

Aberystwyth University

A linear partitioning algorithm for Hybrid Lyndons using V-order

Daykin, David E.; Daykin, Jacqueline W.; Smyth, William F.

Published in:

Theoretical Computer Science

DOI:

[10.1016/j.tcs.2012.02.001](https://doi.org/10.1016/j.tcs.2012.02.001)

Publication date:

2013

Citation for published version (APA):

Daykin, D. E., Daykin, J. W., & Smyth, W. F. (2013). A linear partitioning algorithm for Hybrid Lyndons using V-order. *Theoretical Computer Science*, 483, 149-161. <https://doi.org/10.1016/j.tcs.2012.02.001>

General rights

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400
email: is@aber.ac.uk



A linear partitioning algorithm for Hybrid Lyndons using V -order[☆]

David E. Daykin^a, Jacqueline W. Daykin^{b,c,*}, W.F. Smyth^{c,d}

^a Department of Mathematics, University of Reading, UK

^b Department of Computer Science, Royal Holloway, University of London, UK

^c Department of Informatics, King's College, University of London, UK

^d Algorithms Research Group, Department of Computing & Software, McMaster University, Hamilton ON L8S 4K1, Canada

ARTICLE INFO

Keywords:

Algorithm
Alphabet
Circ-UMFF
Concatenate
Factor
Hybrid Lyndon
Lexicographic order
Lyndon
Maximal
RAM
String
Total order
UMFF
 V -order
Word
 V -word

ABSTRACT

In this paper we extend previous work on unique maximal factorization families (UMFFs) and a total (but non-lexicographic) ordering of strings called V -order. We present new combinatorial results for V -order, in particular concatenation under V -order. We propose linear-time RAM algorithms for string comparison in V -order and for Lyndon-like factorization of a string into V -words. This asymptotic efficiency thus matches that of the corresponding algorithms for lexicographical order. Finally, we introduce Hybrid Lyndon words as a generalization of standard Lyndon words, and hence propose extensions of factorization algorithms to other forms of order.

© 2012 Elsevier B.V. All rights reserved.

1. Introduction

This paper explores a generalization of Lyndon words known as circ-UMFFs; that is, families of strings which permit the unique maximal factorization of any given string [9,10,12]. For over half a century, both the combinatorics and algorithmics of Lyndon words have been studied extensively [4,16,19,22]. Lyndon words have been applied in tackling a surprisingly wide range of problems: the Burrows–Wheeler Transform and data compression [3,17], musicology [1], bioinformatics [15], cryptanalysis [20], string combinatorics [12,18,22], and free Lie algebras [21]. We therefore study circ-UMFFs as generalizations of Lyndon words in order to extend the range of applications.

A **circ-UMFF** \mathcal{W} (see below, Definitions 2.2 and 2.5) is an infinite collection of strings on some alphabet Σ , such that every string in Σ^* can be expressed as a unique concatenation of maximum-length elements of \mathcal{W} – what we call a **max factorization** – with the additional requirement that no two elements of \mathcal{W} can be rotations of each other. It turns out [9,10] that for a circ-UMFF \mathcal{W} , the strings w_i of \mathcal{W} satisfy an **UMFF-order** $<_{\mathcal{W}}$, a total order over all the elements of \mathcal{W} such that every max factorization $w_1 w_2 \cdots w_k$ of a string x over \mathcal{W} satisfies

$$w_1 \geq_{\mathcal{W}} w_2 \geq_{\mathcal{W}} \cdots \geq_{\mathcal{W}} w_k. \quad (1)$$

[☆] A preliminary version of these results appeared in Daykin et al. (2011) [13]. The authors acknowledge the contribution of two anonymous referees whose suggestions and ideas have materially improved this paper. The work of the third author was supported in part by the Natural Sciences & Engineering Research Council of Canada.

* Corresponding author at: Department of Informatics, King's College, University of London, UK.

E-mail addresses: J.Daykin@cs.rhul.ac.uk, jackie.daykin@kcl.ac.uk (J.W. Daykin), smyth@mcmaster.ca (W.F. Smyth).

Thus, by virtue of its defining properties, a circ-UMFF satisfies, in terms of UMFF-order, an analogue or generalization of the Lyndon factorization theorem [4], which guarantees that the Lyndon words \mathcal{L} , in terms of lexicographic order (**lexorder**), provide a unique maximal factorization (1) of every string \mathbf{x} . In other words, \mathcal{L} is an *example* of a circ-UMFF \mathcal{W} , suggesting that other orderings of the strings $\mathbf{x} \in \Sigma^+$ may yield alternative circ-UMFFs distinct from \mathcal{L} . Various alternative orderings were subsequently introduced in [9]:

- binary circ-UMFFs including *B*-words from the Block total order and *T*-words from the Type total order;
- the *V*-word circ-UMFF derived from *V*-order on an arbitrary alphabet [6,8].

The latter is our focus here.

In this paper we first extend combinatorial insights into *V*-order and the associated circ-UMFF factors, *V*-words. For instance, we show that, unlike lexorder, any subsequence of a string precedes the string in *V*-order. Surprisingly, although *V*-order is not immediately intuitive, it gives rise to *V*-words which exhibit similar structure to Lyndon words. One similarity is that when factoring into Lyndon words, a lesser letter than the first denotes the start of a new factor, while with *V*-words this occurs with a larger letter. So, interestingly, *V*-words are analogues of Lyndon words, defined using a simple extension of lexorder, **lex-extension**, a combination of lexorder and *V*-order. These various insights lead to efficient new RAM algorithms for comparing strings in *V*-order and for factoring them into *V*-words. See [2] for a representation of rhythms in Central African music in terms of such a factorization.

More generally, when comparing strings lexicographically, any secondary ordering method can be applied to prescribed substrings. Based on the general form of lex-extension ordering, we define the concept of **Hybrid Lyndons** and show that *V*-words are an instance of this new class of Lyndon words. Then, since *V*-words are defined using lex-extension, an associated factorization algorithm naturally arises from the modification of existing Lyndon partitioning algorithms [16].

The Lyndon factorization for both RAM and PRAM algorithms has linear complexity [7,11,14,16,22], achievable with or without the use of recursion. We compare substrings \mathbf{u} , \mathbf{v} nonrecursively in *V*-order, also achieving linear time by means of a method based on identification of the longest common suffix of the substrings \mathbf{u} , \mathbf{v} and implemented using a simple doubly-linked list.

The currently identified (and, as far as we know, non-exhaustive) classes of circ-UMFFs are type Flight Deck and type Acrobat [12]. Since both Lyndon words and *V*-words belong to the Flight Deck class, which requires that all letters in a factor are not less than the first in UMFF-order, this paper also advances knowledge of this class of circ-UMFFs. We compare and contrast results for concatenation of Lyndon versus *V*-word factors: concatenation of Lyndon words is always in the same order as the defining lexorder of those words [16], whereas *V*-words demonstrate a variety of orders; indeed these two types of factors even exhibit different UMFF-behavior for a pair of strings according to the ordering of the first letters in each string. Similarly to the co-Lyndon circ-UMFF [12], the *V*-word alphabet is in reverse order to that of the defining *V*-order.

We show that for a given string \mathbf{x} , the lexicographic factorization of \mathbf{x} into Lyndon words is distinct from the *V*-order factorization of \mathbf{x} into *V*-words. Hence we have different methods for string factorization. Once a string has been decomposed over a circ-UMFF, say the Lyndon words, the factors are then maximal and thus final. Given an alternative circ-UMFF, say *V*-words, the Lyndon factors can be re-factored allowing for deeper burrowing into the string, for example yielding a subsequence of *V*-word prefixes, suffixes or substrings of all the Lyndon factors. Thus, algorithms for additional circ-UMFFs may lead to computing further subsequences of interest.

Additionally, the existence of a variety of circ-UMFFs opens the way to new optimization problems, for instance minimizing/maximizing the number of factors in a string partition. Indeed, given linear factorization methods, it costs linear time overall to determine an optimal factorization.

The structure of this paper is as follows. We first review the circ-UMFF theory relevant to this work in Section 2, followed by a discussion of *V*-order and associated new combinatorial results in Section 3. An algorithm for fast *V*-order string comparison is described in Section 4, then applied in Section 5 to achieve linear-time Lyndon-like factorization using *V*-order.

Note that the terms **string** and **word** mean the same thing (see References), hence we use both here. Also, for any examples included which use either the Roman alphabet or the non-negative integers, the natural order will be assumed, that is $\{a < b < c < \dots < z\}$ and $\{0 < 1 < 2 < \dots\}$. All strings are written in **mathbold**: \mathbf{x} , \mathbf{w} , and so on.

2. Unique maximal factorization families (UMFFs)

We commence with a summary of relevant UMFF theory from the fuller version given in [12]. Note that we use standard terminology from automata theory and stringology (see definitions in [22]). For brevity we write **lexorder** for lexicographic order, and **co-lexorder** for complementary lexorder; then, if a string \mathbf{x} is less than a string \mathbf{y} in lexorder, we write $\mathbf{x} < \mathbf{y}$ and $\mathbf{y} > \mathbf{x}$.

A string \mathbf{w} is called a **factor** of a string $\mathbf{x}[1..n]$ if and only if $\mathbf{w} = \mathbf{x}[i..j]$ for $1 \leq i \leq j \leq n$. Note that a factor is necessarily nonempty. If $\mathbf{x} = \mathbf{w}_1\mathbf{w}_2 \dots \mathbf{w}_k$, $1 \leq k \leq n$, then $\mathbf{w}_1\mathbf{w}_2 \dots \mathbf{w}_k$ is said to be a **factorization** of \mathbf{x} ; moreover, when every factor \mathbf{w}_j , $1 \leq j \leq k$, belongs to a specified set \mathcal{W} , this is a **factorization of \mathbf{x} over \mathcal{W}** , denoted by $F_{\mathcal{W}}(\mathbf{x})$.

Definition 2.1. 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 $F_{\mathcal{W}}(\mathbf{x})$.

For some string \mathbf{x} and some FF \mathcal{W} , suppose $\mathbf{x} = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_k$, where $\mathbf{w}_j \in \mathcal{W}$ for every $1 \leq j \leq k$. For some $1 \leq k' \leq k$, write $\mathbf{x} = \mathbf{u}\mathbf{w}_{k'}\mathbf{v}$, where $\mathbf{u} = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_{k'-1}$ (empty if $k' = 1$) and $\mathbf{v} = \mathbf{w}_{k'+1}\mathbf{w}_{k'+2} \cdots \mathbf{w}_k$ (empty if $k' = k$). Suppose that there does not exist a suffix \mathbf{u}' of \mathbf{u} nor a prefix \mathbf{v}' of \mathbf{v} such that $\mathbf{u}'\mathbf{w}_{k'}\mathbf{v}' \neq \mathbf{w}_{k'}$ and $\mathbf{u}'\mathbf{w}_{k'}\mathbf{v}' \in \mathcal{W}$; then $\mathbf{w}_{k'}$ is said to be a **max factor** of \mathbf{x} . If every factor $\mathbf{w}_{k'}$ is max, then the factorization $F_{\mathcal{W}}(\mathbf{x})$ is itself said to be **max**. Observe that a max factorization must be unique: there exists no other max factorization of \mathbf{x} that uses only elements of \mathcal{W} .

Definition 2.2. Let \mathcal{W} be an FF on an alphabet Σ . Then \mathcal{W} is a **unique maximal factorization family** (UMFF) if and only if there exists a max factorization $F_{\mathcal{W}}(\mathbf{x})$ for every string $\mathbf{x} \in \Sigma^+$.

Lemma 2.3 (The **xyz Lemma** [9]). An FF \mathcal{W} is an UMFF if and only if whenever $\mathbf{xy}, \mathbf{yz} \in \mathcal{W}$ for some nonempty \mathbf{y} , then $\mathbf{xyz} \in \mathcal{W}$.

The next result establishes that a max factorization is indeed a partition of a string with no overlapping factors. It is also relevant to the parallel PRAM Lyndon algorithm in [14] where recursion is applied to refactor previously factored substrings.

Corollary 2.4 ([12]). Suppose $\mathbf{x} = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_m$ and \mathcal{W} is an UMFF, where for every $1 \leq j \leq m$, $\mathbf{u}_j \in \mathcal{W}$. Then we can write the max factorization $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_k$, where

$$\mathbf{w}_1 = \mathbf{u}_{j_0+1} \cdots \mathbf{u}_{j_1}, \quad \mathbf{w}_2 = \mathbf{u}_{j_1+1} \cdots \mathbf{u}_{j_2}, \quad \dots, \quad \mathbf{w}_k = \mathbf{u}_{j_{k-1}+1} \cdots \mathbf{u}_{j_k},$$

$$0 = j_0 < j_1 < j_2 < \cdots < j_{k-1} < j_k = m.$$

Our algorithms execute on circ-UMFFs, with an associated rule for concatenation. To define these ideas, we first give a few basic stringological definitions [22]:

- If for some string \mathbf{x} we can write $\mathbf{x} = \mathbf{uv} = \mathbf{wu}$ for some nonempty \mathbf{u} , then we say that \mathbf{x} has **border \mathbf{u}** ; if no such \mathbf{u} exists, then \mathbf{x} is said to be **border-free**.
- If we can write $\mathbf{x} = \mathbf{u}^k$ for some nonempty \mathbf{u} and some integer $k > 1$, we say that \mathbf{x} is a **repetition**; otherwise, we say that \mathbf{x} is **primitive**.
- If a string $\mathbf{x} = \mathbf{uv}$, then \mathbf{vu} is said to be a **rotation** (cyclic shift) of \mathbf{x} , specifically the $|\mathbf{u}|^{\text{th}}$ rotation $R_{|\mathbf{u}|}(\mathbf{x})$, where $|\mathbf{u}| \in \{0, \dots, |\mathbf{x}|\}$. Note that $R_0(\mathbf{x}) = R_{|\mathbf{x}|}(\mathbf{x})$.

Definition 2.5. An UMFF \mathcal{W} over Σ^+ is a **circ-UMFF** if and only if it contains exactly one rotation of every primitive string $\mathbf{x} \in \Sigma^+$.

The set of strings consisting of increasing integers is an example of an UMFF that is not a circ-UMFF (see [10] for further examples).

Definition 2.6. If a circ-UMFF \mathcal{W} contains strings \mathbf{u}, \mathbf{v} and \mathbf{uv} , we write $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$ (called the **\mathcal{W} -order**).

The next theorem provides characterizations of circ-UMFFs that will be frequently used in the following sections:

Theorem 2.7 ([10]). Let \mathcal{W} be a circ-UMFF.

- (1) If $\mathbf{u} \in \mathcal{W}$ then \mathbf{u} is border-free.
- (2) If $\mathbf{u}, \mathbf{v} \in \mathcal{W}$ and $\mathbf{u} \neq \mathbf{v}$ then \mathbf{uv} is primitive.
- (3) If $\mathbf{u}, \mathbf{v} \in \mathcal{W}$ and $\mathbf{u} \neq \mathbf{v}$ then $\mathbf{uv} \in \mathcal{W}$ or $\mathbf{vu} \in \mathcal{W}$ (but not both).
- (4) If $\mathbf{u}, \mathbf{v}, \mathbf{uv} \in \mathcal{W}$ then $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$ and $<_{\mathcal{W}}$ is a total order of \mathcal{W} .
- (5) If $\mathbf{w} \in \mathcal{W}$ and $|\mathbf{w}| \geq 2$ then there exist $\mathbf{u}, \mathbf{v} \in \mathcal{W}$ with $\mathbf{w} = \mathbf{uv}$.

An example of \mathcal{W} -order is lexorder which gives the Lyndon circ-UMFF, denoted here as \mathcal{L} . Recall [16] that for Lyndon words, the \mathcal{W} -order of the factors is just lexorder, the same as the order normally applied to the strings in Σ^* .

Suppose now that the elements of Σ^* are totally ordered under $<_{\mathcal{T}}$, say, and that \mathcal{W} is a circ-UMFF. It may then happen that for distinct strings \mathbf{u} and \mathbf{v} , $\mathbf{u} <_{\mathcal{T}} \mathbf{v}$ while $\mathbf{v} <_{\mathcal{W}} \mathbf{u}$ (see [12]). Thus if $<_{\mathcal{T}}$ is used to select a rotation as in Definition 2.5 (for example, lexorder for Lyndon words), then we may need to distinguish whether the generated \mathcal{W} -order is coincidental or otherwise to the order $<_{\mathcal{T}}$.

Definition 2.8. Suppose that the strings of Σ^* are ordered according to a total order $<_{\mathcal{T}}$, and that \mathcal{W} is a circ-UMFF on Σ^* . If, for every $\mathbf{u}, \mathbf{v} \in \mathcal{W}$ such that $\mathbf{u} <_{\mathcal{W}} \mathbf{v}$,

- (1) it follows that $\mathbf{u} <_{\mathcal{T}} \mathbf{v}$, then we write $\mathcal{W} \equiv \mathcal{T}$;
- (2) it follows that $\mathbf{v} <_{\mathcal{T}} \mathbf{u}$, then we write $\mathcal{W} \equiv \overline{\mathcal{T}}$.

A simple example of $\mathcal{W} \equiv \overline{\mathcal{T}}$ arises if we define $<_{\mathcal{T}}$ as lexorder while \mathcal{W} is defined as **co-lexorder** (that is, lexorder of reversed strings [12]). Note that defining \mathcal{W} as lexorder on some other ordering of Σ would yield $\mathcal{W} \neq \mathcal{T}$ and $\mathcal{W} \neq \overline{\mathcal{T}}$. See below, lexicographic extension, Definition 3.9 and Lemma 3.16.

The following result generalizes the Lyndon factorization theorem [4] to circ-UMFFs, and is relevant to the algorithm of Section 5 (cf. Corollary 2.4).

Theorem 2.9 ([12]). *Let \mathcal{W} be a circ-UMFF and let $\mathbf{x} = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_m$ with each $\mathbf{u}_j \in \mathcal{W}$. Then $F_{\mathcal{W}}(\mathbf{x}) = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_m$ if and only if*

$$\mathbf{u}_1 \geq_{\mathcal{W}} \mathbf{u}_2 \geq_{\mathcal{W}} \cdots \geq_{\mathcal{W}} \mathbf{u}_m.$$

The factors in a circ-UMFF exhibit nested structural properties:

Lemma 2.10 ([10]). *Suppose that \mathbf{w} is an element of a circ-UMFF \mathcal{W} . If $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k_1}$ are all the proper prefixes of \mathbf{w} in increasing order of length that belong to \mathcal{W} , and if $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k_2}$ are all the proper suffixes of \mathbf{w} in decreasing order of length that belong to \mathcal{W} , then*

$$\mathbf{u}_1 <_{\mathcal{W}} \mathbf{u}_2 <_{\mathcal{W}} \cdots <_{\mathcal{W}} \mathbf{u}_{k_1} <_{\mathcal{W}} \mathbf{w} <_{\mathcal{W}} \mathbf{v}_1 <_{\mathcal{W}} \mathbf{v}_2 <_{\mathcal{W}} \cdots <_{\mathcal{W}} \mathbf{v}_{k_2}.$$

On the other hand, consider the variation in ordering that may occur in circ-UMFFs:

Lemma 2.11 ([12]). *Let \mathcal{W} be a circ-UMFF and suppose $\mathbf{xy}, \mathbf{yz} \in \mathcal{W}$ for nonempty $\mathbf{x}, \mathbf{y}, \mathbf{z}$ (hence $\mathbf{x} \neq \mathbf{z}$). Then $\mathbf{xyz} \in \mathcal{W}$, $\mathbf{xyyz} \in \mathcal{W}$, and*

- (1) $\mathbf{xy} <_{\mathcal{W}} \mathbf{xyz} <_{\mathcal{W}} \mathbf{yz}$;
- (2) $\mathbf{xy} <_{\mathcal{W}} \mathbf{xyyz} <_{\mathcal{W}} \mathbf{yz}$;
- (3) either $\mathbf{xyyzxyz} \in \mathcal{W}$ or $\mathbf{xyzxyyz} \in \mathcal{W}$ (but not both).

The case $\mathbf{xyyz} <_{\mathcal{W}} \mathbf{xyz}$ of Lemma 2.11(3) occurs for the Lyndon circ-UMFF based on lexicographic ordering [12]. As we shall see at the end of Section 3, the other choice in Lemma 2.11(3), $\mathbf{xyz} <_{\mathcal{W}} \mathbf{xyyz}$, arises in the V -order circ-UMFF, in which moreover both choices actually occur.

In [12] we introduced Flight Deck and Acrobat circ-UMFFs, and in this paper we are working with the former:

Definition 2.12. A circ-UMFF \mathcal{W} is said to be **type Flight Deck** if and only if $\mathbf{w}[1 \dots n] \in \mathcal{W}$ implies that for every $1 \leq i \leq n$, $\mathbf{w}[1] \leq_{\mathcal{W}} \mathbf{w}[i]$.

3. Combinatorics of V -order

We begin with the definition for totally ordering a set of strings according to V -order, followed first by some preliminary lemmas and then by the definition of V -order factors called V -words. Hence in this section \mathcal{W} -order is denoted as \mathcal{V} -order. (See [8,9] for related results and discussion on V -order and V -words; see [6] for an order equivalent to V -order due to Strehl and Winkelmann.)

Let Σ be a totally ordered alphabet, and let $\mathbf{u} = u_1u_2 \dots u_n$ be a string over Σ . Define $h \in \{1, \dots, n\}$ by $h = 1$ if $u_1 \leq u_2 \leq \dots \leq u_n$; otherwise, by the unique value such that $u_{h-1} > u_h \leq u_{h+1} \leq u_{h+2} \leq \dots \leq u_n$. Let $\mathbf{u}^* = u_1u_2 \dots u_{h-1}u_{h+1} \dots u_n$, where the star $*$ indicates deletion of the “ V ” letter u_h . Write \mathbf{u}^{s*} for $(\dots(\mathbf{u}^*)^* \dots)^*$ with $s \geq 0$ stars.¹ Let $\mu\mathbf{u} = \max\{u_1, u_2, \dots, u_n\}$, say $\mu\mathbf{u} = g$, and let $k = \nu\mathbf{u}$ be the number of occurrences of g in \mathbf{u} . Then the sequence $\mathbf{u}, \mathbf{u}^*, \mathbf{u}^{2*}, \dots$ ends $g^k, \dots, g^2, g^1, g^0 = \epsilon$. In the *star tree* each string \mathbf{u} over Σ labels a vertex, and there is a directed edge from \mathbf{u} to \mathbf{u}^* , with ϵ as the root.

Definition 3.1. We define V -order $<$ between distinct strings \mathbf{u}, \mathbf{v} . First $\mathbf{v} < \mathbf{u}$ if \mathbf{v} is in the path $\mathbf{u}, \mathbf{u}^*, \mathbf{u}^{2*}, \dots, \epsilon$. If \mathbf{u}, \mathbf{v} are not in a path, there exist smallest s, t such that $\mathbf{u}^{(s+1)*} = \mathbf{v}^{(t+1)*}$. Put $\mathbf{c} = \mathbf{u}^{s*}$ and $\mathbf{d} = \mathbf{v}^{t*}$; then $\mathbf{c} \neq \mathbf{d}$ but $|\mathbf{c}| = |\mathbf{d}| = m$ say. Let j be the greatest i in $1 \leq i \leq m$ such that $\mathbf{c}[i] \neq \mathbf{d}[i]$. If $\mathbf{c}[j] < \mathbf{d}[j]$ in Σ then $\mathbf{u} < \mathbf{v}$. Clearly $<$ is a total order.

We will be ordering the rotations of a string in V -order $<$, illustrated by:

Example 3.2. In V -order, $a < b < ab < bbb < c < abc < abbc < cbc < bcc$. For the rotations of $abbc$, we have $cabb < bcab < bbca < abbc$.

In order to describe our combinatorial results we introduce the **V -form** of a string \mathbf{x} , a unique representation of \mathbf{x} defined as

$$V_k(\mathbf{x}) = \mathbf{x} = \mathbf{x}_0g\mathbf{x}_1g \cdots \mathbf{x}_{k-1}g\mathbf{x}_k, \tag{2}$$

for possibly empty \mathbf{x}_i , $i = 0, 1, \dots, k$, where g is the largest letter in \mathbf{x} . (Thus we suppose that g occurs exactly k times.) Observe that if g_1 is the second largest letter in \mathbf{x} , then V -form can be recursively (and independently) applied to each \mathbf{x}_i that includes g_1 with frequency $k_1 > 0$:

$$V_{k_1}(\mathbf{x}_i) = \mathbf{x}_{i,0}g_1\mathbf{x}_{i,1} \cdots \mathbf{x}_{i,k_1-1}g_1\mathbf{x}_{i,k_1}.$$

¹ Note that this star operator, as defined in [8,9,13], is distinct from the Kleene star operator.

We now clarify how the sequence $\mathbf{u}, \mathbf{u}^*, \mathbf{u}^{2*}, \dots$ ends $g^k, \dots, g^2, g^1, g^0 = \boldsymbol{\varepsilon}$ in the star tree. Denote by $A^T(\mathbf{x})$ the set of ancestors of \mathbf{x} in the star tree, and by $A_j(\mathbf{x})$ the j th ancestor, $0 \leq j \leq |\mathbf{x}|$ (the zeroth ancestor is \mathbf{x} itself, the $|\mathbf{x}|$ th ancestor is the empty string $\boldsymbol{\varepsilon}$).

Lemma 3.3. *If $V_k(\mathbf{x}) = \mathbf{x}_0 g \mathbf{x}_1 g \dots \mathbf{x}_{k-1} g \mathbf{x}_k$, then*

$$A_{t_j}(\mathbf{x}) = \mathbf{x}_0^T g \mathbf{x}_1^T \dots \mathbf{x}_{k-1}^T g \mathbf{x}_k^T,$$

where $j \in 0..k$, $t_j = |\mathbf{x}_j| + |\mathbf{x}_{j+1}| + \dots + |\mathbf{x}_k|$, and

$$\begin{aligned} \mathbf{x}_h^T &= \mathbf{x}_h \quad \text{for } 0 \leq h < j, \\ &= \boldsymbol{\varepsilon} \quad \text{for } j \leq h \leq k. \end{aligned}$$

Thus $A_{t_0}(\mathbf{x}) = g^k$, while for $j \geq 1$, $A_{t_j}(\mathbf{x})$ has suffix g^{k-j+1} and $A_{t_{j-1}}(\mathbf{x}) = A_{|\mathbf{x}_{j-1}|}(A_{t_j}(\mathbf{x}))$.

Our interest is to uniquely factor a string by applying V -order rather than lexorder. The following lemma forms a basis for the design of the factorization algorithm in Section 5.

Lemma 3.4 ([9]). *Given distinct strings \mathbf{v} and \mathbf{x} with corresponding V -forms*

$$\mathbf{v} = \mathbf{v}_0 g \mathbf{v}_1 g \mathbf{v}_2 \dots \mathbf{v}_{j-1} g \mathbf{v}_j, \quad \mathbf{x} = \mathbf{x}_0 g \mathbf{x}_1 g \mathbf{x}_2 \dots \mathbf{x}_{k-1} g \mathbf{x}_k,$$

let $h \in 0.. \max(j, k)$ denote the least integer such that $\mathbf{v}_h \neq \mathbf{x}_h$. Then $\mathbf{v} < \mathbf{x}$ (respectively, $\mathbf{v} > \mathbf{x}$) if and only if one of the following conditions holds:

- (1) $\mu \mathbf{v} < \mu \mathbf{x}$ (respectively, $\mu \mathbf{v} > \mu \mathbf{x}$);
- (2) $\mu \mathbf{v} = \mu \mathbf{x}$ and $j < k$, that is $\mathbf{v} \mathbf{v} < \mathbf{v} \mathbf{x}$ (respectively, $j > k$);
- (3) $\mu \mathbf{v} = \mu \mathbf{x}$ and $j = k$ and $\mathbf{v}_h < \mathbf{x}_h$ (respectively, $\mathbf{v}_h > \mathbf{x}_h$).

As remarked by a reviewer, this lemma leads directly to a simple recursive $O(\sigma n)$ -time algorithm to compare two strings \mathbf{v} and \mathbf{x} , where $\sigma = |\Sigma|$ and $n = |\mathbf{v}| + |\mathbf{x}|$:

1. Scan \mathbf{v} and \mathbf{x} separately to determine their largest letters $g_{\mathbf{v}}$ and $g_{\mathbf{x}}$, respectively, as well as corresponding frequencies $j = \nu \mathbf{v}$ and $k = \nu \mathbf{x}$.
2. If $g_{\mathbf{v}} \neq g_{\mathbf{x}}$ or $j \neq k$, then return “ $<$ ” or “ $>$ ” according to Lemma 3.4(1) & (2), and exit.
3. Scan \mathbf{v} and \mathbf{x} together to determine the least h such that $\mathbf{v}_h \neq \mathbf{x}_h$.
4. Set $\mathbf{v} \leftarrow \mathbf{v}_h$, $\mathbf{x} \leftarrow \mathbf{x}_h$ and goto (1).

The worst case $O(n^2)$ for this algorithm would be achieved by strings such as $\mathbf{v} = n n - 1 \dots 3 2$ and $\mathbf{x} = n n - 1 \dots 3 1$ with $\sigma = n$.

The next lemma is an extension of Lemma 3.4(3). It provides the interesting insight that if \mathbf{v} is any proper subsequence of \mathbf{x} , then $\mathbf{v} < \mathbf{x}$, which of course is not the case for lexorder. (This is a property shared by the usual ordering of the integers, written in the usual way without leading zeros: any proper subsequence of an integer i , taken as an integer, is less than i .)

Lemma 3.5. *Given a string \mathbf{x} of length n , let $\mathbf{v} = \mathbf{x}[i_1] \mathbf{x}[i_2] \dots \mathbf{x}[i_r]$ for $1 \leq i_1 < i_2 < \dots < i_r$, $0 \leq r < n$. Then $\mathbf{v} < \mathbf{x}$.*

Proof. Write $V_k(\mathbf{x})$ as in V -form equation (2), and observe that if \mathbf{v} includes j occurrences of g , $0 \leq j < k$, then \mathbf{v} must have ancestor g^j in the star tree but no ancestor g^k . Thus in this case, from Lemma 3.4(2), $\mathbf{v} < \mathbf{x}$.

Suppose then that $j = k$ and write $V_k(\mathbf{v}) = \mathbf{v}_0 g \mathbf{v}_1 g \dots \mathbf{v}_{k-1} g \mathbf{v}_k$, also in V -form (2). Since $|\mathbf{v}| < |\mathbf{x}|$, it follows that there is some least integer $k' \in 0..k$ such that $|\mathbf{v}_{k'}| < |\mathbf{x}_{k'}|$, hence that \mathbf{v} has ancestor

$$A_{s_{k'}}(\mathbf{v}) = \mathbf{x}_0^T g \mathbf{x}_1^T g \dots \mathbf{x}_{k-1}^T g \mathbf{x}_k^T, \tag{3}$$

where $s_{k'} = |\mathbf{v}_{k'}| + |\mathbf{v}_{k'+1}| + \dots + |\mathbf{v}_k|$, and

$$\begin{aligned} \mathbf{x}_h^T &= \mathbf{x}_h \quad \text{for } 0 \leq h < k', \\ &= \boldsymbol{\varepsilon} \quad \text{for } k' \leq h \leq k. \end{aligned}$$

Then $A_{s_{k'}}(\mathbf{v}) = \mathbf{z} g^{k-k'}$ for some string \mathbf{z} that either is empty or has suffix g . On the other hand, we know from Lemma 3.3 that \mathbf{x} has ancestor

$$A_{t_{k'}}(\mathbf{x}) = \mathbf{x}_0^T g \mathbf{x}_1^T g \dots \mathbf{x}_{k-1}^T g \mathbf{x}_k^T, \tag{4}$$

where $t_{k'} = |\mathbf{x}_{k'}| + |\mathbf{x}_{k'+1}| + \dots + |\mathbf{x}_k| > s_{k'}$. Comparing (3) and (4), we see that $A_{s_{k'}}(\mathbf{v}) = A_{t_{k'}}(\mathbf{x})$ is a common ancestor of \mathbf{v} and \mathbf{x} , whereas $\bar{\mathbf{v}} = A_{s_{k'}-|\mathbf{v}_{k'}|}(\mathbf{v}) = \mathbf{z} \mathbf{v}_{k'} g^{k-k'}$ and $\bar{\mathbf{x}} = A_{t_{k'}-|\mathbf{x}_{k'}|}(\mathbf{x}) = \mathbf{z} \mathbf{x}_{k'} g^{k-k'}$ are distinct. It follows from Lemma 3.4(3) that, $\mathbf{v} < \mathbf{x}$ (with $|\mathbf{v}| < |\mathbf{x}|$ and largest letter g) if and only if $\mathbf{v}_{k'} < \mathbf{x}_{k'}$ (with $|\mathbf{v}_{k'}| < |\mathbf{x}_{k'}|$ and largest letter $g_1 < g$), and so to prove the lemma we need only consider the proper substrings $\mathbf{v}_{k'}$ and $\mathbf{x}_{k'}$. We can therefore continue recursively, at each stage replacing \mathbf{v}/\mathbf{x} -strings by proper \mathbf{v}/\mathbf{x} -substrings in which every letter is less than the largest letter occurring in the preceding stage. Since \mathbf{x} is finite, and since the \mathbf{v} -substring is always shorter than the \mathbf{x} -substring, it must at some stage be true that the occurrences of the current largest letter in the \mathbf{v} -substring are fewer than its occurrences in the \mathbf{x} -substring. Hence $\mathbf{v} < \mathbf{x}$, as required. \square

Although for strings $u < v$ and w , it does not follow that $uw < v$, we can however extend to the right:

Lemma 3.6. For any strings $u < v$ and w , $u < vw$.

Proof. If $w = \varepsilon$, the result holds trivially; otherwise, applying Lemma 3.5, we have $u < v < vw$. \square

We now introduce a natural analogue of Lyndon words, derived from V -order, which we call V -words.

Definition 3.7 ([9]). A string w over Σ is a V -word if it is the unique minimum in V -order $<$ among the set of rotations of w .

Observe that the word *cabb* in Example 3.2 is a V -word; note also that a V -word cannot be a repetition. We give examples of V -words on the alphabet $1 < 2 < 3 < 4$ that are not Lyndon words on the same ordering:

Example 3.8. (1) The V -ordering of all rotations of the string 13142 gives $42131 < 14213 < 31421 < 13142 < 21314$. Hence 42131 is a V -word, while 13142 is a Lyndon word.

(2) 3211312 and 44124232 are both V -words, while the corresponding Lyndon words are the rotations 1131232 and 12423244, respectively.

In order to describe the structure of V -words, we introduce the idea of **lexicographic extension** $LEX(F)$ based on a specified factorization F of strings.

Definition 3.9 (Modified from [9]). Suppose that according to some factorization F , two strings $u, v \in \Sigma^+$ are expressed in terms of nonempty factors:

$$u = u_1u_2 \cdots u_m, \quad v = v_1v_2 \cdots v_n.$$

Then $u <_{LEX(F)} v$ if and only if one of the following holds:

- (1) u is a proper prefix of v (that is, $u_i = v_i$ for $1 \leq i \leq m < n$); or
- (2) for some $i \in 1.. \min(m, n)$, $u_j = v_j$ for $j = 1, 2, \dots, i - 1$, and $u_i < v_i$.

Of course the choice of the order that is extended from lexorder can be varied: in Definition 3.9 we have used V -order $<$, but for this hybrid method, lexorder with co-lexorder or any other total order could be used instead. Also the definition depends absolutely on the choice of the factorization F ; for example, if $u = v = ab$, but we perversely choose F such that $u_1 = a, u_2 = b, v_1 = ab$, then $u <_{LEX(F)} v$! In order to make the definition useful, we choose the factors associated with every x to be x_1, x_2, \dots, x_k determined by writing x in its unique V -form (2). The pairwise comparisons of Definition 3.9 are then restricted to rotations $x_i x_{i+1} \cdots x_k x_1 \cdots x_{i-1}$, $1 \leq i \leq k$. In this case, we write LEX instead of $LEX(F)$. Consequently, using the fact that $x_i \leq x_j$ implies $gx_i \leq gx_j$, we can establish

Theorem 3.10 ([9]). Let $x \in \Sigma^+$. Then x is a V -word if and only if its V -form (2) has $x_0 = \varepsilon$ with $x_1 x_2 \cdots x_k$ a Lyndon word under lexicographic extension $<_{LEX}$.

An application of Theorem 3.10 and Lemma 2.3 yields the following result:

Theorem 3.11 ([9]). The set of V -words forms a circ-UMFF \mathcal{V} .

From now on we will denote the V -word circ-UMFF by \mathcal{V} .

Notice that the V -form (2) of a V -word x must, by the properties of V -order, begin with the largest letter g (over Σ) in x . Thus in (2) $x_0 = \varepsilon$, and a V -word must take the form

$$gx_1 \cdots gx_{k-1}gx_k. \tag{5}$$

Accordingly the k rotations of (5) (those commencing gx_1, gx_2, \dots, gx_k) are the only ones that need to be considered in order to determine the unique rotation of primitive x which is its V -word. Hence to test if a string w written as in (5) is a V -word it suffices to check that $x_1 x_2 \cdots x_k$ is Lyndon under lex-extension. Further, if $w_1 w_2 \cdots w_m$ is any maximal factorization by V -words, then by Theorem 2.9, $w_1 \geq_v w_2 \geq_v \cdots \geq_v w_m$ while $\mu w_1 \leq \mu w_2 \leq \cdots \leq \mu w_m$. These observations suggest that Duval’s classic Lyndon factorization algorithm [16] can be modified into a V -word factorization algorithm in a fairly straightforward way, as we show later.

We can extend Definition 3.9 in an obvious way to any total order \mathcal{T} , such that $u_i \ll_{\mathcal{T}} v_i$ means u_i comes before v_i in \mathcal{T} , which in turn gives rise to a natural generalization of Lyndon words:

Definition 3.12. Let \mathcal{T} be a total order of Σ^* with the order relation $\ll_{\mathcal{T}}$. Then a string x in V -form with $x_0 = \varepsilon$ is a **Hybrid Lyndon** if and only if it is Lyndon under lexicographic extension $\ll_{\mathcal{T}}$.

Clearly, Hybrid Lyndons satisfy the **xyz Lemma 2.3**, but applied now to prescribed substrings u_i, v_j . It follows that V -words are an instance of Hybrid Lyndons. Further, any Hybrid Lyndon of type Flight Deck can be factored using a modification of Duval’s algorithm, subject to optimizations based on particular features of the given order \mathcal{T} . Observe further that regular Lyndon words are also Hybrid Lyndons since, as in Theorem 3.10, we can ignore all instances of the first letter occurring in the string and compare the remaining substrings in lexorder. Moreover, we can lift the requirement that $x_0 = \varepsilon$ and choose another letter (by property like smallest/largest, or by position) in the string as a marker to distinguish the x_i .

With reference to [Definition 2.12](#), clearly the set of Lyndon words is type Flight Deck, and we now show that the set of V -words is also type Flight Deck but on an inverted alphabet.

Lemma 3.13. *The circ-UMFF \mathcal{V} is type Flight Deck.*

Proof. Let $w \in \mathcal{V}$, so that w is in V -form, with $w[1] = g$. Then g is maximal in w with respect to V -order (that is, in Σ), but applying [Lemma 3.4\(1\)](#), g is minimal in w with respect to \mathcal{V} -order. \square

Further, notice that the Flight Deck condition includes trivial cases where $w[1 \dots n] \in \mathcal{W}$ and, for every $1 < i \leq n$, $w[1] <_{\mathcal{V}} w[i]$; for instance, the Lyndon word $aedcb$ and the V -word $eabcd$. These kinds of simple V -words will be called V -letters, which we assume to be of maximum length (for instance $eabcd$ rather than the prefix ea). Note that similarly to the Lyndon method, a factorization algorithm can just sweep across V -letters looking for an instance of $w[1]$.

We will now show that when a V -word is in the form of a V -letter, then the set of its rotations must be in V -order (see [Examples 3.2](#) and [3.8](#)). But note that this property does not hold for the Lyndon circ-UMFF, as shown by the words abc , acb ; nor is it the case for general V -words, as illustrated by 3132.

Lemma 3.14. *Let w be a V -letter of length n . If $r_i = R_i(w)$, $i = 1, 2, \dots, n$, then $r_1 > r_2 > \dots > r_n = w$.*

Proof. Consider the rotations r_i and r_{i+1} . Since there is one maximal letter g in w , writing these in V -form yields $r_i = x_{i,0}g x_{i,1}$ and $r_{i+1} = x_{i+1,0}g x_{i+1,1}$. Then [Lemma 3.4\(3\)](#) applies, and so we just need to compare $x_{i,0}$ and $x_{i+1,0}$ in V -order. By the rotation operation, $x_{i+1,0}$ is a subsequence (suffix) of $x_{i,0}$, and hence [Lemma 3.5](#) establishes that $r_i > r_{i+1}$. \square

The concept of *shuffling* Lyndon words – that is, interleaving their letters – was introduced in [[19,5](#)], where it was shown that the ordered shuffle of two Lyndon words of the same length is a Lyndon word. Subsequently this operation was considered for V -words in [[9](#)], but not completed there, a problem to which we now return. We clarify that, whereas Lyndon shuffling interleaves letters, for V -words the V -shuffle is the process of interleaving V -letters. So we decompose a V -word which is in V -form into its V -letters; that is, strings of the form $g x_i$.

Lemma 3.15. *Let $v = g v_1 g \dots v_{k-1} g v_k$ and $w = g w_1 g \dots w_{k-1} g w_k$ be V -words with $v <_{\mathcal{V}} w$. Then the ordered V -shuffle $x = g v_1 g w_1 g v_2 g w_2 \dots g v_k g w_k$ is a V -word.*

Proof. This follows from the Lyndon result, along with [Definition 3.9](#) and [Theorem 3.10](#). \square

The Lyndon words $l_1 = aab$ and $l_2 = b$ with $l_1 <_{\mathcal{L}} l_2$ show the necessity for words to be the same length in the shuffle. We can translate this example into ordered V -words to also show this requirement for V -words: $v = cacacb$ and $w = cb$. Similarly, consider the example from [[9](#)],

$$v = 110110000, \quad w = 100010000, \quad v <_{\mathcal{V}} w,$$

whose V -shuffle is not a V -word. Writing

$$v = 1\epsilon 10 1\epsilon 10000, \quad w = 1000 10000$$

in terms of their V -letters, we can see that the number of V -letters is not the same, and so [Lemma 3.15](#) does not apply to their V -shuffle.

With reference to [Definition 2.8](#), we now show that the order of \mathcal{V} is in some cases the same as V -order while in other cases it is reversed. In particular, the alphabets are in opposite order.

Lemma 3.16. *Suppose V -order and \mathcal{V} are defined over an alphabet Σ , with distinct $v, w \in \mathcal{V}$.*

- (1) *If $\mu v < \mu w$, then $v < w$ and $v >_{\mathcal{V}} w$.*
- (2) *If $\mu v = \mu w$ and $\nu v < \nu w$ then $v < w$ and $v <_{\mathcal{V}} w$ if νw is a Hybrid Lyndon.*
- (3) *If $\mu v = \mu w$ and $\nu v = \nu w$, and if $v < w$ then $v <_{\mathcal{V}} w$.*

Proof. (1) The V -order $v < w$ comes from [Lemma 3.4\(1\)](#). From [Theorem 2.7\(2–3\)](#), νw is primitive and either $v <_{\mathcal{V}} w$ or $w <_{\mathcal{V}} v$. Consider all rotations of νw , and the associated V -form of every rotation, $x_0 g x_1 g \dots x_{k-1} g x_k$. Suppose rotation $R_r(\nu w) = r$ starts with μw and rotation $R_{r'}(\nu w) = r'$ starts with a smaller letter. Then in their respective V -forms, $r = \epsilon g x_1 g \dots x_{k-1} g x_k$, $r' = x'_0 g x'_1 g \dots x'_{k-1} g x'_k$ with $x'_0 \neq \epsilon$, and so from [Lemma 3.4\(3\)](#) we have $r < r'$. So the minimum rotation of νw must start with μw , and similarly w starts with μw . Hence $w v \in \mathcal{V}$, and [Definition 2.6](#) applies. When v, w are single letters, we see that these two orders on Σ are opposite.

(2) The V -order $v < w$ is [Lemma 3.4\(2\)](#). The order on \mathcal{V} is determined by lex-extension, [Definition 3.9](#).

(3) This is [Lemma 3.4\(3\)](#) together with lex-extension, [Definition 3.9](#). \square

The pairs ($v = 4142$, $w = 43$) and ($v = 4243$, $w = 41$) illustrate the weak statement of [Lemma 3.16](#) above.

The following example demonstrates the difference between factoring a string into Lyndon words and into V -words.

Example 3.17. Let $w = 33132421$. Then the Lyndon factorization of w is $3 \geq_{\mathcal{L}} 3 \geq_{\mathcal{L}} 13242 \geq_{\mathcal{L}} 1$, while the factorization of w into V -words is $33132 \geq_{\mathcal{V}} 421$.

Clearly, if $\lambda \in \Sigma$ then the repetition λ^k factors into k factors $\lambda \geq_{\mathcal{W}} \lambda \geq_{\mathcal{W}} \dots \geq_{\mathcal{W}} \lambda$ in any \mathcal{W} -order. By contrast in the more general case, we establish distinctness of factorizations considered here. So for a string w , let $L = L(w)$ be the Lyndon factorization of w and let $V = V(w)$ be the V -word factorization of w .

Lemma 3.18. *If a string w on an ordered alphabet Σ contains at least two distinct letters, then $L(w) \neq V(w)$.*

Proof. By hypothesis, w has prefix $w[1..k] = \lambda_1^{k-1}\lambda_2$ for some integer $k > 1$ and letters λ_1 and $\lambda_2 \neq \lambda_1$ in Σ . If $\lambda_1 < \lambda_2$ in Σ , then $\lambda_1^{k-1}\lambda_2$ is a prefix of the first factor of w in L , while the first factor in V is λ_1 . If on the other hand $\lambda_1 > \lambda_2$ in Σ , then λ_1 is the first factor of w in L , while the prefix of the first factor in V is $\lambda_1^{k-1}\lambda_2$. \square

An immediate consequence of this lemma is that for any pair of words (ℓ, v) , $\ell \in \mathcal{L}$, $v \in \mathcal{V}$, and not both single letters, it follows that $\ell \neq v$. For the particular case of a string that decomposes into a repetition of Lyndon words, the distinct partition into V -words consists (almost) of a cyclic shift of those Lyndon words:

Lemma 3.19. *Let w be a string over Σ^+ containing at least two distinct letters. Suppose that $L(w) = \mathbf{l}^j$, where $\mathbf{l} \in \mathcal{L}$ and $j > 1$. Then $V(w) = V(\mathbf{l}_p)\mathbf{v}^kV(\mathbf{l}_s)$, where $\mathbf{v} \in \mathcal{V}$, and $V(\mathbf{l}_p), V(\mathbf{l}_s)$ are the V -word factorizations of a proper prefix, suffix of \mathbf{l} respectively, and $k \geq 1$.*

Proof. Let $g = \mu\mathbf{l}$, then since L is a repetition (not equal to g^j), \mathbf{l} starts with some non-empty prefix \mathbf{l}' where $\mu\mathbf{l}' < g$. We now write $w = \mathbf{l}_1\mathbf{l}_2 \dots \mathbf{l}_j$ as a sequence of V -letters, namely $\mathbf{l}'\mathbf{g}_1\mathbf{g}_2 \dots \mathbf{g}_t$, where each \mathbf{g}_i is a V -letter starting with g . Let $\mathbf{x} = \mathbf{g}_1\mathbf{g}_2 \dots \mathbf{g}_h$ be the substring of V -letters such that \mathbf{g}_h is in \mathbf{l}_1 , and \mathbf{g}_{h+1} is in \mathbf{l}_2 . Then the non-empty suffix of $\mathbf{g}_h, \mathbf{l}'$, is the prefix of \mathbf{l}_2 . Since the Lyndon word \mathbf{l} is primitive, \mathbf{x} is primitive. So let π be the rotation of $\{1, 2, \dots, h\}$ such that $\mathbf{x}^r = \mathbf{g}_{\pi_1}\mathbf{g}_{\pi_2} \dots \mathbf{g}_{\pi_h}$ is a maximal substring of w which satisfies lex-extension. Then $w = \mathbf{x}_s^r(\mathbf{g}_{\pi_1}\mathbf{g}_{\pi_2} \dots \mathbf{g}_{\pi_h})^k\mathbf{x}_p^r = \mathbf{x}_s^r\mathbf{v}^k\mathbf{x}_p^r$, where $\mathbf{x}_s^r, \mathbf{x}_p^r$ are the suffix, prefix of \mathbf{x}^r respectively, $\mathbf{x}_p^r\mathbf{x}_s^r = \mathbf{x}^r$, and $\mathbf{x}_s^r\mathbf{x}_p^r = \mathbf{l}$. From Lemma 3.18, $\mathbf{l} \neq \mathbf{v}$, so we then factor the non-empty $\mathbf{x}_s^r = \mathbf{l}_p, \mathbf{x}_p^r = \mathbf{l}_s$ in the obvious way, giving $V(\mathbf{l}_p), V(\mathbf{l}_s)$. \square

Using the above lemma, we see that the Lyndon factorization $(113232)^3$ has the following V -word decomposition: $(1)(1)(323211)^2(32)(32)$. Naturally we get a reciprocal result if we start with a repetition of a V -word, and deduce a repetition of a Lyndon word with a similar factorization of a prefix and suffix of the V -word.

Alternatively, by considering the string w to be a necklace, that is a circular string, if $L(w) = \mathbf{l}^j$, then trivially by Lemma 3.18 we have the distinct $V(w) = \mathbf{v}^k$.

One of the attractive features of lexorder is that if we have $v < w$, then we can always add a common prefix which does not alter the ordering, such as $uv < uw$. Also, for Lyndon words, if $v <_{\mathcal{L}} w$, both with minimum letter not less than λ , then $\lambda v <_{\mathcal{L}} \lambda w$. We get a similar result for V -words, where we add a V -letter prefix:

Lemma 3.20. *For any strings u and $v < w$ defined on Σ , all with maximum letter less than g , $guv <_{\mathcal{V}} guw$.*

Proof. By Theorem 3.10 both guv and guw are V -words. The result is immediate from Lemma 3.3 for $u = \epsilon$, and for $v = \epsilon$ it follows from the corresponding result for prefixes of general circ-UMFFs (Lemma 2.10 or Lemma 3.5). It suffices therefore to show that $uv < uw$ for nonempty u and v . The proof is by induction. Suppose first that u is a single letter $\lambda < g$, and consider the sequence of deletions in the star tree that yields a common ancestor z of λv and λw . Three outcomes are possible: z has prefix $\mu > \lambda$, $z = \epsilon$, and z has prefix λ . In the first of these cases, the prefix λ has been deleted on both upward paths (from λv and from λw) in the star tree, and so therefore at some point λ was the strictly least letter (but not necessarily unique) remaining in each of the reduced strings. Thus apart from the deletion of λ , the deletions from λv and λw are the same as they would have been from v and w , and so $\lambda v < \lambda w$. Consider the second case. Since $z = \epsilon$, the previous deletions were from strings consisting of single letters, α, β say, where distinct α, β are maximal letters in $\lambda v, \lambda w$ respectively. Suppose that one of α, β is equal to λ , for otherwise the previous argument applies. If $\beta = \lambda$, then $\alpha > \lambda$ contradicting $v < w$. If $\alpha = \lambda$, then $\beta \geq \lambda$, and due to distinctness, $\alpha < \beta$ as required. In the third case, $z = \lambda z^r$, where z^r is the root for v and w , and therefore at the deciding positions on the two paths below the root, the same letters determine $\lambda v < \lambda w$ that also determine $v < w$. Thus the lemma holds for a single letter λ .

Suppose now that the lemma holds for every prefix u of length $\ell \geq 1$, and consider a prefix \bar{u} of length $\ell + 1$. Here there are two cases to consider: either the common ancestor z of $\bar{u}v$ and $\bar{u}w$ has prefix \bar{u} , or not. If so, then \bar{u} has played no role in the decision, and $v < w \implies \bar{u}v < \bar{u}w$. If not, there exist ancestors $u'v'$ and $u''w'$ of $\bar{u}v$ and $\bar{u}w$, respectively, such that $|u'| = |u''| = \ell$. In fact, since the first letter deleted from \bar{u} must be the same in both cases, we see moreover that $u' = u''$, and we can conclude therefore from the inductive assumption that $u'v' < u''w'$, hence that $uv < uw$, as required. \square

Naturally, known results for Lyndon words can also be considered for Hybrid Lyndon words. The following lemma shows a straightforward generalization of the statement that, if \mathbf{p}, \mathbf{s} are a proper prefix, suffix respectively of a Lyndon word L , then $\mathbf{p} <_{\mathcal{L}} L <_{\mathcal{L}} \mathbf{s}$. Note that in the Lyndon case we only require that \mathbf{p} is also Lyndon (cf. Lemma 2.10).² Similarly, for V -words the prefix needs to consist of V -letters, but the suffix is unrestricted.

Lemma 3.21. *Let $w \in \mathcal{V}$ where $w = g\mathbf{x}_1g\mathbf{x}_2 \dots g\mathbf{x}_k$ in V -form (5). If $\mathbf{p} = g\mathbf{x}_1g\mathbf{x}_2 \dots g\mathbf{x}_j$, $1 \leq j < k$, and \mathbf{s} is any proper suffix of w , then $\mathbf{p} <_{\mathcal{V}} w <_{\mathcal{V}} \mathbf{s}$ and $\mathbf{p}\mathbf{s} \in \mathcal{V}$.*

² We apologize for the typo 'prefix' instead of 'suffix' prior to Lemma 2.5 in [12].

Proof. The prefix case is immediate from [Theorem 3.10](#), Lyndon properties and [Definition 2.6](#). Similarly the suffix case is immediate if the suffix consists of V -letters; that is, it starts with the letter g . Otherwise, assume the suffix does not consist only of V -letters; that is, $\mathbf{s} = \mathbf{x}_j^s g \mathbf{x}_{j+1} \cdots g \mathbf{x}_k$, where $1 \leq j \leq k$ and \mathbf{x}_j^s is a proper suffix of \mathbf{x}_j . Then $\mathbf{ws} = g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_k \mathbf{x}_j^s g \mathbf{x}_{j+1} \cdots g \mathbf{x}_k$. To show that $\mathbf{ws} \in \mathcal{V}$, we will consider its rotations, and use the property that if $\mathbf{a}, \mathbf{b}, \mathbf{c}$ are strings and $\mathbf{a} < \mathbf{b}$, then $\mathbf{a} < \mathbf{bc}$.

We need only consider rotations which start with g , so first compare \mathbf{ws} with the row $g \mathbf{x}_k \mathbf{x}_j^s g \mathbf{x}_{j+1} \cdots g \mathbf{x}_k g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_{k-1}$. Consider the new substring $\mathbf{x}_k \mathbf{x}_j^s$. By Lyndon properties and lex-extension we have $\mathbf{x}_1 < \mathbf{x}_k$, and from [Lemma 3.5](#) we have $\mathbf{x}_k < \mathbf{x}_k \mathbf{x}_j^s$, hence $\mathbf{x}_1 < \mathbf{x}_k \mathbf{x}_j^s$ and so the row \mathbf{ws} is smaller in V -order. Any other rotation of \mathbf{ws} commences with a proper suffix of \mathbf{w} consisting of V -letters, and the result follows. \square

The example of the V -word $\mathbf{x} = 321312$, and the V -letter prefix 321 and suffix 21312 , shows that $\mathbf{pws} = 32132131221312$ is a V -word. However, the prefixes 3213 or 32131 , which are not V -letters, yield 3213321312 or 32131321312 , neither of which are V -words.

We conclude this section by studying [Lemma 2.11](#) for V -words, which yields a pseudo-partition of the circ-UMFF \mathcal{V} (it is ‘pseudo’ as each part must contain every letter of the alphabet). From [Theorem 2.7\(5\)](#) we know that we can generally split any factor into two factors. In [12] we described the structure of Lyndon words which yield a unique split: a V -word counterpart is 41424143 . In contrast to factors with a unique split, [Lemma 2.11](#) involves overlapping factors in the sense of [Lemma 2.3](#). Comparing [Lemma 3.22](#) with [Lemma 2.11](#) and its following commentary also shows that not all Lyndon results carry over to Hybrid Lyndons in a straightforward way.

Lemma 3.22. Let $\mathbf{xy}, \mathbf{yz} \in \mathcal{V}$ for some nonempty $\mathbf{x}, \mathbf{y}, \mathbf{z}$.

(1) If $\mathbf{x}[1] = \mathbf{y}[1]$ then $\mathbf{xy} <_{\mathcal{V}} \mathbf{xyyz} <_{\mathcal{V}} \mathbf{xyz} <_{\mathcal{V}} \mathbf{yz}$.

(2) If $\mathbf{x}[1] <_{\mathcal{V}} \mathbf{y}[1]$ then $\mathbf{xy} <_{\mathcal{V}} \mathbf{xyz} <_{\mathcal{V}} \mathbf{xyyz} <_{\mathcal{V}} \mathbf{yz}$.

Proof. First note that since they are the same length, neither $\mathbf{xyyzxyz}$ nor $\mathbf{xyzxyyz}$ can be in each other’s paths in the star tree of V -order. From [Lemma 2.11](#), $\mathbf{xyz}, \mathbf{xyyz} \in \mathcal{V}$, and since \mathbf{y} is nonempty, we can apply [Theorem 2.7\(3\)](#), and consider when either $\mathbf{xyzxyyz}$ or $\mathbf{xyyzxyz}$ is a V -word.

(1) In view of [Lemma 2.11](#), we need only verify that $\mathbf{xyyz} <_{\mathcal{V}} \mathbf{xyz}$. We write $\mathbf{xyyzxyz}$ in suitable V -form ($\mathbf{x}_0 = \epsilon$) where there are two cases, depending on whether or not the suffix \mathbf{z} of \mathbf{yz} starts with g . If \mathbf{z} does start with g , then the V -form is given by

$$\mathbf{w} = g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i (g \mathbf{y}_1 g \mathbf{y}_2 \cdots g \mathbf{y}_j)^2 g \mathbf{z}_1 g \mathbf{z}_2 \cdots g \mathbf{z}_k g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i g \mathbf{y}_1 g \mathbf{y}_2 \cdots g \mathbf{y}_j g \mathbf{z}_1 g \mathbf{z}_2 \cdots g \mathbf{z}_k,$$

where $i, j, k \geq 1$. It is required to show that $\mathbf{xyyzxyz} < \mathbf{xyzxyyz}$, so we compare \mathbf{w} with the rotation

$$\mathbf{w}^r = g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i g \mathbf{y}_1 g \mathbf{y}_2 \cdots g \mathbf{y}_j g \mathbf{z}_1 g \mathbf{z}_2 \cdots g \mathbf{z}_k g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i (g \mathbf{y}_1 g \mathbf{y}_2 \cdots g \mathbf{y}_j)^2 g \mathbf{z}_1 g \mathbf{z}_2 \cdots g \mathbf{z}_k.$$

Using lex-extension ordering, since there is a common prefix \mathbf{xy} in \mathbf{w} and \mathbf{w}^r , we compare the substrings $g \mathbf{y}_1 g \mathbf{y}_2 \cdots g \mathbf{y}_j$ and $g \mathbf{z}_1 g \mathbf{z}_2 \cdots g \mathbf{z}_k$. Since $\mathbf{yz} \in \mathcal{V}$, by [Theorem 2.7\(1\)](#) it is border-free, hence $g \mathbf{y}_1 g \mathbf{y}_2 \cdots g \mathbf{y}_j < g \mathbf{z}_1 g \mathbf{z}_2 \cdots g \mathbf{z}_k$. Otherwise, if \mathbf{z} does not contain the letter g , the comparison of the rows \mathbf{w} and \mathbf{w}^r depends on the comparison of the substrings $\mathbf{y}_j < \mathbf{y}_j \mathbf{z}$. In this case, since \mathbf{z} is nonempty, we apply [Lemma 3.5](#), establishing that $\mathbf{y}_j < \mathbf{y}_j \mathbf{z}$. Hence $\mathbf{xyyzxyz}$ satisfies the lex-extension [Theorem 3.10](#), as required.

(2) Here we need only verify that $\mathbf{xyz} <_{\mathcal{V}} \mathbf{xyyz}$. We know that $\mathbf{x}[1] = \mu(\mathbf{xyzxyyz}) = \mu(\mathbf{xyyzxyz})$, and in particular $\mu = \mathbf{x}[1] <_{\mathcal{V}} \mathbf{y}[1]$. We write $\mathbf{xyzxyyz}$ in suitable V -form, that is, $\mathbf{w} = g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i \mathbf{y} \mathbf{z} g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i \mathbf{y} \mathbf{z}$, $i \geq 1$. So we compare \mathbf{w} with the rotation $g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i \mathbf{y} \mathbf{z} g \mathbf{x}_1 g \mathbf{x}_2 \cdots g \mathbf{x}_i \mathbf{y} \mathbf{z}$. Applying [Lemma 3.5](#) with nonempty \mathbf{y} , we see that $\mathbf{x}_i \mathbf{y} \mathbf{z} < \mathbf{x}_i \mathbf{y} \mathbf{y} \mathbf{z}$, and similarly $\mathbf{xyzxyyz}$ satisfies the lex-extension [Theorem 3.10](#) for this case also.

Recalling that the case $\mathbf{x}[1] >_{\mathcal{V}} \mathbf{y}[1]$ is not possible in a Flight Deck circ-UMFF completes the proof. \square

Example 3.23. Let $\Sigma = \{1 < 2 < 3 < \cdots\}$, and so $\{1 >_{\mathcal{V}} 2 >_{\mathcal{V}} 3 >_{\mathcal{V}} \cdots\}$.

(1) Let $\mathbf{xy} = 4142$ and $\mathbf{yz} = 4243$, hence $\mathbf{xyz} = 414243$ and $\mathbf{xyyz} = 41424243$, where $\mathbf{xy}, \mathbf{yz}, \mathbf{xyz}, \mathbf{xyyz} \in \mathcal{V}$. Then $4142 <_{\mathcal{V}} 41424243 <_{\mathcal{V}} 414243 <_{\mathcal{V}} 4243$.

(2) Let $\mathbf{xy} = 4142$ and $\mathbf{yz} = 21$, hence $\mathbf{xyz} = 41421$ and $\mathbf{xyyz} = 414221$, where $\mathbf{xy}, \mathbf{yz}, \mathbf{xyz}, \mathbf{xyyz} \in \mathcal{V}$. Then $4142 <_{\mathcal{V}} 41421 <_{\mathcal{V}} 414221 <_{\mathcal{V}} 21$.

Notice that in [Lemma 3.22\(1\)](#) we have $\mathbf{xyyz} <_{\mathcal{V}} \mathbf{xyz}$, while [Lemma 3.22\(2\)](#) yields $\mathbf{xyz} <_{\mathcal{V}} \mathbf{xyyz}$, so the conditions on $\mathbf{x}[1], \mathbf{y}[1]$ imply that it is not possible to apply the [xyz Lemma 2.3](#) in such a way as to generate a bordered word. Hence we have a pseudo-partition of the \mathbf{xyz} type strings in \mathcal{V} .

4. Linear string $<$ comparison

In this section we describe a new algorithm for V -order $<$ comparison of finite strings in worst-case time linear in their lengths, thus asymptotically the same as comparison in lexorder. We suppose throughout that any string \mathbf{x} is defined on a finite ordered alphabet Σ of size $\sigma = |\Sigma|$.

```

– Left-extend the suffix of nondecreasing letters, then delete.
– Position  $i$  is next right of the deleted position  $\delta$ .
function delete( $\mathbf{x}$ ,  $i$ ) :  $\mathbf{x}$ ,  $i$ ,  $\delta$ 
  prev  $\leftarrow \mathbf{x}[i].\text{left}$ 
  while  $\mathbf{x}[i].\lambda \geq \mathbf{x}[\text{prev}].\lambda$  do
     $i \leftarrow \text{prev}$ ; prev  $\leftarrow \mathbf{x}[i].\text{left}$ 
   $\delta \leftarrow i$ ;  $i \leftarrow \mathbf{x}[\delta].\text{right}$ 
   $\mathbf{x}[i].\text{left} \leftarrow \text{prev}$ ;  $\mathbf{x}[\text{prev}].\text{right} \leftarrow i$ 

– Perform  $n_2 - n_1$  successive deletions.
function reduce( $\mathbf{x}_2$ ,  $n_2$ ,  $n_1$ ) :  $\mathbf{x}_2$ ,  $i_2$ 
   $i_2 \leftarrow n_2 + 1$ 
  for  $j \leftarrow 1$  to  $n_2 - n_1$  do
    ( $\mathbf{x}_2$ ,  $i_2$ ,  $\delta$ )  $\leftarrow$  delete( $\mathbf{x}_2$ ,  $i_2$ )

```

Fig. 1. Reduction of \mathbf{x}_2 to n_1 letters.

In order to simplify the description of our algorithm, we suppose without loss of generality that the letters $\lambda \in \Sigma$ are drawn from the first σ natural numbers. We represent strings $\mathbf{x} = \mathbf{x}[1..n]$ using added sentinel positions $\mathbf{x}[0]$ and $\mathbf{x}[n+1]$ for processing convenience. At each position $i \in 1..n$, we define a triple $(\lambda, \text{left}, \text{right})$, where λ is the letter at position i , $\text{left} = i-1$, and $\text{right} = i+1$; for $i = 0$, $(\lambda, \text{left}, \text{right}) = (\sigma+1, -1, 1)$ and for $i = n+1$, $(\lambda, \text{left}, \text{right}) = (\sigma+1, n, n+2)$. The elements of the triple at position i are denoted $\mathbf{x}[i].\lambda$, $\mathbf{x}[i].\text{left}$, and $\mathbf{x}[i].\text{right}$. Thus we represent \mathbf{x} as a linked list to facilitate the deletion of positions in \mathbf{x} required to locate parent nodes in the star tree used to define V -order (Definition 3.1). Note that the linked list operations consist of: comparison, reading and writing to list cells, which can be implemented in constant time.

We suppose that two unequal strings \mathbf{x}_1 and \mathbf{x}_2 are given, both defined on Σ , of lengths $n_1 > 0$ and $n_2 \geq n_1$ respectively. We describe a $\Theta(n_2)$ -time algorithm to solve the following problem:

(P1) Determine whether $\mathbf{x}_1 < \mathbf{x}_2$ or $\mathbf{x}_1 > \mathbf{x}_2$.

Recall that in order to compare two strings in V -order, it is necessary first to identify a common ancestor \mathbf{v} (perhaps the empty string ϵ) in the star tree, and then to compare the children of \mathbf{v} . The common ancestor is found by performing iterated deletes in the given strings according to the following rule:

Delete the letter at the rightmost position $i \in \{1, \dots, n\}$ of \mathbf{x} such that $\mathbf{x}[i] \leq \mathbf{x}[i+1]$ and $\mathbf{x}[i] < \mathbf{x}[i-1]$. (Note that the sentinel positions 0 and $n+1$ ensure the correct application of this rule.)

Since no comparison can be made until \mathbf{x}_2 is reduced to the same length as \mathbf{x}_1 , we begin by performing $n_2 - n_1$ deletions from \mathbf{x}_2 , as shown in Fig. 1. Then we continue by performing single deletions from both \mathbf{x}_1 and \mathbf{x}_2 , checking at each step to determine whether or not the current deletion has made the reduced \mathbf{x}_1 equal to the reduced \mathbf{x}_2 . To facilitate this checking, we introduce the idea of the **longest matching suffix** (LMS): the longest suffix of the reduced \mathbf{x}_1 that matches a suffix of the reduced \mathbf{x}_2 . Of course LMS is a subsequence (possibly empty) of both \mathbf{x}_1 and \mathbf{x}_2 ; it is useful to identify the leftmost position in LMS as LMS- l , and the position immediately to the left of LMS- l as LMS- ll . Clearly, when LMS(\mathbf{x}_1 , \mathbf{x}_2) equals the length of the reduced strings, the common ancestor of \mathbf{x}_1 and \mathbf{x}_2 in the star tree has been found (ϵ when LMS is empty). We make the following observations:

Remark 4.1. A position (not LMS- l) of \mathbf{x}_1 (respectively, \mathbf{x}_2) is deleted within LMS(\mathbf{x}_1 , \mathbf{x}_2) if and only if the corresponding position (also not LMS- l) within \mathbf{x}_2 (respectively, \mathbf{x}_1) is also deleted. In this case LMS cannot be extended to the left, but is reduced in length by one.

Remark 4.2. LMS cannot be extended to the left if either LMS- l is deleted in either \mathbf{x}_1 or \mathbf{x}_2 or positions left of LMS- ll are deleted in both \mathbf{x}_1 and \mathbf{x}_2 . In the latter case, LMS is unchanged; in the former it is reduced in length by one.

Remark 4.3. The LMS may possibly (but not necessarily) be extended to the left only if either LMS- ll is deleted in both \mathbf{x}_1 and \mathbf{x}_2 or LMS- ll is deleted in \mathbf{x}_1 (respectively, \mathbf{x}_2) and a position left of LMS- ll is deleted in \mathbf{x}_2 (respectively, \mathbf{x}_1).

These are the essential facts exploited in the function *prec* (Fig. 2) that returns TRUE or FALSE depending on whether or not $\mathbf{x}_1 < \mathbf{x}_2$.

Lemma 4.4. Procedure *prec* executes in time and space linear in $|\mathbf{x}_2|$.

Proof. We consider the efficiency of function *prec*. In terms of space, *prec* requires only three storage locations for each position in \mathbf{x}_1 and \mathbf{x}_2 , namely the triple $(\lambda, \text{left}, \text{right})$, in addition to constant space for variables and program storage; thus its space requirement is $\Theta(n_1 + n_2) = \Theta(n_2)$. To estimate time usage, first consider the function *delete*. This function may left-extend the current suffix of nondecreasing letters, but the deleted position is always at the left end of the suffix. It follows that over all the executions of *delete* (at most $n_1 + n_2$ of them), each position in \mathbf{x}_1 and \mathbf{x}_2 is visited at most twice,

```

– Recompute the positions  $\ell_1, \ell_2$  that define LMS- $\ell$ .
function LMS( $\mathbf{x}_1, \ell_1, \mathbf{x}_2, \ell_2$ ) :  $\ell_1, \ell_2$ 
 $i_1 \leftarrow \mathbf{x}_1[\ell_1].\text{left}; i_2 \leftarrow \mathbf{x}_2[\ell_2].\text{left}$ 
while  $\mathbf{x}_1[i_1].\lambda = \mathbf{x}_2[i_2].\lambda$  and  $i_1 > 0$  do
     $i_1 \leftarrow \mathbf{x}_1[i_1].\text{left}; i_2 \leftarrow \mathbf{x}_2[i_2].\text{left}$ 
 $\ell_1 \leftarrow \mathbf{x}_1[i_1].\text{right}; \ell_2 \leftarrow \mathbf{x}_2[i_2].\text{right}$ 

– Determine whether or not  $\mathbf{x}_1 < \mathbf{x}_2$ .
function prec( $\mathbf{x}_1, n_1, \mathbf{x}_2, n_2$ ) : boolean
( $\mathbf{x}_2, i_2$ )  $\leftarrow$  reduce( $\mathbf{x}_2, n_2, n_1$ )
( $\ell_1, \ell_2$ )  $\leftarrow$  LMS( $\mathbf{x}_1, n_1 + 1, \mathbf{x}_2, n_2 + 1$ )
– Since  $\mathbf{x}_1 \neq \mathbf{x}_2$ ,  $\mathbf{x}_1[\ell_1].\text{left} = 0$  implies that  $\mathbf{x}_1$  lies
– on the upward path from  $\mathbf{x}_2$  in the star tree.
if  $\mathbf{x}_1[\ell_1].\text{left} = 0$  then return TRUE
 $i_1 \leftarrow n_1 + 1$  –  $i_2$  already set by reduce( $\mathbf{x}_2, n_2, n_1$ )
repeat
    ( $\mathbf{x}_1, i_1, \delta_1$ )  $\leftarrow$  delete( $\mathbf{x}_1, i_1$ ); ( $\mathbf{x}_2, i_2, \delta_2$ )  $\leftarrow$  delete( $\mathbf{x}_2, i_2$ )
    – Remark 4.2: possibly LMS- $\ell$  was deleted in at least
    – one of  $\mathbf{x}_1, \mathbf{x}_2$ ; if so, shift LMS- $\ell$  right in both  $\mathbf{x}_1, \mathbf{x}_2$ .
    if  $\delta_1 = \ell_1$  or  $\delta_2 = \ell_2$  then
         $\ell_1 \leftarrow \mathbf{x}_1[\ell_1].\text{right}; \ell_2 \leftarrow \mathbf{x}_2[\ell_2].\text{right}$ 
        – Remark 4.3: possibly LMS- $\ell$  was deleted in one
        – of  $\mathbf{x}_1, \mathbf{x}_2$ , while LMS- $\ell$  was not deleted in either.
    elseif  $\ell_1 = i_1$  or  $\ell_2 = i_2$  then
        ( $\ell_1, \ell_2$ )  $\leftarrow$  LMS( $\mathbf{x}_1, \ell_1, \mathbf{x}_2, \ell_2$ )
until  $\mathbf{x}_1[\ell_1].\text{left} = 0$  – ( $\mathbf{x}_1[\ell_1].\text{left} = 0 \Leftrightarrow \mathbf{x}_2[\ell_2].\text{left} = 0$ )
    – V-order is determined by the lexorder of the last deleted letters.
if  $\mathbf{x}_1[\delta_1].\lambda < \mathbf{x}_2[\delta_2].\lambda$  then
    return TRUE
else
    return FALSE

```

Fig. 2. Match the reduced \mathbf{x}_1 and \mathbf{x}_2 to compare $\mathbf{x}_1 : \mathbf{x}_2$.

with constant-time processing corresponding to each visit. These visits comprise either the scanning of letters from right to left, or, while re-scanning from left to right, a deletion. (Any deletions within LMS are achieved by the function *delete* that updates the linked list in constant time.) Similarly, the calls to function LMS (at most n_1 altogether) go from right to left without backtracking (although LMS- ℓ can be decremented), and so each position is visited at most twice, each with constant-time processing. We conclude that the time requirement for *prec* is $\Theta(n_1 + n_2)$; that is, linear in the lengths of the compared strings. \square

5. Lyndon-like factorization using V-order

In this section we describe an asymptotically optimal algorithm for Lyndon-like factorization using V-order. Using the basic structure of Duval's Lyndon factorization algorithm [16] (while replacing letter by substring comparison), we implement efficiencies based on features specific to V-order (Lemma 3.4), while also applying the fast function for $<$ string comparison of Section 4.

Thus, denoting by $VF(\mathbf{x})$ the unique sequence of V-words of an input string \mathbf{x} we have:

(P2) Given a nonempty string \mathbf{x} on a finite ordered alphabet Σ , compute $VF(\mathbf{x})$.

In order to compute $VF(\mathbf{x})$, we need to identify the unique subset of V-words $\mathbf{w}_j, j = 1, 2, \dots, J$, such that $\mathbf{x} = \mathbf{w}_1\mathbf{w}_2 \cdots \mathbf{w}_J$ and $\mathbf{w}_j \geq_v \mathbf{w}_{j+1}$ for every $j \in \{1, \dots, J - 1\}$, thus satisfying the necessary factorization criteria given by Corollary 2.4 and Theorem 2.9. Observe too that when applying lex-extension (Definition 3.9) to V-word Eq. (5), for any comparison yielding $\mathbf{x}_i = \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon}$ is the least string in V-order (the root of the star tree).

Fig. 3 gives pseudocode for a linear-time algorithm VF to solve (P2), based on Duval's algorithm [16] as presented in [22]. Algorithm VF essentially considers two adjacent substrings of $\mathbf{x}[1..n]$: a left hand one \mathbf{u} (that may be a repetition repeated rep times, and whose positions are tracked by i) of length $rep \times \ell_1$ with a prefix of k_1 g's; and a right hand one \mathbf{v} (tracked by j) of length ℓ_2 with a prefix of k_2 g's. At each step of the algorithm j is incremented by one, and then tests are performed to determine whether as a result \mathbf{u}^{rep} and \mathbf{v} should be concatenated into a single V-word, or whether each of rep occurrences of \mathbf{u} (and perhaps also \mathbf{v} , if g has increased) should be output as the current V-word(s), or whether no decision can currently be taken. The variable h gives the rightmost position of the last V-word output (thus one less than the starting position of

```

procedure VF( $\mathbf{x}$ )
  – Given the input string  $\mathbf{x}\$,$  output in ascending sequence
  the rightmost positions of every V-word in VF( $\mathbf{x}$ )
   $h \leftarrow 0$  – the total length of VF already output
  while  $h < n$  do – continue as long as total VF  $< n$ 
    RESET
    while  $\mathbf{x}[j] \leq g$  do
      if  $\mathbf{x}[i] = \mathbf{x}[j] = g$  then
        – Extend the  $g$ -prefix of either  $\mathbf{u}$  or  $\mathbf{v}$ .
        if  $j - i = k_1$  then  $k_1 \leftarrow k_1 + 1$ ;  $rep \leftarrow rep + 1$ ;  $\ell_1 \leftarrow \ell_1 + 1$ 
        else
          if  $i - h \leq k_1$  then  $k_2 \leftarrow k_2 + 1$ 
           $\ell_2 \leftarrow \ell_2 + 1$ ;  $i \leftarrow i + 1$ 
      elseif  $\mathbf{x}[j] = g$  then
        – The copy  $\mathbf{v}$  of  $\mathbf{u}$  truncated by a mismatched  $g$ :
        output  $\mathbf{u}^{rep}$  and restart  $\mathbf{u} \leftarrow \mathbf{v}$ .
        output ( $h, rep, \ell_1$ ); RESET;  $j \leftarrow j - 1$ 
      elseif  $\mathbf{x}[i] = g$  then
        –  $\mathbf{u}$  contains a mismatched  $g$ : if  $k_2 = 0$ , continue  $\mathbf{u}$ 
        by incrementing  $\ell_1$ ; otherwise, concatenate  $\mathbf{u} \leftarrow \mathbf{u}^{rep}\mathbf{v}$ .
        if  $j = k_1 + 1$  then  $\ell_1 \leftarrow k_1 + 1$ 
        else  $\ell_1 \leftarrow (rep \times \ell_1) + 1$ 
         $rep \leftarrow 1$ ;  $i \leftarrow h + 1$ 
        if  $k_2 > 0$  then  $\ell_1 \leftarrow \ell_1 + \ell_2$ ;  $k_2, \ell_2 \leftarrow 0$ 
      elseif  $\mathbf{x}[i] = \mathbf{x}[j]$  then
        – Corresponding non- $g$  positions in  $\mathbf{u}$  &  $\mathbf{v}$  are equal.
         $\ell_2 \leftarrow \ell_2 + 1$ 
        if  $\ell_2 = \ell_1$  then  $rep \leftarrow rep + 1$ ;  $k_2, \ell_2 \leftarrow 0$ ;  $i \leftarrow h + 1$ 
        else  $i \leftarrow i + 1$ 
      else
        – Unequal non- $g$  positions: compare  $\mathbf{u}$  & extended  $\mathbf{v}$ .
         $j' \leftarrow j$ ; while  $\mathbf{x}[j + 1] < g$  do  $j \leftarrow j + 1$ 
        – If  $\mathbf{u} > \mathbf{v}$ , output  $\mathbf{u}^{rep}$ ; otherwise, concatenate  $\mathbf{u} \leftarrow \mathbf{u}^{rep}\mathbf{v}$ .
        if not  $prec(\mathbf{x}[h + 1..h + \ell_1], \ell_1, \mathbf{x}[j' - \ell_2..j], j - j' + \ell_2 + 1)$  then
          output ( $h, rep, \ell_1$ )
           $\ell_1 \leftarrow j - h$ ;  $rep \leftarrow 1$ ;  $k_2, \ell_2 \leftarrow 0$ ;  $i \leftarrow h + 1$ 
         $j \leftarrow j + 1$ 
    output ( $h, rep, \ell_1$ ); output ( $h, 1, \ell_2$ )

```

Fig. 3. Compute the $O(n)$ partition VF(\mathbf{x}).

\mathbf{u}^{rep}). VF invokes three other routines: the linear function $prec$ of Section 4, an initialization procedure RESET:

$$i \leftarrow h + 1; j \leftarrow h + 2; g \leftarrow \mathbf{x}[i]$$

$$k_1, \ell_1, rep \leftarrow 1; k_2, \ell_2 \leftarrow 0$$

and a procedure **output** (h, rep, ℓ) that for $\ell > 0$ outputs the rightmost positions h of V-word prefixes of $\mathbf{u}^{rep}\mathbf{v}$ (and updates the global variable h):

$$\mathbf{if} \ell > 0 \mathbf{then}$$

$$\mathbf{for} i \leftarrow 1 \mathbf{to} rep \mathbf{do}$$

$$h \leftarrow h + \ell; \mathbf{output} h$$

For processing convenience, the input string to procedure VF is actually $\mathbf{x}\$,$ where $\$$ is a sentinel letter greater than any letter of Σ (chosen to be $\sigma + 1$ in Section 4).

Lemma 5.1. Procedure VF(\mathbf{x}) executes correctly in time and space linear in $|\mathbf{x}|$.

Proof. Analogous to Lyndon factorization, V-word factorization maintains current candidates for factors (corresponding to the current maximum letter g) in the form $\mathbf{u}^{rep}\mathbf{v}$, subject to the following rules:

- $\mathbf{u} = g^{k_1}\mathbf{u}'$, where \mathbf{u}' contains no letter $\lambda > g$ and no substring g^k , $k \geq k_1$;
- $\mathbf{v} = g^{k_2}$, $0 \leq k_2 < k_1$, or $\mathbf{v} = g^{k_1}\mathbf{v}'$, where \mathbf{v}' is a proper prefix of \mathbf{u}' ;
- if it becomes true that $\mathbf{v}' = \mathbf{u}'$, set $rep \leftarrow rep + 1$, $\mathbf{v} \leftarrow \epsilon$.

Only in the case that distinct letters $\lambda_1 < g$ and $\lambda_2 < g$ are found at corresponding positions i in \mathbf{u} and j in \mathbf{v} , respectively, does it become necessary to call *prec* to determine whether or not $\mathbf{u} < \mathbf{v}$ – with \mathbf{v} first extended to the right until some letter $\lambda \geq g$ (possibly the sentinel \$) is found. We remark that this case can occur only if $k_2 = k_1$, with the current $\ell_2 < \ell_1$.

To establish the linearity of VF's execution time, note first that all the outputs taken together require at most $O(n)$ time. Also, the restarts that occur when $\mathbf{x}[j] = g \neq \mathbf{x}[i]$ can total at most $O(n)$ time. Otherwise, all operations performed at each increment of j require only constant time except for the call to *prec*. Observe that the total number of positions tested in *prec* over all calls is at most the sum of all the terms $j - j' + \ell_1 + \ell_2 + 1$. If a call to *prec* in procedure VF yields $\mathbf{u} > \mathbf{v}$, then an output occurs, and $\mathbf{u} \leftarrow \mathbf{v}$; if not, however, then $\mathbf{u} \leftarrow \mathbf{uv}$ (Theorem 2.9). In the latter case, \mathbf{u} necessarily increases in size, so that the next call to *prec*, if it occurs, requires that correspondingly more positions need to be tested. We see that repeated consecutive invocations of *prec*, without output, must yield a worst case for total time, since each call must include all the positions already tested in previous calls. Note however that the result $\mathbf{u} < \mathbf{v}$ can occur only if the number of g 's in \mathbf{v} is at most equal to the number of g 's in \mathbf{u} , subject to the constraint noted above that $k_2 = k_1$; thus the lengths of both \mathbf{u} and \mathbf{v} must approximately double at each step, as in the example

$$\mathbf{x} = g1g2g1g3g1g2g1g4 \dots,$$

where we first find $g1 < g2$, then $g1g2 < g1g3$, $g1g2g1g3 < g1g2g1g4$, and so on. This example in fact must constitute a worst case for the number of positions tested in *prec*; assuming without loss of generality that $|\mathbf{x}| = 2^m$, we find that the total number of positions tested will be

$$2^m + 2^{m-1} + \dots + 2 = 2(|\mathbf{x}| - 1). \quad \square \tag{6}$$

6. Conclusion

In this paper we have designed new algorithms for V -words, which are analogous structures to the well-known Lyndon words. Linear V -order string comparison supports a V -word factorization algorithm with linear time and space complexities. Hence we have established that important Lyndon and V -word operations have the same asymptotic complexity.

Additionally, we have presented new combinatorial results on V -order and V -words, thus identifying similarities and differences between their inherent orders and that of lexicographic order and Lyndon words. Lexicographic order is primarily a positional ordering method, whereas V -order has the potential of defining order in terms of semantics or patterns rather than position. We suggest that applications of V -order and V -words may arise in music analysis (see [1,2] for applications of circ-UMFFs in musicology).

Future research into string factorization is suggested by the novel Hybrid Lyndons and generalized lex-extension method described here, as they are adaptable to a variety of ordering techniques. We propose a suite of binary factorization algorithms to be engineered for the 32 Block-like UMFFs defined in [9], along with binary Lyndon and binary V -word factorization. En route to more general circ-UMFFs, a ternary circ-UMFF conjectured in [11] may also be suitable for algorithmics. The existence of alternative factoring methods opens avenues for optimization problems on string partitioning.

References

- [1] M. Chemillier, Periodic musical sequences and Lyndon words, *Soft Comput.* (8–9) (2004) 611–616. ISSN 1432–7643 (Print) 1433–7479 (Online).
- [2] M. Chemillier, C. Truchet, Computation of words satisfying the rhythmic oddity property (after Simha Arom's works), *Inform. Process. Lett.* 86 (2003) 255–261.
- [3] M. Crochemore, J. Désarménien, D. Perrin, A note on the Burrows-Wheeler transformation, *Theoret. Comput. Sci.* 332 (1–3) (2005) 567–572.
- [4] K.T. Chen, R.H. Fox, R.C. Lyndon, Free differential calculus IV—the quotient groups of the lower central series, *Ann. Math.* 68 (1958) 81–95.
- [5] L.J. Cummings, M.E. Mayes, Shuffled Lyndon words, *ARS Combin.* 33 (1992) 47–56.
- [6] D.E. Daykin, Ordered ranked posets, representations of integers and inequalities from extremal poset problems, in: *Graphs and Order*, in: I. Rival (Ed.), *Proceedings of a Conference in Banff, Canada (1984)*, NATO Advanced Sciences Institutes Series C: Mathematical and Physical Sciences, vol. 147, Reidel, Dordrecht-Boston, 1985, pp. 395–412.
- [7] D.E. Daykin, Algorithms for the Lyndon unique maximal factorization, *J. Combin. Math. Combin. Comput.* 77 (2011) 65–74.
- [8] T.-N. Danh, D.E. Daykin, The structure of V -order for integer vectors, in: A.J.W. Hilton (Ed.), *Congressus Numerantium*, vol. 113, Utilitas Mat. Pub. Inc., Winnipeg, Canada, 1996, pp. 43–53.
- [9] D.E. Daykin, J.W. Daykin, Lyndon-like and V -order factorizations of strings, *J. Discrete Algorithms* 1 (2003) 357–365.
- [10] D.E. Daykin, J.W. Daykin, Properties and construction of unique maximal factorization families for strings, *Internat. J. Found. Comput. Sci.* 19–4 (2008) 1073–1084.
- [11] D.E. Daykin, J.W. Daykin, C.S. Iliopoulos, W.F. Smyth, Generic Algorithms for Factoring Strings (submitted for publication).
- [12] D.E. Daykin, J.W. Daykin, W.F. Smyth, Combinatorics of Unique Maximal Factorization Families (UMFFs), in: R. Janicki, S.J. Puglisi, M.S. Rahman (Eds.), *Special String Masters Issue*, *Fund. Inform.* 97 (3) (2009), 295–309.
- [13] D.E. Daykin, J.W. Daykin, W.F. Smyth, String comparison and Lyndon-like factorization using V -order in linear time, in: R. Giancarlo, G. Manzini (Eds.), *Proc. 22nd Annual Symposium on Combinatorial Pattern Matching*, in: LNCS, vol. 6661, 2011, pp. 65–76.
- [14] J.W. Daykin, C.S. Iliopoulos, W.F. Smyth, Parallel RAM algorithms for factorizing words, *Theoret. Comput. Sci.* 127 (1) (1994) 53–67.
- [15] O. Delgrange, E. Rivals, STAR: an algorithm to search for tandem approximate repeats, *Bioinformatics* 20–16 (2004) 2812–2820.
- [16] J.P. Duval, Factorizing words over an ordered alphabet, *J. Algorithms* 4 (1983) 363–381.
- [17] J. Gil, D.A. Scott, A bijective string sorting transform, <http://bijective.dogma.net/00yyy.pdf> (submitted for publication).
- [18] C.S. Iliopoulos, W.F. Smyth, Optimal algorithms for computing the canonical form of a circular string, *Theoret. Comput. Sci.* 92 (1) (1992) 87–105.
- [19] M. Lothaire, *Combinatorics on Words*, 2nd Edition, Cambridge University Press, Cambridge, 1997.
- [20] L. Perret, A chosen ciphertext attack on a public key cryptosystem based on Lyndon words, in: *Proceedings of International Workshop on Coding and Cryptography*, January 2005, pp. 235–244.
- [21] C. Reutenauer, *Free Lie algebras*, in: *London Math. Soc. Monographs New Ser.*, vol. 7, Oxford University Press, 1993, 288 pp.
- [22] Bill Smyth, *Computing Patterns in Strings*, Pearson, 2003, 423 pp.