, RANDOM 2004, Cambridge, MA, USA, August 22-24, 2004, Proceedings,*
546 *pages 195–206, 2004. doi:10.1007/978-3-540-27821-4_18.*
- 547 **32** Jaroslav Opatrny. Total ordering problem. *SIAM J. Comput.*, 8(1):111–114, 1979. doi:
548 10.1137/0208008.
- 549 **33** Prasad Raghavendra. Optimal algorithms and inapproximability results for every csp? In
550 *Proceedings of the 40th Annual ACM Symposium on Theory of Computing, Victoria, British*
551 *Columbia, Canada, May 17-20, 2008, pages 245–254, 2008. doi:10.1145/1374376.1374414.*
- 552 **34** Kazuya Tsuruta, Dominik Köppl, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and
553 Masayuki Takeda. Grammar-compressed self-index with lyndon words. *CoRR*, abs/2004.05309,
554 2020. URL: <https://arxiv.org/abs/2004.05309>, arXiv:2004.05309.
- 555 **35** Yuki Urabe, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. On the
556 size of overlapping lempel-ziv and lyndon factorizations. In Nadia Pisanti and Solon P. Pissis,
557 editors, *30th Annual Symposium on Combinatorial Pattern Matching, CPM 2019, June 18-20,*
558 *2019, Pisa, Italy*, volume 128 of *LIPICs*, pages 29:1–29:11. Schloss Dagstuhl - Leibniz-Zentrum
559 für Informatik, 2019. doi:10.4230/LIPICs.CPM.2019.29.