

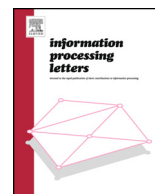


ELSEVIER

Contents lists available at ScienceDirect

Information Processing Letters

www.elsevier.com/locate/ipl

Enhanced string factoring from alphabet orderings^{☆, ☆☆}Amanda Clare^{a,*}, Jacqueline W. Daykin^{a,b,c}^a Department of Computer Science, Aberystwyth University, SY23 3DB, UK^b Department of Informatics, King's College London, WC2B 4BG, UK^c Department of Information Science, Stellenbosch University, South Africa

ARTICLE INFO

Article history:

Received 15 June 2018

Received in revised form 10 October 2018

Accepted 21 October 2018

Available online xxxx

Communicated by Kun-Mao Chao

Keywords:

Alphabet ordering

Combinatorial problems

Design of algorithms

Greedy algorithm

Lyndon factorization

ABSTRACT

In this note we consider the concept of alphabet ordering in the context of string factoring. We propose a greedy algorithm that produces Lyndon factorizations with small numbers of factors which can be modified to produce large numbers of factors. For the technique we introduce the Exponent Parikh vector. Applications and research directions derived from circ-UMFFs are discussed.

© 2018 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Factoring strings is a powerful form of the divide and conquer problem-solving paradigm for strings or words. Notably, the Lyndon factorization [1] is both efficient to compute [5] and useful in practice [7]. We study the effect of an alphabet's order on the number of factors in a Lyndon factorization and propose a greedy algorithm that assigns an ordering to the alphabet. In addition, we formalize the distinction between the sets of Lyndon and co-Lyndon words [3] as avenues for alternative string factorizations. More generally, circ-UMFFs provide the opportunity for achieving further diversity with string factors [2,3].

[☆] The authors were part-funded by the European Regional Development Fund through the Welsh Government, grant 80761-AU-137 (West).

^{☆☆} A preliminary version of this paper was accepted as a poster in IWOCA 2018 (International Workshop on Combinatorial Algorithms).

* Corresponding author.

E-mail addresses: afc@aber.ac.uk (A. Clare), jwd6@aber.ac.uk, jackie.daykin@kcl.ac.uk (J.W. Daykin).

<https://doi.org/10.1016/j.ipl.2018.10.011>

0020-0190/© 2018 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1.1. Notation

Given an integer $n \geq 1$ and a nonempty set of symbols Σ (bounded or unbounded), a **string of length n** , equivalently **word**, over Σ takes the form $\mathbf{x} = x_1 \dots x_n$ with each $x_i \in \Sigma$. For brevity, we write $\mathbf{x} = \mathbf{x}[1..n]$ with $\mathbf{x}[i] = x_i$. The length n of a string \mathbf{x} is denoted by $|\mathbf{x}|$. The set Σ is called an **alphabet** whose members are **letters** or **characters**, and Σ^+ denotes the set of all nonempty finite strings over Σ . The **empty string** of length zero is denoted ϵ ; we write $\Sigma^* = \Sigma^+ \cup \{\epsilon\}$ and let $|\Sigma| = \sigma$. We use exponents to denote repetition, for instance if $\alpha \in \Sigma$ then α^3 means $\alpha\alpha\alpha$. If $\mathbf{x} = \mathbf{u}\mathbf{w}\mathbf{v}$ for strings $\mathbf{u}, \mathbf{w}, \mathbf{v} \in \Sigma^*$, then \mathbf{u} is a **prefix**, \mathbf{w} is a **substring** or **factor**, and \mathbf{v} is a **suffix** of \mathbf{x} ; we say $\mathbf{u} \neq \mathbf{x}$ is a **proper prefix** and similarly for the other terms. If $\mathbf{x} = \mathbf{u}\mathbf{v}$, then $\mathbf{v}\mathbf{u}$ is said to be a **rotation** (**cyclic shift** or **conjugate**) of \mathbf{x} . A string \mathbf{x} is said to be a **repetition** if and only if it has a factorization $\mathbf{x} = \mathbf{u}^k$ for some integer $k > 1$; otherwise, \mathbf{x} is said to be **primitive**. For a string \mathbf{x} , the reversed string $\bar{\mathbf{x}}$ is defined as $\bar{\mathbf{x}} = \mathbf{x}[n]\mathbf{x}[n-1] \dots \mathbf{x}[1]$. A string which is both a proper prefix and a proper suffix of a string $\mathbf{x} \neq \epsilon$ is called a **border** of \mathbf{x} ; a string is **border-free** if the only border it has is the empty string ϵ .

If Σ is a totally ordered alphabet then **lexicographic ordering** (**lexorder**) $\mathbf{u} < \mathbf{v}$ with $\mathbf{u}, \mathbf{v} \in \Sigma^+$ means that either \mathbf{u} is a proper prefix of \mathbf{v} , or $\mathbf{u} = \mathbf{ras}$, $\mathbf{v} = \mathbf{rbt}$ for some $a, b \in \Sigma$ such that $a < b$ and for some $\mathbf{r}, \mathbf{s}, \mathbf{t} \in \Sigma^*$. We call the ordering $<$ based on lexorder of reversed strings **co-lexicographic ordering** (**co-lexorder**). Using the Roman alphabet: $ca < cat < dog$ while $dog < cat$.

2. Unique Maximal Factorization Families (UMFFs)

A subset $\mathcal{W} \subseteq \Sigma^+$ is a **factorization family** (FF) if and only if for every nonempty string \mathbf{x} on Σ there exists a factorization $F_{\mathcal{W}}(\mathbf{x})$ of \mathbf{x} over \mathcal{W} . If every factor of $F_{\mathcal{W}}(\mathbf{x})$ is maximal (**max**) in length over \mathcal{W} then the factorization is said to be max, and hence must be unique. So if \mathcal{W} is an FF on an alphabet Σ then \mathcal{W} is a **unique maximal factorization family** (UMFF) if there exists a max factorization $F_{\mathcal{W}}(\mathbf{x})$ for every string $\mathbf{x} \in \Sigma^+$ – for this theory see [2,3].

An UMFF \mathcal{W} is a **circ-UMFF** if it contains exactly one rotation of every primitive string $\mathbf{x} \in \Sigma^+$. The classic and foundational circ-UMFF is the set of Lyndon words, which we denote \mathcal{L} , where the rotation chosen is the one that is strictly least in the lexorder derived from an ordering of the letters of the alphabet Σ ([1,5,2]). The co-Lyndon circ-UMFF, $\text{co-}\mathcal{L}$, was introduced in [3], where a co-Lyndon word is strictly least amongst its rotations in co-lexorder.

Every circ-UMFF \mathcal{W} yields a strict order relation, the \mathcal{W} -**order**: if \mathcal{W} contains strings \mathbf{u}, \mathbf{v} and \mathbf{uv} then $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$. For the Lyndon circ-UMFF, its specific \mathcal{W} -order is lexorder:

Theorem 1. (Duval [5]) Suppose $\mathbf{u}, \mathbf{v} \in \mathcal{L}$, then $\mathbf{uv} \in \mathcal{L}$ if and only if $\mathbf{u} < \mathbf{v}$.

It was observed in [3] that the analogue of Theorem 1 does not hold for every circ-UMFF – we will establish this phenomenon for the co-Lyndon circ-UMFF.

Lemma 1. Suppose $\mathbf{u}, \mathbf{v} \in \text{co-}\mathcal{L}$, then $\mathbf{uv} \in \text{co-}\mathcal{L}$ if and only if $\mathbf{v} < \mathbf{u}$.

Proof. Since $\mathbf{u}, \mathbf{v} \in \text{co-}\mathcal{L}$ then $\bar{\mathbf{u}}, \bar{\mathbf{v}} \in \mathcal{L}$. If $\mathbf{v} < \mathbf{u}$ in co-lexorder then $\bar{\mathbf{v}} < \bar{\mathbf{u}}$ in lexorder. Applying Theorem 1 we have $\bar{\mathbf{v}}\bar{\mathbf{u}} \in \mathcal{L}$ and hence $\mathbf{uv} \in \text{co-}\mathcal{L}$. Next if $\mathbf{uv} \in \text{co-}\mathcal{L}$ then it must be primitive and border-free [2]. Thus $\mathbf{u} \neq \mathbf{v}$ which gives rise to two cases. Suppose first that $\mathbf{u} < \mathbf{v}$. If \mathbf{u} is a proper suffix of \mathbf{v} then $\mathbf{uv} = \mathbf{uwu}$ for some $\mathbf{w} \neq \varepsilon$ contradicting the border-free property. Otherwise, with $|\mathbf{u}| = n$ there is some largest j , $1 \leq j \leq n$, such that $\mathbf{u}[j] \neq \mathbf{v}[j]$. If $\mathbf{u}[j] < \mathbf{v}[j]$ then $\mathbf{vu} < \mathbf{uv}$ contradicting $\mathbf{uv} \in \text{co-}\mathcal{L}$. We conclude that $\mathbf{u}[j] > \mathbf{v}[j]$, and so $\mathbf{v} < \mathbf{u}$ as required. \square

The sets of Lyndon and co-Lyndon words are distinct and almost disjoint.

Lemma 2. For a given Σ , $\mathcal{L} \neq \text{co-}\mathcal{L}$ and $\mathcal{L} \cap \text{co-}\mathcal{L} = \Sigma$.

Proof. Let $\mathbf{v} \in \mathcal{L}$ and $\mathbf{w} \in \text{co-}\mathcal{L}$ with $|\mathbf{v}|, |\mathbf{w}| \geq 2$. Then \mathbf{v} starts with some letter α which is minimal in \mathbf{v} . Since \mathbf{v} is border-free [2] then it ends with some β where $\alpha < \beta$. Similarly, \mathbf{w} starts γ and ends δ , where $\gamma > \delta$. Therefore

$\mathbf{v} \neq \mathbf{w}$. Finally, every circ-UMFF contains the alphabet Σ as expressed in [2,3]. \square

The following result generalizes the Lyndon factorization theorem [1] and is a key to further applications of string decomposition.

Theorem 2. [2] Let \mathcal{W} be a circ-UMFF and suppose $\mathbf{x} = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_m$, with each $\mathbf{u}_j \in \mathcal{W}$. Then $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_m$ iff $\mathbf{u}_1 \geq_{\mathcal{W}} \mathbf{u}_2 \geq_{\mathcal{W}} \cdots \geq_{\mathcal{W}} \mathbf{u}_m$.

3. Alphabet ordering

Suppose the goal is to optimize a Lyndon factorization by minimizing or maximizing the number of factors. For this we consider choosing the order of the letters in the – assumed unordered – alphabet so as to influence the number of factors. To illustrate, consider the string $\mathbf{x} = \text{abcabcdabcaba}$. If $\Sigma = \{a < b < c < d\}$, then $F_{\mathcal{L}}(\mathbf{x}) = \text{abcabcd} \geq \text{abc} \geq \text{ab} \geq a$. Whereas, if we choose the alphabet ordering to be $\{b < c < a < d\}$, the Lyndon factorization of \mathbf{x} becomes $a \geq \text{bcabcdabcaba}$.

Towards this goal we now describe a greedy algorithm for producing small numbers of factors which has performed well in practice on the biological $\{A, C, G, T\}$ alphabet – the experimentation compared results with those for the $4!$ letter permutations. For a string $\mathbf{v} = v_1 \dots v_n$, we suppose that the number of distinct characters in \mathbf{v} is $\delta \leq \sigma$; for practical purposes we can assume σ is at most $O(n)$.

The proposed method requires an extension to a Parikh vector, $p(\mathbf{v})$, of a finite word \mathbf{v} , where $p(\mathbf{v})$ enumerates the occurrences of each letter of the alphabet in \mathbf{v} . Our modification is that for each distinct letter we will record its individual RLE (run length encoding) exponent pattern (left-to-right sequence of exponents for a letter) – so the sum of these exponents is the Parikh entry for that letter. We call this the **Exponent Parikh vector**, or EP vector, implemented as an EP array. For example, over the alphabet $\Sigma = \{b < c < d < f\}$, if $\mathbf{v} = \text{bbbffbbcf}$ then $p(\mathbf{v}) = [5, 1, 0, 3]$; whereas, for the EP vector we record the sequences $[(3, 2), (2, 1), (1)]$. So usually the letters are processed in the alphabet order with a Parikh vector while in the EP case we process them in order of first occurrence.

An overview of the method is that we use the fact that in a Lyndon factorization the first factor is the longest prefix which is a Lyndon word. Then the heuristic is that the left-most letter, α say, in the given string whose exponents, when read as a string, form a Lyndon word with the minimal number of factors is chosen as the least letter in the alphabet ordering. In order to construct a Lyndon word using the exponents of letters we require the order of the exponent integer alphabet to be inverted, that is, let $\bar{\Sigma} = \{\dots 3 < 2 < 1\}$. Next, the algorithm attempts to assign order to letters in the substrings between runs of α characters, where these substrings are denoted X_i – if it gets stuck it tries backtracking. Finally, if there is a nonempty prefix prior to the first α then it is processed.

So note that with this algorithm the required property for the exponents of α is that they form a Lyndon

word over $\bar{\Sigma}$ and in conjunction a requirement for assigning letters to the X_i substrings is that the ordering will be cycle-free. The algorithm can be modified to generate large numbers of factors which involves assigning distinct letters to be in decreasing order.

3.1. Greedy algorithm

The pseudocode in Algorithm 1 greedily assigns an alphabet order to letters.

Algorithm 1: Order the alphabet so as to reduce the number of factors in a Lyndon factorization.

```

With a linear scan record the Exponent Parikh (EP) vector of the
string for  $\delta$  distinct letters –  $O(n)$ ;
Compute  $F_{\mathcal{L}}(\mathbf{p}_r)$  of each exponent string  $\mathbf{p}_r$  over  $\bar{\Sigma}$  and record its
number of factors –  $O(n)$ ;
bool  $\leftarrow$  true; // initiate alphabet ordering
while bool = true do
  Select the next leftmost  $\mathbf{p}_r, \mathbf{p}_i$  say, with minimal number of
  factors,  $t$  say –  $O(n)$ ;
  // assign alphabet order to the  $t$  factors of
   $F_{\mathcal{L}}(\mathbf{p}_i) = \mathbf{f}_1 \geq \dots \geq \mathbf{f}_t$ 
  // where  $\mathbf{f}_j = \alpha^{j_1} X_1 \alpha^{j_2} X_2 \dots \alpha^{j_q} X_q$ , and  $\alpha \notin X_h$ ,
   $1 \leq h \leq q$ , with  $j_1 j_2 \dots j_q \in \mathcal{L}$  over  $\bar{\Sigma}$ 
   $\alpha = \lambda_1$ ; // assign first letter to be minimal
  in  $\Sigma$ ; if  $q=1$  assign each new letter in  $X_1$ 
  successively in  $\Sigma$ 
  for  $h = 2$  to  $q$  do
    if  $j_h = j_1$  then // same exponents so assign
    alphabet in order to letters in  $X_1$  and  $X_h$ 
    substrings
       $d \leftarrow 1$ ;
      while  $X_1[d] = X_h[d]$  do
        assign each new letter successively in  $\Sigma$ ;  $d++$ ;
      if  $X_1[d] = \alpha$  &  $X_h[d] \neq \alpha$  then
        assign  $X_h[d]$  to be next successive letter;
      else if  $X_h[d] = \alpha$  &  $X_1[d] \neq \alpha$  then
        bool  $\leftarrow$  false; // not Lyndon
      else if assignment would not make inconsistency then
        //  $X_1[d] \neq X_h[d]$ 
        assign  $X_h[d] > X_1[d]$ ;
      else
        bool  $\leftarrow$  false; // inconsistent
    if bool then
      attempt assignment process for the  $t$  factors of  $F_{\mathcal{L}}(\mathbf{p}_i)$ ;
    if bool then
      if string prefix  $\mathbf{u}$  (prior to  $\mathbf{f}_1$ ) is non-empty then //  $\alpha \notin \mathbf{u}$ 
        repeat process on  $\mathbf{u}$  starting with next successive letters
        in  $\Sigma$ ;
        // lookup EP vector
        complete assignment of any remaining letters;
        if letters in prefix  $\mathbf{u}$  do not occur in suffix then re-assign
        all letters starting from prefix;
      else
        arbitrarily choose next leftmost  $\mathbf{p}_r$  with minimal number of
        factors and attempt new assignment;

```

The following example, which uses the notation of Algorithm 1, illustrates how backtracking can lead the algorithm from an inconsistent ordering to a successful assignment and associated factorization.

Example 1. Assume $\Sigma = \{a, b, c, d\}$ and

$$\mathbf{x} = a^2 b d c a^2 c d a^2 b d b a^1 b a^2 b c a^2 c a^1 b.$$

Only the letter a has an exponent greater than 1, and $F_{\mathcal{L}}(EP(a)) = 2221 \geq 2221$ has the minimal number of factors with $\mathbf{f}_1 = \mathbf{f}_2 = 2221$. After assigning a to be the first letter in Σ , processing \mathbf{f}_1 causes inconsistency (and similarly \mathbf{f}_2), since the substrings X_1 and X_2 give $b < c$ while X_1 and X_3 would require $c < b$. So the algorithm then backtracks through the EP array and chooses the letter with the least number of factors (albeit singletons), namely d – the result is $\Sigma = \{d < c < a < b\}$ with $F_{\mathcal{L}}(\mathbf{x}) = aab \geq dcaacdaabdbabaabcaacaacab$. Note the order $\{a < b < c < d\}$ would have given 3 factors.

4. Experimentation: factorization of DNA strings

We chose as an example the 120 prokaryotic reference genomes from RefSeq,¹ to investigate the results of the algorithm in practice.² Most of these genomes are provided as a single contiguous sequence but some of them have additional smaller pieces representing plasmids or other information. The longest contiguous sequence was chosen for each genome in these cases, and smaller pieces were discarded. The retained sequences ranged from 640,681 letters to 10,236,715 in length, with a mean of 3,629,792.

In order to determine how often our greedy algorithm found a good or optimal alphabet reordering in practice, we calculated the Lyndon factorizations resulting from all possible ($4! = 24$) alphabet reorderings of the characters A, C, G and T across this collection of genomes. The improvement that could potentially be made to the factorization by reordering is substantial, with at least a halving of the number of factors in most cases and an improvement reducing 25 factors down to 3 in one case. For each genome we ranked the results of all possible reorderings by the number of factors produced and compared the reordering produced by the algorithm. The algorithm found the optimal reordering for 21/120 genomes and the second-most optimal in 31/120 genomes.

The EP vector is used to determine the least letter in the reordering. If the first choice leads to inconsistency (and hence small factors), backtracking to inspect other possible choices can be helpful. However, in many cases, the initial choice is still better than the next possible consistent solution found via backtracking through the EP vector. Without backtracking, the algorithm found 23/120 optimal orderings and a further second-most optimal orderings in 31/120 genomes. Histograms of the full results, with and without backtracking, can be seen in the Supplemental Information.

5. Applications

In many cases, such as natural language text processing, the order of the alphabet is prescribed, and hence

¹ RefSeq: <https://www.ncbi.nlm.nih.gov/refseq/about/prokaryotes>.

² Code available at <https://github.com/amandaclare/lyndon-factors>.

the Lyndon factors of an input text cannot be manipulated. On the other hand, bioinformatics alphabets have no inherent ordering suggested by biological systems and applications involving Lyndon words, such as the Burrows–Wheeler transform (BWT), will allow for useful manipulation of the Lyndon factors. The co-BWT is the regular BWT of the reversed string, or the BWT with co-lexorder, which has been applied in the highly successful Bowtie sequence alignment program [6]. Integral with the BWT transform is the computation of suffix arrays via induced suffix-sorting. We also propose that pattern matching can be implemented with the Lyndon factorization in big data applications, such as sequence alignment, and further enhanced by fortuitous arrangements of the alphabet.

In [7] a new method is presented for constructing the suffix array of a text by using its Lyndon factorization advantageously. Partitioning the text according to its Lyndon properties allows tackling the problem in local portions of the text, local suffixes, prior to extending the solution globally, to achieve the suffix array of the entire text. It is stated that the algorithm’s time complexity is not competitive for the construction of the overall suffix array – we propose that manipulating the factors by alphabet ordering could improve the efficiency.

6. Research problems

- As a complementary structure to the Lyndon array we introduce and propose studies of the Lyndon factorization array. The **Lyndon array** $\lambda = \lambda_x[1..n]$ of a given $\mathbf{x} = \mathbf{x}[1..n]$ gives at each i the length of the maximal Lyndon word starting at i – reverse engineering in [4] includes a linear-time test for whether an integer array is a Lyndon array. We define the **Lyndon factorization array** $\mathbf{F} = \mathbf{F}_x[1..n]$ of \mathbf{x} to give at each position i the number of factors in the Lyndon factorization starting at i .
- The greedy algorithm does not necessarily produce an optimal solution hence natural problems are to design algorithms for Lyndon factorizations with a guaranteed minimal/maximal number of Lyndon factors.
- Using Duval’s Lyndon factorization algorithm [5] as a benchmark, modify the alphabet order so as to increase/decrease the number of factors.
- Theorem 2 supports the following optimization problem from [3]: Determine the circ-UMFF(s) which factors a string \mathbf{x} into the minimal/maximal number of factors, possibly combined with alphabet ordering.

Appendix A. Supplementary material

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.ipl.2018.10.011>.

References

- [1] K.T. Chen, R.H. Fox, R.C. Lyndon, Free differential calculus, IV. The quotient groups of the lower central series, *Ann. Math.* 68 (1) (1958) 81–95.
- [2] D.E. Daykin, J.W. Daykin, Properties and construction of unique maximal factorization families for strings, *Int. J. Found. Comput. Sci.* 19 (4) (2008) 1073–1084.
- [3] D.E. Daykin, J.W. Daykin, W.F. Smyth, Combinatorics of unique maximal factorization families (UMFFs), in: *Special Issue on Stringology*, *Fundam. Inform.* 97 (3) (2009) 295–309.
- [4] J.W. Daykin, F. Franek, J. Holub, A.S.M.S. Islam, W.F. Smyth, Reconstructing a string from its Lyndon arrays, *Theor. Comput. Sci.* 710 (2018) 44–51.
- [5] J.-P. Duval, Factorizing words over an ordered alphabet, *J. Algorithms* 4 (4) (1983) 363–381.
- [6] B. Langmead, C. Trapnell, M. Pop, S.L. Salzberg, Ultrafast and memory-efficient alignment of short DNA sequences to the human genome, *Genome Biol.* 10 (3) (2009) R25.
- [7] S. Mantaci, A. Restivo, G. Rosone, M. Sciortino, Suffix array and Lyndon factorization of a text, *J. Discret. Algorithms* 28 (2014) 2–8.