

Aberystwyth University

Computation of the suffix array, burrows-wheeler transform and FM-index in V-order

Daykin, Jacqueline; Mhaskar, Neerja; Smyth, W. F.

Published in:
Theoretical Computer Science

DOI:
[10.1016/j.tcs.2021.06.004](https://doi.org/10.1016/j.tcs.2021.06.004)

Publication date:
2021

Citation for published version (APA):
Daykin, J., Mhaskar, N., & Smyth, W. F. (2021). Computation of the suffix array, burrows-wheeler transform and FM-index in V-order. *Theoretical Computer Science*, 880, 82-96. <https://doi.org/10.1016/j.tcs.2021.06.004>

Document License CC BY-NC-ND

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

Computation of the Suffix Array, Burrows-Wheeler Transform and FM-Index in V -Order*

Jacqueline W. Daykin^{e,a,b}, Neerja Mhaskar^{c,*}, W. F. Smyth^{e,c,d}

^a*Department of Computer Science, Aberystwyth University, Wales*

^b*Department of Information Science, Stellenbosch University, South Africa*

^c*Department of Computing and Software, McMaster University, Canada*

^d*School of Engineering & Information Technology, Murdoch University, Australia*

^e*Department of Informatics, King's College London, UK*

Abstract

V -order is a total order on strings that determines an instance of Unique Maximal Factorization Families (UMFFs), a generalization of Lyndon words. The fundamental V -comparison of strings can be done in linear time and constant space. V -order has been proposed as an alternative to lexicographic order (lex-order) in the computation of suffix arrays and in the suffix-sorting induced by the Burrows-Wheeler transform (BWT). In line with the recent interest in the connection between suffix arrays and Lyndon factorization, in this paper we obtain similar results for the V -order factorization. Indeed, we show that the results describing the connection between suffix arrays and Lyndon factorization are matched by analogous V -order processing. We also describe a methodology for efficiently computing the FM-Index in V -order, as well as V -order substring pattern matching using backward search.

Keywords: Burrows-Wheeler transform, Combinatorics, FM-index, Lexorder, String comparison, Substring pattern matching, Suffix sorting, V -BWT, V -order.

*The paper is an outgrowth of a contribution to WALCOM 2019, published in [1].

*Corresponding author.

Email addresses: jackie.daykin@gmail.com (Jacqueline W. Daykin),
pophlin@mcmaster.ca (Neerja Mhaskar), smyth@mcmaster.ca (W. F. Smyth)

1. Introduction

This paper extends current knowledge on applications of a non-lexicographic global order known as V -order [2] (see Definition 2). It is an intriguing question, now investigated for more than 20 years, whether such a counter-intuitive ordering might nevertheless yield algorithmic efficiencies, or other theoretical/computational benefits, compared to lexicographic order (*lexorder*). We begin by mentioning some applications already considered.

Analogous to lexicographic comparison, the central problem of efficient V -ordering of strings was considered in [3, 4, 5, 6, 7], culminating in a remarkably simple, linear time, constant space comparison algorithm [8], further improved in [9]. In work closely related to ideas in this paper, [3, 4] also described efficient Lyndon-like factorization of a string $\mathbf{x} = \mathbf{x}[1..n]$ into V -words (Definition 11).

Other V -order applications include a variant (V -BWT or V -transform) of the standard lexicographic Burrows-Wheeler transform, introduced in [10], where instances of enhanced data clustering were demonstrated. Also described in [10] was a linear-time algorithm for V -sorting all the conjugates of a string, based on the work in [3, 4] and the linear-time suffix-sorting algorithm of Ko-Aluru [11]. In particular, [10] showed how to compute the V -transform of a V -word \mathbf{x} in $\Theta(n)$ time and space; in addition, inverting the V -transform to recover the original \mathbf{x} was achieved in time $O(n^2 \log k)$, using $O(n + k)$ words of additional storage, where k is the number of sequences of the largest letter in \mathbf{x} .

In this paper, we modify ideas given in [12] that relate Lyndon factorization to suffix arrays and the Burrows-Wheeler transform in order to obtain similar results for the V -order factorization (Section 4). We go on to introduce FM-index type functions in V -order (Section 5), thus raising the possibility of directly applying them to pattern matching for letters in V -order, as done for *lexorder* in [12]. We conclude by applying the backward search technique to substring pattern matching in V -order (Section 6). The differences between *lexorder* and V -order are intriguing: while the latter generally appears trickier to work with,

on the other hand computing suffix arrays in V -order is trivial.

More generally, several recent papers have revealed deep connections between suffix arrays and Lyndon decomposition that are still not well understood. For example, in the first step of his linear-time non-recursive suffix array algorithm, Baier[13] actually computes a sorted version of the *Lyndon array* $LA_{\mathbf{x}}$ of string \mathbf{x} , where $LA_{\mathbf{x}}[i]$ identifies the longest Lyndon word occurring at $\mathbf{x}[i]$. The same structure arises in the recent proof [14] that the number of maximal periodicities (“runs”) in a string is less than its length. For further commentary see [15, 16, 17].

In this paper we begin the process of exploring analogous connections in the context of V -order.

2. Preliminaries

We are given a finite totally ordered set Σ of cardinality $\sigma = |\Sigma|$, called the *alphabet*, whose elements are *characters* (equivalently *letters*). A *string* (equivalently *word*) is a sequence of zero or more characters over Σ . (In our examples, we shall often suppose that $\Sigma = \{a, b, c, \dots\}$, the Roman alphabet in its natural order, or $\Sigma = \{1, 2, 3, \dots\}$, the natural numbers.) A string $\mathbf{x} = x_1x_2 \cdots x_n$ of *length* $|\mathbf{x}| = n$ is represented by $\mathbf{x}[1..n]$, where $x_i = \mathbf{x}[i] \in \Sigma$ for $1 \leq i \leq n$. The set of all nonempty strings over the alphabet Σ is denoted by Σ^+ . The *empty string* of zero length is denoted by ε , with $\Sigma^* = \Sigma^+ \cup \varepsilon$. If $\mathbf{x} = \mathbf{u}\mathbf{v}\mathbf{w}$ for strings $\mathbf{u}, \mathbf{w}, \mathbf{v} \in \Sigma^*$, then \mathbf{u} is a *prefix*, \mathbf{w} a *substring* or *factor*, and \mathbf{v} a *suffix* of \mathbf{x} . We denote the suffix starting at position i , $1 \leq i \leq n$ by S_i . For a given string $\mathbf{x} = \mathbf{x}[1..n]$ and an integer sequence $0 < i_1 < i_2 < \cdots < i_k < n+1$, the string $\mathbf{y} = \mathbf{x}[i_1]\mathbf{x}[i_2] \cdots \mathbf{x}[i_k]$ is said to be a *subsequence* of \mathbf{x} , *proper* if $k < n$.

If $\mathbf{x} = \mathbf{u}^k$ (a concatenation of k copies of \mathbf{u}) for some nonempty string \mathbf{u} and some integer $k > 1$, then \mathbf{x} is said to be a *repetition*; otherwise, \mathbf{x} is *primitive*. A string $\mathbf{y} = R_i(\mathbf{x})$ is the i^{th} *conjugate* (or *rotation*) of $\mathbf{x} = \mathbf{x}[1..n]$ if $\mathbf{y} = \mathbf{x}[i+1..n]\mathbf{x}[1..i]$ for some $0 \leq i < n$ (so that $R_0(\mathbf{x}) = \mathbf{x}$).

Note that all strings are written in mathbold: \mathbf{x}, \mathbf{w} instead of x, w .

We shall be considering various standard arrays associated with strings. The *suffix array* $SA = SA_{\mathbf{x}}[1..n]$ of \mathbf{x} is an array of integers $1..n$ such that for every $i \in 1..n$, $SA[i] = j$ if and only if S_j is the i^{th} suffix in lexorder. For example, the suffixes of $\mathbf{x} = acab$ are $ab, acab, b, cab$ in lexorder, so that $SA_{\mathbf{x}} = (3, 1, 4, 2)$.

The *Burrows-Wheeler transform* $BWT = BWT_{\mathbf{x}}$ of a given string \mathbf{y} [18] is a permutation formed in two equivalent ways from the string $\mathbf{x} = \mathbf{y}\1 , where $\$$ is a *sentinel* letter assumed, for processing convenience, to be smaller than any letter in Σ :

- (1) Form a matrix $M_{\mathbf{x}}^{lex}$ of the rotations of \mathbf{x} sorted in increasing lexorder, then select the final column; or
- (2) Compute the suffix array $SA_{\mathbf{x}}$, then for $i \leftarrow 1$ to $n+1$, set $BWT[i] \leftarrow \mathbf{x}[SA[i]-1]$ whenever $SA[i] > 1$; $BWT[i] \leftarrow \mathbf{x}[n+1] = \$$, otherwise.

Thus, given $\mathbf{x} = acab\$$, we effectively work from

$$\begin{array}{cccccc}
 M_{\mathbf{x}}^{lex} = & \underline{\$} & a & c & a & \mathbf{b} \\
 & \underline{a} & \underline{b} & \underline{\$} & a & c \\
 & \underline{a} & \underline{c} & \underline{a} & \underline{b} & \underline{\$} \\
 & \underline{b} & \underline{\$} & a & c & \mathbf{a} \\
 & \underline{c} & \underline{a} & \underline{b} & \underline{\$} & \mathbf{a}
 \end{array} \tag{1}$$

choosing $BWT_{\mathbf{x}} = bc\$aa$ from the last column. The underlined prefixes in the table correspond to the suffix array entries $SA[1..5] = (5, 3, 1, 4, 2)$.

In practice, method (2) is always used, since it requires only $\Theta(n)$ time, compared to $\Theta(n^2)$ for method (1). As in the example, the BWT often has the useful property that equal letters are grouped together.

The *FM-Index* [19] is defined in terms of two functions:

- $C_{\mathbf{x}}(\ell)$ gives, for each letter $\ell \in \Sigma$ that occurs in \mathbf{x} , the number of letters in \mathbf{x} smaller in lexorder than ℓ . For $\mathbf{x} = acab$, $C_{\mathbf{x}}(a) = 0, C_{\mathbf{x}}(b) = 2, C_{\mathbf{x}}(c) = 3$.

¹The introduction of the $\$$ ensures that the two methods are indeed equivalent.

- $\text{RANK}_{\mathbf{x}}(p, \ell)$ gives the number of occurrences of letter ℓ in the prefix of \mathbf{x} of length p . Thus RANK_{acab} takes the following values:

p/ℓ	<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>
a	1	1	2	2
b	0	0	0	1
c	0	1	1	1

The functions C and RANK can be efficiently computed and are useful in the processing of the BWT, as discussed in Section 5.

Finally, a *Lyndon word* [20] is a primitive word that is minimum in lexorder $<$ over its conjugacy class. The following Lyndon factorization (LF) theorem is fundamental in stringology and will be important in many of our applications:

Theorem 1. [20] *Any word \mathbf{x} can be written uniquely as a product $LF\mathbf{x} = \mathbf{x} = \mathbf{u}_1\mathbf{u}_2 \cdots \mathbf{u}_k$ of Lyndon words, with $(\mathbf{u}_1 \geq \mathbf{u}_2 \geq \cdots \geq \mathbf{u}_k)$.*

For further stringological definitions and theory, see [21, 22].

3. Results on V -Order

We now define a non-lexorder called V -order and describe some of its notable properties. Let $\mathbf{x} = x_1x_2 \cdots x_n$ be a string over Σ . Define $h \in \{1, \dots, n\}$ by $h = 1$ if $x_1 \leq x_2 \leq \cdots \leq x_n$; otherwise, by the unique value such that $x_{h-1} > x_h \leq x_{h+1} \leq x_{h+2} \leq \cdots \leq x_n$. Let $\mathbf{x}^* = x_1x_2 \cdots x_{h-1}x_{h+1} \cdots x_n$, where the star $*$ indicates deletion of x_h . Write $\mathbf{x}^{s*} = (\dots(\mathbf{x}^*)^* \dots)^*$ with $s \geq 0$ stars. Let $g = \max\{x_1, x_2, \dots, x_n\}$, and let k be the number of occurrences of g in \mathbf{x} . Then the sequence $\mathbf{x}, \mathbf{x}^*, \mathbf{x}^{2*}, \dots$ ends $g^k, \dots, g^1, g^0 = \varepsilon$. From all strings \mathbf{x} over Σ we form the *star tree*, where each string \mathbf{x} labels a vertex and there is a directed edge upward from \mathbf{x} to \mathbf{x}^* , with the empty string ε as the root.

Definition 2 (V -order [2]). *We define V -order \prec between distinct strings \mathbf{x} ,*

\mathbf{y} , where \mathbf{x}, \mathbf{y} do not contain $\$$ ². First $\mathbf{x} \prec \mathbf{y}$ if in the star tree \mathbf{x} is in the path $\mathbf{y}, \mathbf{y}^*, \mathbf{y}^{2*}, \dots, \varepsilon$. If \mathbf{x}, \mathbf{y} are not in a path, there exist smallest s, t such that $\mathbf{x}^{(s+1)*} = \mathbf{y}^{(t+1)*}$. Let $\mathbf{s} = \mathbf{x}^{s*}$ and $\mathbf{t} = \mathbf{y}^{t*}$; then $\mathbf{s} \neq \mathbf{t}$ but $|\mathbf{s}| = |\mathbf{t}| = m$ say. Let $j \in [1..m]$ be the greatest integer such that $\mathbf{s}[j] \neq \mathbf{t}[j]$. If $\mathbf{s}[j] < \mathbf{t}[j]$ in Σ then $\mathbf{x} \prec \mathbf{y}$; otherwise, $\mathbf{y} \prec \mathbf{x}$. Clearly \prec is a total order on all strings in Σ^* .

Example 3. [Star tree] We illustrate the star tree for the V -order comparison of the words $\mathbf{x} = \text{lexorder}$ and $\mathbf{y} = \text{matrix}$. The subscript h indicates the V letter to be deleted (defined above as $x_{h-1} > x_h \leq x_{h+1} \leq x_{h+2} \leq \dots \leq x_n$). The circled letters are those compared in alphabetic order (defined above as $\mathbf{s}[j] \neq \mathbf{t}[j]$).

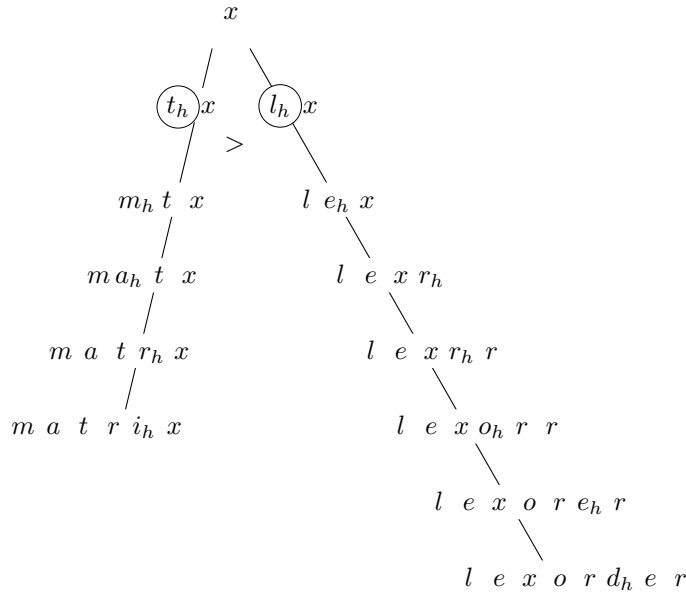


Figure 1: $\text{lexorder} \prec \text{matrix}$

An instance where \mathbf{x} and \mathbf{y} occur in the same path in the star tree is given by $\mathbf{x} = 42527$, $\mathbf{y} = 57$, yielding $\mathbf{x}^* = 4257$, $\mathbf{x}^{2*} = 457$, $\mathbf{x}^{3*} = 57$ and so $57 \prec 42527$.

²It turns out that, in V -order, it is useful to require, not only that $\$$ be the smallest letter, but moreover that $\$ \prec \varepsilon$. To avoid this situation, we transform strings containing $\$$ into equivalent $\$$ -free strings to which Definition 2 applies. See Definition 17.

Definition 4 (*V-form* [2, 23, 3, 4]). The *V-form* of a string \mathbf{x} is defined as

$$V_k(\mathbf{x}) = \mathbf{x} = \mathbf{x}_0 g \mathbf{x}_1 g \cdots \mathbf{x}_{k-1} g \mathbf{x}_k$$

for strings \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. Note that, in this representation, any \mathbf{x}_i may be the empty string ε . For clarity, when more than one string is involved, we use the notation $\mathcal{L}\mathbf{x} = g$, $\mathcal{C}\mathbf{x} = k$.

Lemma 5. [2, 23, 3, 4] Suppose we are given distinct strings \mathbf{x} and \mathbf{y} with corresponding *V-forms* as follows:

$$\begin{aligned} \mathbf{x} &= \mathbf{x}_0 \mathcal{L}\mathbf{x} \mathbf{x}_1 \mathcal{L}\mathbf{x} \mathbf{x}_2 \cdots \mathbf{x}_{j-1} \mathcal{L}\mathbf{x} \mathbf{x}_j, \\ \mathbf{y} &= \mathbf{y}_0 \mathcal{L}\mathbf{y} \mathbf{y}_1 \mathcal{L}\mathbf{y} \mathbf{y}_2 \cdots \mathbf{y}_{k-1} \mathcal{L}\mathbf{y} \mathbf{y}_k, \end{aligned}$$

where $j = \mathcal{C}\mathbf{x}$, $k = \mathcal{C}\mathbf{y}$.

Let $h \in 0.. \max(j, k)$ be the least integer such that $\mathbf{x}_h \neq \mathbf{y}_h$. Then $\mathbf{x} \prec \mathbf{y}$ if and only if one of the following conditions holds:

- (C1) $\mathcal{L}\mathbf{x} < \mathcal{L}\mathbf{y}$
- (C2) $\mathcal{L}\mathbf{x} = \mathcal{L}\mathbf{y}$ and $\mathcal{C}\mathbf{x} < \mathcal{C}\mathbf{y}$
- (C3) $\mathcal{L}\mathbf{x} = \mathcal{L}\mathbf{y}$, $\mathcal{C}\mathbf{x} = \mathcal{C}\mathbf{y}$ and $\mathbf{x}_h \prec \mathbf{y}_h$.

Observe the recursive nature of determining \prec in (C3); that is, each substrings pair \mathbf{x}_h , \mathbf{y}_h can likewise be decomposed into *V-forms*.

Example 6. [C1-C3] The following examples illustrate the three conditions in Lemma 5:

- Using C1, if $\mathbf{x} = 345$ and $\mathbf{y} = 456$, we have $\mathbf{x} \prec \mathbf{y}$ as $\mathcal{L}\mathbf{x} = 5 < 6 = \mathcal{L}\mathbf{y}$.
- Using C2, if $\mathbf{x} = 1818181$ and $\mathbf{y} = 78787$, we have $\mathbf{x} \succ \mathbf{y}$ as $\mathcal{L}\mathbf{x} = \mathcal{L}\mathbf{y} = 8$ and $\mathcal{C}\mathbf{x} = 3 > 2 = \mathcal{C}\mathbf{y}$.

- Using C3, if $\mathbf{x} = 9616921$ and $\mathbf{y} = 9616912$, we have $\mathbf{x} \prec \mathbf{y}$ because $\mathbf{x}_0 = \mathbf{y}_0 = \varepsilon$, $\mathbf{x}_1 = \mathbf{y}_1 = 616$ while $\mathbf{x}_2 = 21 \prec 12 = \mathbf{y}_2$, since the V -form \mathbf{x}_0 prefix of 21 is ε while for 12 it is 1 and $\varepsilon \prec 1$. It is instructive to extend this example to $\mathbf{x}' = 961692198888$ and $\mathbf{y}' = 961691294$, where $\mathbf{x}' \prec \mathbf{y}'$ since the order is still determined by $\mathbf{x}'_2 = 21$ and $\mathbf{y}'_2 = 12$.

Lemma 7. [3, 4] For given strings \mathbf{x} and \mathbf{y} , if \mathbf{y} is a proper subsequence of \mathbf{x} , then $\mathbf{y} \prec \mathbf{x}$.

For example, using this lemma we immediately have $9374 \prec 93748336$, $937336 \prec 93748336$, and $9786 \prec 93748336$.

The remarkable Lemma 7 has many consequences, not least the trivial ordering of suffixes in V -order discussed in Section 4.

Lemma 8. [6, 8] For any strings $\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}$: $\mathbf{x} \prec \mathbf{y} \Leftrightarrow \mathbf{uxv} \prec \mathbf{uyv}$.

By Lemma 8, a comparison of two strings in V -order can ignore equal prefixes (or suffixes). Therefore, on this basis, Algorithm COMPARE(\mathbf{x}, \mathbf{y}), first presented in [8] and later upgraded in [9], ignores common prefixes of strings \mathbf{x} and \mathbf{y} while comparing them. COMPARE also makes use of the following:

Lemma 9. [8] For any two strings \mathbf{x}, \mathbf{y} with $\mathbf{x} \prec \mathbf{y}$ and any two letters λ, μ such that $\mathbf{y} \prec \mathbf{x}\lambda$:

- (i) if $\lambda \leq \mu$, then $\mathbf{x}\lambda \prec \mathbf{y}\mu$;
- (ii) if $\lambda > \mu$, then $\mathbf{y}\mu \prec \mathbf{x}\lambda$.

In contrast to comments made in the Conclusion of [4], Lemma 9 implies that V -order comparison can be conducted in a positional manner just like lexorder comparison. In fact, COMPARE(\mathbf{x}, \mathbf{y}) is an on-line algorithm that requires only a one-character window!

The improved V -order comparison algorithm given in [9], called COMPARE-Sensitive, was obtained by simple modifications of COMPARE. COMPARE-Sensitive considers the structure of the input strings in order to achieve improved

running times. A probabilistic argument shows that COMPARE-Sensitive runs faster than COMPARE in almost all cases of interest, at least twice as fast according to the experiments described in [9].

Lemma 10. [3, 4, 5, 7, 8, 9] *V-order comparison of given strings \mathbf{x} and \mathbf{y} requires linear time and constant space.*

We now introduce the V -order equivalent of the lexorder Lyndon word:

Definition 11 (V -Word [23]). *A string \mathbf{x} over an ordered alphabet Σ is said to be a V -word if it is the unique minimum in V -order \prec in the conjugacy class of \mathbf{x} .*

Thus, like a Lyndon word, a V -word is necessarily primitive.

Example 12. [\prec] *We can apply Definition 2, equivalently the structure given by Lemma 5, to conclude that*

$$7173 \prec 7371 \prec 1737 \prec 3717,$$

so that 7173 is a V -word, while on the other hand 1737 is a Lyndon word. Similarly, 71727174 and 818382 are V -words, while conjugates 17271747 and 183828 are Lyndon words.

Now consider two V -words \mathbf{x}_1 and \mathbf{x}_2 on an alphabet Σ . It may be that $\mathbf{x}_1\mathbf{x}_2$ is also a V -word, a fact we record by writing $\mathbf{x}_1 \prec_V \mathbf{x}_2$; if not, then we write $\mathbf{x}_1 \geq_V \mathbf{x}_2$. Thus, corresponding to the Lyndon factorization into Lyndon words using \geq (Theorem 1), we arrive at a V -order factorization expressed in terms of an ordering \geq_V of V -words:

Lemma 13. [3, 4] *Using only linear time and space (see Algorithm VF in [4]³), a string \mathbf{x} can be factored uniquely into V -words $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \cdots \mathbf{x}_m$, where by definition $\mathbf{x}_1 \geq_V \mathbf{x}_2 \geq_V \cdots \geq_V \mathbf{x}_m$.*

³Although no claim is made in this reference, VF is an on-line algorithm, since it does no backtracking and outputs the V -word factors in order from left to right.

Example 14. [\geq_V] For $\mathbf{x} = 33132421$, the Lyndon decomposition is $3 \geq 3 \geq 13242 \geq 1$, while the V -order factorization identifies nonextendible V -words 33132 and 421 with $33132 \geq_V 421$. (Note however that $33132 \prec 421$! See [24] for more background on this phenomenon.) Similarly, from Example 12, the string

$$\mathbf{x} = \mathbf{uvw} = (7173)(71727174)(818382)$$

has the unique V -order factorization $\mathbf{u} \geq_V \mathbf{v} \geq_V \mathbf{w}$, even though $\mathbf{u} \prec \mathbf{v} \prec \mathbf{w}$.

As we now discover, it will also be useful to order strings \mathbf{x} , \mathbf{y} based on a lexicographic approach to their V -order factorizations. We call this ordering, denoted $\prec_{LEX(V)}$, *lex-extension* order, expressed here with respect to V -order factorization, while another choice of factorization is given by the unique V -form of a string.

Definition 15 (Lex-extension order [23, 4, 10]). *Suppose that the V -order factorization of distinct strings $\mathbf{x}, \mathbf{y} \in \Sigma^+$ yields nonempty factors*

$$\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \cdots \mathbf{x}_p, \quad \mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \cdots \mathbf{y}_q.$$

If either

- (1) \mathbf{x} is a proper prefix of \mathbf{y} (that is, $\mathbf{x}_i = \mathbf{y}_i$ for $1 \leq i \leq p < q$) or else
- (2) for some $i \in 1..\min(p, q)$, $\mathbf{x}_j = \mathbf{y}_j$ for every $1 \leq j < i$, and $\mathbf{x}_i \prec \mathbf{y}_i$,

then we write $\mathbf{x} \prec_{LEX(V)} \mathbf{y}$. Otherwise, $\mathbf{x} \succ_{LEX(V)} \mathbf{y}$ ($\mathbf{y} \prec_{LEX(V)} \mathbf{x}$).

Intuitively, in $\prec_{LEX(V)}$ ordering, we first factor the two strings to be compared into V -factors. Then we treat each V -factor as a letter, and compare strings by comparing the V -factors at the same position. Consider Example 16, where $\mathbf{y} \prec \mathbf{x}$, but $\mathbf{x} \prec_{LEX(V)} \mathbf{y}$.

Example 16. [$\prec_{LEX(V)}$] *Given V -order factorizations*

$$\mathbf{x} = \mathbf{x}_1\mathbf{x}_2\mathbf{x}_3 = (33132)(422)(5), \quad \mathbf{y} = \mathbf{y}_1\mathbf{y}_2 = (33132)(413),$$

we find $\mathbf{x}_1 = \mathbf{y}_1$, but since $\mathbf{x}_2 \prec \mathbf{y}_2$, therefore $\mathbf{x} \prec_{LEX(V)} \mathbf{y}$ (even though $\mathbf{x} \succ \mathbf{y}$).

4. Computing the Extended Suffix Array in V -Order

In this section we discuss the V -order version, V -BWT, of the Burrows-Wheeler transform [10]. As in the lexorder case, described in Section 2, it is convenient to append a special sentinel letter $\$$, less in lexorder than every letter in Σ , to each string being processed. However, it turns out that, in V -order, it could become necessary also to require that $\$$ be less than the empty string ε . For example, this requirement would arise if, while executing a V -order comparison of strings containing $\$$, an element \mathbf{x}_i corresponding to $\$$ in the V -form representation (Definition 4) were empty.

In order to avoid this situation, we introduce Definition 17, which transforms \mathbf{x} , \mathbf{y} , each containing at most a single $\$$, into $\$$ -free equivalent strings \mathbf{x}' , \mathbf{y}' such that $\mathbf{x}' \prec \mathbf{y}' \Rightarrow \mathbf{x} \prec \mathbf{y}$. This is achieved by introducing new letters α and $\alpha' > \alpha$, both less than any letter in Σ , but *not* less than ε , to replace $\$$ in \mathbf{x} and \mathbf{y} .

Definition 17. *[Transforming strings containing $\$$ to equivalent $\$$ -free strings.] Given distinct strings \mathbf{x} , \mathbf{y} , each containing at most one occurrence of $\$$, we construct new strings \mathbf{x}' , \mathbf{y}' as defined below, by introducing new letters α and $\alpha' > \alpha$, both less than any letter in Σ , to replace $\$$; that is, we replace $\Sigma \cup \{\$\}$ by $\Sigma' = \Sigma \cup \{\alpha, \alpha'\}$. The following rules are applied:*

1. If $\mathbf{x}[r] = \mathbf{y}[t] = \$$, $1 \leq r < t \leq n$:
 - (a) $\mathbf{x}' \leftarrow \mathbf{x}[1..r-1]\alpha\mathbf{x}[r+1..|\mathbf{x}|]$ (α replaces $\$$),
 $\mathbf{y}' \leftarrow \mathbf{y}[1..r-1]\alpha'\mathbf{y}[r..|\mathbf{y}|]$ (insert α');
 - (b) $\mathbf{y}' \leftarrow \mathbf{y}'[1..t]\alpha\mathbf{y}'[t+2..|\mathbf{y}'|]$ (α replaces $\$$),
if $|\mathbf{x}'| \geq t$ then $\mathbf{x}' \leftarrow \mathbf{x}'[1..t]\alpha'\mathbf{x}'[t+1..|\mathbf{x}'|]$ (insert α').
2. If $\mathbf{x}[r] = \$$, $1 \leq r \leq |\mathbf{x}|$, and $\forall i, 1 \leq i \leq |\mathbf{y}| : \mathbf{y}[i] \neq \$$:
 - $\mathbf{x}' \leftarrow \mathbf{x}[1..r-1]\alpha\mathbf{x}[r+1..|\mathbf{x}|]$, (α replaces $\$$),
if $|\mathbf{y}| \geq r-1$ then $\mathbf{y}' \leftarrow \mathbf{y}[1..r-1]\alpha'\mathbf{y}[r..|\mathbf{y}|]$ (insert α').

3. If $\mathbf{x}[r] = \mathbf{y}[r] = \$$, $1 \leq r \leq |\mathbf{x}|, |\mathbf{y}|$:
 $\mathbf{x}' \leftarrow \mathbf{x}[1..r-1]\mathbf{x}[r+1..|\mathbf{x}|]$, $\mathbf{y}' \leftarrow \mathbf{y}[1..r-1]\mathbf{y}[r+1..|\mathbf{y}|]$
(remove $\$$ from both \mathbf{x} and \mathbf{y}).

Then $\mathbf{x}' \prec \mathbf{y}' \Rightarrow \mathbf{x} \prec \mathbf{y}$.

Since the ordering of letters between the strings \mathbf{x}, \mathbf{y} and \mathbf{x}', \mathbf{y}' is preserved, and because $\alpha < \alpha'$, both less than every element in Σ , the ordering between strings \mathbf{x}, \mathbf{y} and strings \mathbf{x}', \mathbf{y}' , respectively, is preserved. Thus we extend the application of Definition 2 to comparison of distinct strings \mathbf{x}, \mathbf{y} , each containing at most one occurrence of $\$ \prec \varepsilon$. This also ensures the inequalities between strings are well defined in all cases. Furthermore, since strings \mathbf{x}' and \mathbf{y}' are strings over an alphabet Σ' which does not include $\$$, and V -order (Definition 2) is defined for arbitrary alphabets (excluding $\$$), all known results on V -order, including all the lemmas presented in this paper, apply immediately to \mathbf{x}' and \mathbf{y}' . In particular, the transitivity property for V -order on Σ naturally holds on Σ' also.

We apply Definition 17 to the example strings $\mathbf{x} = \$acab$ and $\mathbf{y} = acab\$$ introduced above. From Rule 1(a), we get $\mathbf{x}' = \alpha acab$ and $\mathbf{y}' = \alpha' acab\$$; next, applying 1(b), we find $\mathbf{x}' = \alpha acaba'$ and $\mathbf{y}' = \alpha' acaba\alpha$. Then using the star tree construction of Definition 2 to compare \mathbf{x}' and \mathbf{y}' , we find as above that $s = \mathbf{x}'^{3*} = \alpha ac$ and $t = \mathbf{y}'^{3*} = \alpha' ac$, yielding $\mathbf{x}' \prec \mathbf{y}'$. Hence, by Definition 17, $\mathbf{x} \prec \mathbf{y}$.

Applying this methodology to $\mathbf{x} = acab\$$, we can sort all its conjugates into V -order – $cab\$a \prec \$acab \prec acab\$ \prec b\$aca \prec ab\$ac$ – yielding the V -order Burrows-Wheeler matrix $M_{\mathbf{x}}^{\vee}$ instead of the lexorder version $M_{\mathbf{x}}^{lex}$ given in (1):

$$\begin{aligned}
M_{\mathbf{x}}^{\vee} = & \begin{array}{ccccc} c & a & b & \$ & a \\ \$ & a & c & a & b \\ a & c & a & b & \$ \\ b & \$ & a & c & a \\ a & b & \$ & a & c \end{array} & (2)
\end{aligned}$$

Finally we form $V\text{-BWT}_{\mathbf{x}} = ab\ac from the last column L of $M_{\mathbf{x}}^V$.

We go on now to introduce the extended suffix array and its V -order implementation, making use of the V -BWT in a manner suggested recently by the lexorder methodology of Mantaci *et al.* [12].

We begin by showing how to translate the results of [12] from lexorder into V -order. These authors describe a strategy for obtaining the suffix array $\text{SA}_{\mathbf{x}}$ of a string \mathbf{x} from its Lyndon factorization $\text{LF}_{\mathbf{x}}$. To do the equivalent calculation in V -order — $\text{SA}_{\mathbf{x}}$ from the V -order factorization of \mathbf{x} (Lemma 13) — is not however of interest, because, as remarked earlier, sorting the suffixes of \mathbf{x} in V -order is trivial! By virtue of Lemma 7, the V -order of the suffixes is

$$\mathbf{x}[n] \prec \mathbf{x}[n-1..n] \prec \cdots \prec \mathbf{x}[1..n],$$

so that $\text{SA}_{\mathbf{x}}[i] = n - i + 1$ for all \mathbf{x} .

Thus, in the V -order version of this problem, we replace the suffixes with the “extended suffixes” (conjugates of \mathbf{x}) employed to determine the V -BWT. Then the i -th *extended suffix* is defined by

$$ES_i \equiv \mathbf{x}[i..n]\mathbf{x}[1..i-1]. \quad (3)$$

The following example, from [10], shows that the ordering of the extended suffixes is no longer straightforward:

Example 18. For $\mathbf{x} = 9191919293$, consider suffixes $S_5 = 919293$ and $S_3 = 91919293$, with $S_5 \prec S_3$ by Lemma 7. However, extending these suffixes and applying Lemma 5 (C3) yields the opposite order:

$$\underline{9191929391} \prec \underline{9192939191} \iff ES_3 \prec ES_5.$$

Thus we show here how to use the V -order factorization of \mathbf{x} (in particular, the on-line algorithm VF mentioned in Lemma 13) to generate the Burrows-Wheeler matrix $M_{\mathbf{x}}^V$ and thus the corresponding $V\text{-BWT}_{\mathbf{x}}$. We call the cor-

responding suffix array $\text{EVSA} = \text{EVSA}_{\mathbf{x}}$, equivalent to the rotations of \mathbf{x} as ordered in $M_{\mathbf{x}}^{\vee}$. For example, in Equation (2), for $\mathbf{x}' = \text{acab}\$,$ $\text{EVSA}_{\mathbf{x}'} = (2, 5, 1, 4, 3)$ and $M_{\mathbf{x}'}^{\vee}$ can be derived from each other, each yielding the conclusion that $V\text{-BWT}_{\mathbf{x}'} = \text{ab}\ac .

We remark that the difference between Lyndon and V -order factorization of \mathbf{x} [3, 4] provides more options for efficient string processing. Examples 12 and 14 give an idea of the variation that can occur between lexorder and V -order.

Now we turn to the idea of “compatibility” of sorted suffixes as introduced in [12]. Let $\mathbf{x} = \mathbf{x}[1..n]$ be a string and $\mathbf{u} = \mathbf{x}[i..j]$, $1 \leq i \leq j \leq n$, a substring of \mathbf{x} . Then the sorting of suffixes $\mathbf{s}_p = \mathbf{x}[p..j]$, $\mathbf{s}_q = \mathbf{x}[q..j]$ of \mathbf{u} is *compatible* with the sorting of the corresponding suffixes $\mathbf{x}[p..n]$, $\mathbf{x}[q..n]$ of \mathbf{x} if these two (p, q) pairs have the same order in both \mathbf{u} and \mathbf{x} . In lexorder, compatibility of arbitrarily chosen \mathbf{u} and \mathbf{v} does not always hold [12], but does hold when they are substrings of Lyndon factors in the Lyndon factorization of \mathbf{x} . However, in V -order, by Lemma 7, compatibility holds for every choice of p, q . Moreover, the shorter suffix is always lesser, thus allowing comparison in terms of indexes:

Lemma 19. *Let $\mathbf{x} \in \Sigma^+$ and \mathbf{u} be a substring of \mathbf{x} with \mathbf{s}_1 a suffix of \mathbf{u} . If \mathbf{s}_2 is a proper suffix of \mathbf{s}_1 then $\mathbf{s}_2 \prec \mathbf{s}_1$ with respect to both \mathbf{u} and \mathbf{x} .*

Proof. Consider the suffixes $\mathbf{s}_1\mathbf{t}_1$ and $\mathbf{s}_2\mathbf{t}_2$ of \mathbf{x} for possibly empty $\mathbf{t}_1, \mathbf{t}_2$. Applying Lemma 7, we find that $\mathbf{s}_2 \prec \mathbf{s}_1$ with respect to \mathbf{u} and $\mathbf{s}_2\mathbf{t}_2 \prec \mathbf{s}_1\mathbf{t}_1$ with respect to \mathbf{x} , as required. \square

As we have seen, Lemma 19 is not sufficient for extended suffixes, but since each conjugate has the same number of maximum g 's, condition (C3) implicitly applies. However, as a special case of this lemma, we obtain the following V -order analogy of the result (Theorem 3.2) given in [12] for Lyndon decomposition:

Corollary 20. *Let $\mathbf{x} \in \Sigma^+$ with V -order factorization $\mathbf{x} = \mathbf{v}_1 \cdots \mathbf{v}_k$, and let $\mathbf{u} = \mathbf{v}_i \cdots \mathbf{v}_j$, for $1 \leq i \leq j \leq k$. Then the sorting of the suffixes of \mathbf{u} is compatible with the sorting of the corresponding suffixes of \mathbf{x} .*

The intricate differences between lexorder and V -order, along with the superficial nature of the V -order suffix array, lead to the following lemma that aims toward incrementally building the extended V -order suffix array, EVSA, in a manner similar to that employed in [12] for building SA. Thus, for $k = 2$, the next lemma partly describes the structure of $M_{\mathbf{u}}^{\mathcal{V}}$, hence of $\text{EVSA}_{\mathbf{u}}$:

Lemma 21. *Let $\mathbf{u} = \mathbf{v}_1\mathbf{v}_2\$ \in \Sigma^+$, where $\mathbf{v}_1, \mathbf{v}_2$ are V -words with $\mathbf{v}_1 \geq_{\mathcal{V}} \mathbf{v}_2$ and $\mathcal{L}\mathbf{v}_1 = \mathcal{L}\mathbf{v}_2 = \mathcal{L}$. Let $R_0, R_1, \dots, R_{\mathcal{C}\mathbf{v}_2}$ be the $\mathcal{C}\mathbf{v}_2 + 1$ consecutive rotations of $\mathbf{v}_2\$$ occurring in the V -order Burrows-Wheeler matrix $M_{\mathbf{v}_2\$}^{\mathcal{V}}$ that start with $\$$ or \mathcal{L} . Then in $M_{\mathbf{u}}^{\mathcal{V}}$:*

- (a) *for every $j \in 0..\mathcal{C}\mathbf{v}_2$, there exists a unique row E_r , with (nonempty) prefix P_r and (possibly empty) suffix S_r , such that $P_r S_r = R_j$ and $S_r P_r = \mathbf{v}_2\$$;*
- (b) *E_0 has prefix $\$$ and the remaining $\mathcal{C}\mathbf{v}_2$ entries R_j identified in (a) occur in rows $E_j, 1 \leq j \leq \mathcal{C}\mathbf{v}_2$, in ascending order of j .*

Proof.

For both cases (a) and (b), appending $\$$ to \mathbf{v}_2 ensures that working with suffix arrays is equivalent to working with $M_{\mathbf{u}}^{\mathcal{V}}$ and similarly for $\text{EVSA}_{\mathbf{u}}$ (see commentaries on the sentinel $\$$ in Section 2 and prior to Equation 2).

- (a) This follows from Corollary 20 and the fact that $M_{\mathbf{u}}^{\mathcal{V}}$ includes all the rotations of \mathbf{u} .
- (b) Since the first row of $M_{\mathbf{u}}^{\mathcal{V}}$ must be least in V -order, and $\$ \prec \varepsilon$ and \mathbf{v}_1 being a V -word beginning with \mathcal{L} , it has prefix $\$$ — this unique row is identified as E_0 .

For the sake of rigour, consider converting to $\$$ -free strings. Since \mathbf{u} contains exactly one $\$$, then it can only appear in the first position in any row of $M_{\mathbf{u}}^{\mathcal{V}}$, once. Also since \mathbf{v}_1 is a V -word it must begin with \mathcal{L} . So let r be the row of $M_{\mathbf{u}}^{\mathcal{V}}$ starting $\$\mathcal{L}$. Comparing row r with other rows in $M_{\mathbf{u}}^{\mathcal{V}}$, and applying Definition 17 part 1, row r now starts with α and every other row starts with α' , where $\alpha < \alpha'$. Consider the \mathbf{x}_0 prefix for the V -form

of each row in $M_{\mathbf{u}}^{\mathcal{V}}$, then for row r we have $\mathbf{x}_0 = \alpha$ and hence it is unique out of all the \mathbf{x}_0 substrings. Furthermore, each of the other \mathbf{x}_0 substrings must contain some letter greater than α , so that Lemma 5 (C1) applies, thus establishing that $r = E_0$. We now return to working with $\$$.

Consider those suffixes P_r in the EVSA suffix array of $\mathbf{v}_2\$$ which start with $\$$ or \mathcal{L} , and suppose that they are V -ordered so that their subscripts are in ascending order of j . The rows E_r in $M_{\mathbf{u}}^{\mathcal{V}}$ are then identified as those of the form $P_r \mathbf{v}_1 S_r$.

Suppose that $P_k \prec P_l$ which is determined by applying lex-extension Definition 15 with the factorizations given by the V -forms of P_k and P_l . Since $\mathbf{v}_2\$$ only contains one $\$$, then part 1 of this definition does not apply, for if P_k was a proper prefix of P_l then P_l would contain two $\$$'s. From part 2 there is a pair of substrings $\mathbf{x}_i, \mathbf{y}_i$ with minimal i where \mathbf{x}_i is in P_k and \mathbf{y}_i is in P_l with $\mathbf{x}_i \prec \mathbf{y}_i$. Then $P_k \prec P_l$ implies $P_k \mathbf{v}_1 S_k \prec P_l \mathbf{v}_1 S_l$. In other words, the V -order of $P_k \mathbf{v}_1 S_k$ and $P_l \mathbf{v}_1 S_l$ is determined before the start of \mathbf{v}_1 by the order of P_k and P_l which both end in $\$$. Hence the order of the P_r in $M_{\mathbf{v}_2\$}^{\mathcal{V}}$ is preserved in $M_{\mathbf{u}}^{\mathcal{V}}$.

Overall, we are essentially merging two lists of suffixes each ordered in lex-extension order into a further lex-extension ordered list of suffixes, hence maintaining the lists' original orders. \square

We illustrate Lemma 21 with the following example:

Example 22. Let V -words $\mathbf{v}_1 = 83$ and $\mathbf{v}_2 = 88182$, where $\mathbf{v}_1 \geq_{\mathcal{V}} \mathbf{v}_2$, $\mathcal{L}\mathbf{v}_1 = \mathcal{L}\mathbf{v}_2 = 8$, and $\mathbf{u} = 8388182\$$. The resulting matrices are given below, where the rotations R_j (starting with $\$$ or $\mathcal{L} = 8$) are indicated in $M_{\mathbf{v}_2\$}^{\mathcal{V}}$ and the $P_i S_i$ are

underlined in $M_{\mathbf{u}}^{\mathcal{V}}$.

$$\begin{aligned}
M_{\mathbf{v}_2\$}^{\mathcal{V}} = & R_0 : \$ 8 8 1 8 2 \\
& R_1 : 8 8 1 8 2 \$ \\
& R_2 : 8 1 8 2 \$ 8 \\
& R_3 : 8 2 \$ 8 8 1 \\
& \quad 1 8 2 \$ 8 8 \\
& \quad 2 \$ 8 8 1 8
\end{aligned} \tag{4}$$

$$\begin{aligned}
M_{\mathbf{u}}^{\mathcal{V}} = & \underline{\$} 8 3 \underline{8} \underline{8} \underline{1} \underline{8} \underline{2} \\
& \underline{8} \underline{8} \underline{1} \underline{8} \underline{2} \underline{\$} 8 3 \\
& \underline{8} \underline{1} \underline{8} \underline{2} \underline{\$} 8 3 \underline{8} \\
& \underline{8} \underline{2} \underline{\$} 8 3 \underline{8} \underline{8} \underline{1} \\
& 8 3 8 8 1 8 2 \$ \\
& 1 8 2 \$ 8 3 8 8 \\
& 2 \$ 8 3 8 8 1 8 \\
& 3 8 8 1 8 2 \$ 8
\end{aligned} \tag{5}$$

Equipped with Corollary 20 and Lemma 21, we can modify the clever incremental suffix-sorting/BWT strategy of [12] to construct the extended V -order suffix array, EVSA, of $\mathbf{x} = \mathbf{x}[1..n]$. We outline the steps:

Step 1: Identify the first factor \mathbf{v}_1 in the V -order factorization $\mathbf{v}_1 \geq_{\mathcal{V}} \cdots \geq_{\mathcal{V}} \mathbf{v}_k$ of \mathbf{x} in linear time [3, 4]; we assume the maximum letter in each factor is \mathcal{L} .

Step 2: Compute the lex-extension order suffix array $\text{EVSA}_{\mathbf{v}_1\$}$ of $\mathbf{v}_1\$$ in linear time [10]. In practice we will only need the suffixes starting with \mathcal{L} in our pattern matching application in Section 6.

Step 3: Extract $\text{BWT}_{\mathbf{v}_1\$}$ from $\text{EVSA}_{\mathbf{v}_1\$}$.

Step 4: For factor $\mathbf{v}_i[1..t]$, $1 < i \leq k$, insert each suffix $\mathbf{v}_i[t]$, $\mathbf{v}_i[t-1..t]$, \dots , $\mathbf{v}_i[1..t]$ into the current EVSA in its V -order — so for $i = 2$, one by one the

suffixes of $\mathbf{v}_2\$$ are inserted into $\text{EVSA}_{\mathbf{v}_1\$}$ giving the EVSA for $\mathbf{v}_1\mathbf{v}_2\$$. Further, as each \mathbf{v}_i is processed we can extract $\text{BWT}_{\mathbf{v}_1\dots\mathbf{v}_i\$}$: for a suffix $\mathbf{x}[i..j]$ the BWT letter is $\mathbf{x}[i - 1]$.

Note: After the above steps all the V -word factors are incrementally processed and $V\text{-BWT}_{\mathbf{x}}$ computed. To reduce complexity, the method avoids merging the suffix arrays of both the extended suffix array computed to date and that of the current factor being processed.

We now describe suffix insertion (Step 4) in more detail. For this we apply Lemma 3.16 in [4] which states that, unlike Lyndon words, the order of the set \mathcal{V} of V -words is in some cases the same as V -order while in other cases it is reversed. This leads to two cases for Step 4: first, the next factor to be processed has a larger maximum letter \mathcal{L} than those in all factors processed so far, and second, the maximum letters are the same. However, the following example shows that the method doesn't work in the first case of distinct maximal letters — hence they are restricted to be equal in Lemma 21.

Example 23. Let $\mathbf{u} = \mathbf{v}_1 \geq_{\mathcal{V}} \mathbf{v}_2$ with $\mathbf{v}_1 = 321$ and $\mathbf{v}_2 = 5152$. The ordered conjugates of \mathbf{v}_1 are: $321 < 132 < 213$, but note that the order of these conjugates is changed after processing \mathbf{v}_2 . We list the ordered conjugates of $\mathbf{u} = 3215152$ below with complete factors shown in square brackets and the conjugates of 321 underlined.

$5152[321]$
 $52[321]51$
 $\underline{1}[5152]32$
 $152[321]5$
 $\underline{21}[5152]3$
 $\underline{321}[5152]$
 $2[321]515$

On the other hand, as we have seen, we are constrained to extended suffixes

relating to a conjugacy class where each conjugate has the same maximal letter \mathcal{L} which occurs with the same frequency. In particular, Lemma 21 and condition (C3) ensure that the order of an existing suffix array is preserved during this iterative processing of factors for those suffixes starting with \mathcal{L} ; therefore we now go on to deal with this case.

Suppose a prefix \mathbf{p} of p factors of \mathbf{x} has been processed resulting in $\text{EVSA}(\mathbf{p})$ with an associated set of ordered conjugates, where all conjugates have the same maximum letter. The task is now to insert the suffixes of factor \mathbf{v}_{p+1} into $\text{EVSA}(\mathbf{p})$. Applying properties established in [10], if $g = \mathcal{L}\mathbf{w}$ and $k = \mathcal{C}\mathbf{w}$ for a string \mathbf{w} (which does not have the sentinel $\$$ appended), then in V -order the first k conjugates of \mathbf{w} start with g . Processing will be as follows:

- First all existing suffixes computed prior to processing \mathbf{v}_{p+1} will have \mathbf{v}_{p+1} appended to them.
- The technique for calculating the index for inserting each proper suffix of \mathbf{v}_{p+1} (increasing one V -form \mathbf{x}_i type substring at a time from right to left) into EVSA is given with the FM-index in Section 5.
- Inserting the last improper suffix, \mathbf{v}_{p+1} , is straightforward. We establish that, analogous to the lexorder case in [12], this suffix is least in V -order and can be entered directly into the extended suffix array. Consider \mathbf{v}_1 : since it is a V -word it is least in $\text{EVSA}(\mathbf{v}_1)$ and likewise for \mathbf{v}_2 . If we suppose that $\mathbf{v}_1 \neq \mathbf{v}_2$, then $\mathbf{v}_2\mathbf{v}_1 \in \mathcal{V}$ (see the generalized form of properties of V -words described by Theorem 2.7(3) in [4]), and $\mathbf{v}_2\mathbf{v}_1$ is least in $\text{EVSA}(\mathbf{v}_2\mathbf{v}_1)$, and further, $\mathbf{v}_2 \prec \mathbf{v}_2\mathbf{v}_1$. Similarly, if $\mathbf{v}_1 = \mathbf{v}_2$, then $\mathbf{v}_2 \prec \mathbf{v}_2\mathbf{v}_2$. So suppose this assumption holds for all the first $p+1$ factors. Then additionally, by Lemma 7, we have $\mathbf{v}_{p+1} \prec \mathbf{v}_1 \cdots \mathbf{v}_p\mathbf{v}_{p+1}$ holding in the extended suffix array. This property holds iteratively as we continue processing the factors of \mathbf{x} .

As expressed in [12] for the Lyndon case, this technique is suitable for integration with the on-line V -order factoring algorithm: suffix-sorting can proceed

in tandem as soon as the first V -factor is identified.

5. Computing the FM-Index in V -Order

Here we describe algorithms to compute the components of the FM-Index in V -order. We assume that for a given string \mathbf{x} , the V -BWT matrix $M_{\mathbf{x}}^{\mathcal{V}}$ (sorted conjugates in V -order) has been computed using the methodology described in Section 4. Since the FM-index stores the first (denoted F) and last (denoted L) columns of the Burrows-Wheeler matrix efficiently, our ultimate goal is to implement, in the context of V -order, BWT-type pattern matching using FM-index functions VC and VRANK analogous to those defined in Section 2.

It will be convenient for this discussion to assume that the given string \mathbf{x} is defined on an integer alphabet $\Sigma = \{0, 1, \dots, g\}$, where the least letter 0 corresponds to the sentinel letter \$, thus ensuring that every rotation of the string is distinct (\mathbf{x} is not a repetition). In order to adapt *RANK* and *C* from lexorder to V -order, we first establish a collection of results, Lemmas 24 – 29, on properties of $M_{\mathbf{x}}^{\mathcal{V}}$, and illustrated by references to the example $\mathbf{x} = \mathbf{w}'$ in Figure 2, which shows $M_{\mathbf{w}'}^{\mathcal{V}}$ corresponding to the following string:

$$\begin{array}{cccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ \mathbf{w}' = & 1 & 9 & 2 & 3 & 9 & 2 & 6 & 5 & 9 & 2 & 3 & 0 \end{array}$$

Lemma 24. *The Last First Mapping holds for V -order.*

Proof. We will establish the claim using the well-known BWT context sorting argument.

Given a string \mathbf{x} , let $\mathbf{r}_i = \mathbf{u}\lambda$ and $\mathbf{r}_j = \mathbf{v}\lambda$ be a pair of distinct rotations of \mathbf{x} occurring in the rows i, j of $M_{\mathbf{x}}^{\mathcal{V}}$ such that $\mathbf{r}_i \prec \mathbf{r}_j$, $\mathbf{u}, \mathbf{v} \neq \varepsilon$ and $\lambda \in \Sigma$. The context therefore is given by \mathbf{u} and \mathbf{v} .

By Lemma 8 we have $\mathbf{u}\lambda \prec \mathbf{v}\lambda \implies \mathbf{u} \prec \mathbf{v}$ and further $\mathbf{u} \prec \mathbf{v} \implies \lambda\mathbf{u} \prec \lambda\mathbf{v}$. That is, the relative order of rows \mathbf{r}_i and \mathbf{r}_j is maintained under the rotation of the letter λ .

G_{-1}	$\begin{bmatrix} 9 & 2 & 3 & 9 & 2 & 6 & 5 & 9 & 2 & 3 & 0 & 1 \\ 9 & 2 & 3 & 0 & 1 & 9 & 2 & 3 & 9 & 2 & 6 & 5 \\ 9 & 2 & 6 & 5 & 9 & 2 & 3 & 0 & 1 & 9 & 2 & 3 \end{bmatrix}$
G_1	$\begin{bmatrix} 0 & \underline{1} & 9 & 2 & 3 & 9 & 2 & 6 & 5 & 9 & 2 & 3 \\ \underline{1} & 9 & 2 & 3 & 9 & 2 & 6 & 5 & 9 & 2 & 3 & 0 \end{bmatrix}$
G_3	$\begin{bmatrix} \underline{3} & 9 & 2 & 6 & 5 & 9 & 2 & 3 & 0 & 1 & 9 & 2 \\ \underline{3} & 0 & 1 & 9 & 2 & 3 & 9 & 2 & 6 & 5 & 9 & 2 \\ 2 & \underline{3} & 9 & 2 & 6 & 5 & 9 & 2 & 3 & 0 & 1 & 9 \\ 2 & \underline{3} & 0 & 1 & 9 & 2 & 3 & 9 & 2 & 6 & 5 & 9 \end{bmatrix}$
G_5	$\begin{bmatrix} \underline{5} & 9 & 2 & 3 & 0 & 1 & 9 & 2 & 3 & 9 & 2 & 6 \end{bmatrix}$
G_6	$\begin{bmatrix} \underline{6} & 5 & 9 & 2 & 3 & 0 & 1 & 9 & 2 & 3 & 9 & 2 \\ 2 & \underline{6} & 5 & 9 & 2 & 3 & 0 & 1 & 9 & 2 & 3 & 9 \end{bmatrix}$

Figure 2: $M_{\mathbf{w}'}^{\mathcal{V}}$, for $\mathbf{w}' = 192392659230$, where $\$ = 0$. This matrix can be partitioned into five nonempty groups $G_{-1}, G_1, G_3, G_5, G_6$ as shown in the table. For each nonempty group G_j , let μ_j denote the maximum letter in the prefix \mathbf{w}'_0 before the first occurrence of $g = 9$ in every row of G_j . Then $\mu_{-1} = \varepsilon, \mu_1 = 1, \mu_3 = 3, \mu_5 = 5$ and $\mu_6 = 6$. For $j > -1$, the μ_j are underlined in each row. The groups G_0, G_2, G_4, G_7, G_8 are empty and not shown in the figure. EVSA = (2, 9, 5, 12, 1, 4, 11, 3, 10, 8, 7, 6), describing the ordering of the strings in $M_{\mathbf{w}'}^{\mathcal{V}}$, while $L = V\text{-BWT}_{\mathbf{w}'} = 153302299629$. Note that although \mathbf{w}' contains no letter repetitions, in contrast $V\text{-BWT}_{\mathbf{w}'}$ does, as is typical with this transformation scheme.

Consider now the sequence $S = \mathbf{r}_{h_1}, \mathbf{r}_{h_2}, \dots, \mathbf{r}_i, \dots, \mathbf{r}_j, \dots, \mathbf{r}_{h_t}$ of all rows in $M_{\mathbf{x}}^{\mathcal{V}}$ having suffix λ . Then the above argument holds for each adjacent pair in the sequence S and therefore the relative order of all rows in the sequence is maintained under the rotation of the letter λ .

Hence the $i - th$ occurrence of the suffix letter λ in the last column L is the same as the $i - th$ occurrence of the prefix letter λ in the first column F . An analogous argument shows the converse holds which establishes the claim. \square

For example, consider rows 3 and 4 of $M_{\mathbf{w}'}^{\mathcal{V}}$, where $\mathbf{u} = 92659230192$, $\mathbf{v} = 01923926592$, and both end with $\lambda = 3$. Since row 3 is above row 4, so that $\mathbf{u}\lambda \prec \mathbf{v}\lambda$, therefore the rotations $\lambda\mathbf{u}$ and $\lambda\mathbf{v}$, which correspond respectively to rows 6 and 7, and which begin with $\lambda = 3$, satisfy $\lambda\mathbf{u} \prec \lambda\mathbf{v}$.

Lemma 25. *Let \mathbf{x} (not containing $\$$) be a string with k maximum letters g and its associated $V\text{-BWT}$ matrix be $M_{\mathbf{x}}^{\mathcal{V}}$. Then, the k letters g occur as the first k letters in the first column F of $M_{\mathbf{x}}^{\mathcal{V}}$.*

Proof. We apply Lemma 5 (C3). Observe that any rotation \mathbf{x}^r of \mathbf{x} with $\mathbf{x}_0^r = \varepsilon$ in V -form must precede any rotation $\mathbf{x}^{r'}$ of \mathbf{x} with $\mathbf{x}_0^{r'} \neq \varepsilon$. Hence the k rotations of \mathbf{x} starting g will occur first in column F of $M_{\mathbf{x}}^{\mathcal{V}}$. (Also see [10] for a statement of this fact.) \square

Definition 26 (Groups). *Consider the rows of $M_{\mathbf{x}}^{\mathcal{V}}$ in their V -forms (Definition 4), and examine the (possibly empty) \mathbf{x}_0 substring of each row. The rows of $M_{\mathbf{x}}^{\mathcal{V}}$ are partitioned into collections of adjacent rows called **groups**, where each group is characterized by a letter $\mu < g$ which is the maximum within every \mathbf{x}_0 substring in the group, and is denoted by G_{μ} . When $\mu = \varepsilon$, we denote its corresponding group by G_{-1} .*

Lemma 27. *Let \mathbf{x} be a string, and $M_{\mathbf{x}}^{\mathcal{V}}$ be its associated V -BWT matrix.*

- (i) *If the first letter of column F is not $\$$, then the first group in $M_{\mathbf{x}}^{\mathcal{V}}$ is G_{-1} ; otherwise, it is G_0 .*
- (ii) *If the first letter of column F is $\$$, then the first two groups must be G_0, G_{-1} respectively, followed by groups in the increasing order of their characteristic letters μ .*

Proof. The proof is a direct consequence of the fact that $\$ \prec \varepsilon$, condition (C3) of Lemma 5 and Definition 26. \square

Lemma 28. *Consider the group G_{-1} in $M_{\mathbf{x}}^{\mathcal{V}}$ of all rows beginning with εg . Let $|G_{-1}| = k$ and denote by $k' \leq k$ the number of distinct letters in column L of G_{-1} that do not equal g .*

- (i) *Then there must exist k' groups in $M_{\mathbf{x}}^{\mathcal{V}}$ corresponding to these distinct letters.*
- (ii) *The k' distinct groups will occur in ascending order of the k' letters.*

Proof.

- (i) Let $g\mathbf{u}_1\lambda_1, g\mathbf{u}_2\lambda_2, \dots, g\mathbf{u}_{k'}\lambda_{k'}$ be the rows in $M_{\mathbf{x}}^{\mathcal{V}}$ ending in the k' suffix letters $\lambda_1, \lambda_2, \dots, \lambda_{k'}$. Since $M_{\mathbf{x}}^{\mathcal{V}}$ contains every rotation of \mathbf{x} , the rotations beginning with the k' suffix letters (which are now prefixes), that is, $\lambda_1g\mathbf{u}_1, \lambda_2g\mathbf{u}_2, \dots, \lambda_{k'}g\mathbf{u}_{k'}$, must also exist in $M_{\mathbf{x}}^{\mathcal{V}}$. Considering these rows in V -form, the x_0 portion in each row is a single letter (possibly equal to \$). Hence, by Definition 26, groups $G_{\lambda_1}, \dots, G_{\lambda_{k'}}$ exist, corresponding to rows $\lambda_1g\mathbf{u}_1, \dots, \lambda_{k'}g\mathbf{u}_{k'}$, respectively. Therefore, k' groups in $M_{\mathbf{x}}^{\mathcal{V}}$ exist corresponding to the k' distinct suffix letters.
- (ii) The proof of Part(ii) is a direct consequence of condition (C3) of Lemma 5 and Lemma 27. \square

For example, in Figure 2, the suffix letters of the first $k' = k = 3$ rows are 1, 5, 3 corresponding to groups G_1, G_3 , and G_5 with $\mu_1 = 1$, $\mu_3 = 3$, and $\mu_5 = 5$.

Observation. It is possible for a letter λ in column F to occur (but not as the maximum letter) in the \mathbf{x}_0 portion of distinct groups in a V -BWT matrix.

For example, in Figure 2, the letter 2 occurs in \mathbf{x}_0 for two groups G_3 and G_6 .

For the following lemma, recall that the V -form of a string \mathbf{x} is $V_k(\mathbf{x}) = \mathbf{x} = \mathbf{x}_0g\mathbf{x}_1g \cdots \mathbf{x}_{k-1}g\mathbf{x}_k$, where, if \mathbf{x} is over the alphabet $\Sigma = \{1 < 2 < \cdots\}$ and the final letter of \mathbf{x}_k is \$, then \$ is conveniently represented by 0.

Lemma 29. *Given a string $\mathbf{x} = \mathbf{x}_0g\mathbf{x}_1g \cdots \mathbf{x}_{k-1}g\mathbf{x}_k$, let $\mathbf{x}' = \mathbf{x}'_0g\mathbf{x}'_1g \cdots \mathbf{x}'_{k-1}g$ be the rotation of \mathbf{x} such that $\mathbf{x}'_0 = \mathbf{x}_k\mathbf{x}_0$ and $\mathbf{x}'_i = \mathbf{x}_i$ for $0 < i < k$. Suppose $\mathbf{z} = \mathbf{z}[1..j_i] = \mathbf{x}'_i$ is nonempty for some $0 \leq i < k$, and consider a letter $\mu = \mathbf{z}[j]$ in \mathbf{z} , $1 \leq j \leq j_i$.*

- (i) *Then the group G_μ exists corresponding to the rotation of \mathbf{x} starting with μ , that is, the row $\mathbf{z}[j..j_i]g\mathbf{x}'_{i+1}g \cdots \mathbf{x}'_{k-1}g\mathbf{x}'_0g\mathbf{x}'_1g \cdots \mathbf{x}'_{i-1}g\mathbf{z}[1..j-1]$ in $M_{\mathbf{x}}^{\mathcal{V}}$, if and only if $\mu \geq \mathbf{z}[j']$ for every $j' \in j+1..j_i$.*
- (ii) *Let μ occur t number of times in \mathbf{x}_i . Let p_1, p_2, \dots, p_t be the lengths of the maximum prefix preceding each of the t occurrences of μ in \mathbf{x}_i*

respectively, such that every letter in the prefix is strictly less than μ .

Then the total number of rows in G_μ corresponding to μ in \mathbf{x}_i is equal to $p_1 + p_2 + \dots + p_t + t$.

(iii) If the suffix letter of \mathbf{x}_k is $\$,$ then G_μ exists corresponding to $\mu = 0$ if and only if $\mathbf{x}_0 = \varepsilon$.

Proof. Part (i) (\Rightarrow). The proof is by contradiction. Suppose, there exists a largest $j' \in j+1..j_i$ such that $\mu < \mathbf{z}[j']$ and $\mathbf{z}[j']$ is maximal in $j+1..j_i$. Let

$$\mathbf{u} = \mathbf{x}'_{i+1}g \cdots \mathbf{x}'_{k-1}g \mathbf{x}'_0 g \mathbf{x}'_1 g \cdots \mathbf{x}'_{i-1}g \mathbf{z}[1..j-1].$$

Then,

$$r_h = \mathbf{z}[j..j'..j_i]g\mathbf{u}, r_{h+1} = \mathbf{z}[j+1..j'..j_i]g\mathbf{u}\mathbf{z}[j], \dots, r_{h+j'-j} = \mathbf{z}[j'..j_i]g\mathbf{u}\mathbf{z}[j..j'-1]$$

are the rows in the $M_{\mathbf{x}}^{\mathcal{V}}$ matrix corresponding to the rotations starting with letters $\mathbf{z}[j], \dots, \mathbf{z}[j']$ respectively, and all these prefixes contain $\mathbf{z}[j']$. By Definition 26, Lemma 27, and the hypothesis $\mu < \mathbf{z}[j']$, the rows $r_h, r_{h+1}, \dots, r_{h+j'-j}$ are in the vicinity of each other forming a group $G_{\mathbf{z}[j']}$ characterized by the maximum letter $\mathbf{z}[j']$. Clearly, $G_{\mathbf{z}[j']}$ contains r_h starting with $\mu = \mathbf{z}[j]$. Therefore, a group G_μ does not exist containing r_h , that is, the rotation of \mathbf{x} starting with $\mu = \mathbf{z}[j]$, as this rotation belongs to $G_{\mathbf{z}[j]}$ – hence a contradiction is established. The converse can be shown analogously.

Part(ii): By Part (i), all t rotations starting with μ in \mathbf{x}'_i occur in G_μ and therefore contribute to t rows. Further, any rotation of \mathbf{x}' starting with a prefix preceding an occurrence (say the r -th occurrence) of μ in \mathbf{x}'_i such that each letter in the prefix is strictly less than μ , also belongs to G_μ by Definition 26 and Lemma 5. Then the total number of such rows contributing to the size of G_μ by the r -th occurrence of μ in \mathbf{x}_i is equal to the maximum length of the associated prefix, that is, p_r . Therefore, the total number of rows in G_μ corresponding to all t occurrences of μ in \mathbf{x}_i is equal to $p_1 + p_2 + \dots + p_t + t$.

Part(iii) (\Rightarrow) Suppose G_0 exists in $M_{\mathbf{x}}^{\mathcal{V}}$. Then by Lemma 27, G_0 is the first

group in $M_{\mathbf{x}}^{\mathcal{V}}$ followed by G_{-1} . Then, there exists a rotation starting with $\$g$ in $M_{\mathbf{x}}^{\mathcal{V}}$. This is possible only when $\mathbf{x}_0 = \varepsilon$. Therefore (\Rightarrow) holds. The converse can be shown analogously. \square

As an example, in Figure 2, consider $\mathbf{w}'_2 = 265$, which contributes to groups with $\mu = 5$ and $\mu = 6$, but NOT $\mu = 2$. Further, as an example of case (iii) of the previous lemma, in the same example, $\mathbf{w}'_0 = 1 \neq \varepsilon$, and so the three rows starting with 9 appear at the start of $M_{\mathbf{w}'}^{\mathcal{V}}$, and the group G_0 does not exist. See Example 22 for an instance where $\mathbf{x}_0 = \varepsilon$.

We make use of Lemmas 24 – 29 to determine the number of rotations of the given string \mathbf{x} that occur within each of the ordered groups G_0, G_1, \dots, G_{g-1} . Using it, we compute $VC = VC[0..g-1]$, the V -order version of the C array defined originally for lexorder:

Definition 30. (Array $VC[0, 1, \dots, g-1]$) Suppose $\mathbf{x} = \mathbf{x}[1..n]$ is defined on an integer alphabet $\Sigma = \{0, 1, \dots, g\}$, where $\mathbf{x}[n] = 0$ corresponds to the sentinel letter $\$$. Then for each letter $\mu \in \Sigma$ that corresponds to a nonempty group G_μ , $0 \leq \mu < g$, $VC[\mu]$ is the sum of the sizes of the groups G_j , $0 \leq j \leq \mu-1$; that is,

$$VC[\mu] = \sum_{j=0}^{\mu-1} |G_j|;$$

for all other letters μ , $VC[\mu] = 0$.

For instance, the VC array for the string $\mathbf{w}' = 192392659230$ is $VC = (0, 0, 0, 2, 0, 6, 7, 0, 0)$.

The computation of the VC array is given by Procedure COMPUTE_VC in Figure 3. Based on Lemmas 24 – 29 and Definitions 26 & 30, we state:

Lemma 31. Procedure COMPUTE_VC correctly computes the VC array.

Proof. The proof is a direct consequence of Lemma 29. \square

Lemma 32. Procedure COMPUTE_VC computes the VC array in $O(n)$ time.

```

procedure COMPUTE_VC( $\mathbf{x}, n$ )
▷ Scan  $\mathbf{x}$  to compute  $g$  and  $k$ 
 $g \leftarrow -1$ 
for  $i \leftarrow 1$  to  $n$  do
    if  $\mathbf{x}[i] > g$  then  $g \leftarrow \mathbf{x}[i]$ ;  $k \leftarrow 1$ 
    elseif  $\mathbf{x}[i] = g$  then  $k \leftarrow k+1$ 
     $i \leftarrow i+1$ 
▷ Compute  $\mathbf{x}' = \text{rotated } \mathbf{x}$  (Lemma 29):  $\mathbf{x}'[n] = g$ 
 $h \leftarrow 0$ 
while  $\mathbf{x}[n-h] < g$  do  $h \leftarrow h+1$ 
 $\mathbf{x}' \leftarrow \mathbf{x}[n-h+1..n]\mathbf{x}[1..n-h]$ 
▷ Compute  $G$  (Lemma 29)
 $G[0..g-1] \leftarrow 0^g$ 
for  $i = n$  downto  $1$  do
    if  $\mathbf{x}'[i] = g$  then  $\mu \leftarrow -1$ 
    else
        if  $\mu < \mathbf{x}'[i]$  then  $\mu \leftarrow \mathbf{x}'[i]$ 
         $G[\mu] \leftarrow G[\mu]+1$ 
     $i \leftarrow i-1$ 
▷ Compute  $VC[0..g-1]$ 
 $VC[0..g-1] \leftarrow 0^g$ ;  $\mu_0 \leftarrow 0$ 
while  $\mu_0 < g$  and  $G[\mu_0] = 0$  do  $\mu_0 \leftarrow \mu_0+1$ 
 $lastmu \leftarrow \mu_0$ 
for  $\mu = \mu_0+1$  to  $g-1$  do
    if  $G[\mu] \neq 0$  then
         $VC[\mu] \leftarrow VC[lastmu]+G[lastmu]$ ;  $lastmu \leftarrow \mu$ 
     $\mu \leftarrow \mu+1$ 

```

Figure 3: Computing the VC array.

Proof. The procedure COMPUTE_VC has **for** and **while** loops that are executed at most $O(n)$ time. Therefore, the running time of the procedure is $O(n)$. \square

We now discuss the FM-index RANK function with respect to V -order; that is, the VRANK function defined as follows:

Definition 33 (VRANK Function). $VRANK_{\mathbf{x}}(I, \ell)$ gives the number of occurrences of letter ℓ in the substring of \mathbf{x} defined by the interval $I = (i, j)$, where $1 \leq i \leq j \leq |\mathbf{x}|$.

As can be observed, the only difference between RANK and VRANK is that RANK returns the rank of the letter ℓ in a certain prefix of \mathbf{x} , while VRANK

		L	Tally Table (T_L)						
			0	1	2	3	5	6	9
	1	1	0	1	0	0	0	0	0
	2	5	0	1	0	0	1	0	0
	3	3	0	1	0	1	1	0	0
	4	0	1	1	0	1	1	0	0
	5	3	1	1	0	2	1	0	0
$i =$	6	2	1	1	1	2	1	0	0
	7	2	1	1	2	2	1	0	0
	8	9	1	1	2	2	1	0	1
$j =$	9	9	1	1	2	2	1	0	2
	10	6	1	1	2	2	1	1	2
	11	2	1	1	3	2	1	1	2
	12	9	1	1	3	2	1	1	3

Figure 4: $\mathbf{w}' = 192392659230$, and $L = V\text{-BWT}_{\mathbf{w}'} = 153302299629$.

returns the rank of ℓ in the substring $\mathbf{x}[i..j]$ of \mathbf{x} . Then VRANK of a letter λ in the substring $\mathbf{x}[i..j]$ of \mathbf{x} can be computed as follows:

$$\text{VRANK}_{\mathbf{x}}((i, j), \lambda) = \begin{cases} \text{VRANK}_{\mathbf{x}}((1, j), \lambda), & \text{if } i = 1 \\ \text{VRANK}_{\mathbf{x}}((1, j), \lambda) - \text{VRANK}_{\mathbf{x}}((1, i-1), \lambda), & \text{if } i > 1 \end{cases} \quad (6)$$

Analogous to RANK, VRANK can be computed easily using the Tally table for $L(T_L)$ — a matrix where each column represents a letter in the string and the i -th row represents the number of occurrences of the corresponding letter in the prefix of length i of L . Then a cell $T_L(i, j)$ in the Tally table stores the number of occurrences of letter j in the prefix of length i of L . Therefore, $T_L(i, j) = \text{VRANK}_L((1, i), j)$. For example, consider the Tally table given in Figure 4: by Equation (6),

$$\text{VRANK}_{\mathbf{w}'}((6, 9), 2) = \text{VRANK}_{\mathbf{w}'}((1, 9), 2) - \text{VRANK}_{\mathbf{w}'}((1, 5), 2) = 2 - 0 = 2.$$

The Tally table can be computed in linear time in the length of the string n , and the space required by it is $\mathcal{O}(\sigma n)$. As discussed in [19], this space can be

significantly reduced by storing only certain rows, starting at the first one, and occurring at a predetermined constant offset (α) in the Tally table. A row that is no longer stored in the Tally table is then computed in constant time by using existing rows and simple arithmetic calculations (for more details see [19]). The total size of the reduced Tally table is $\mathcal{O}(\frac{\sigma n}{\alpha})$.

Application of VC and VRANK to fully implement BWT-type pattern matching in V -order, analogous to the results in [12], remains an open problem.

6. V -Order Substring Matching

In this section we introduce V -order pattern matching using the backward search technique for substrings rather than letters. For clarity of concept, we exclude the use of the sentinel $\$$ in this section.

To achieve substring matching, the last column in the V -BWT matrix will be regarded as a column of substrings rather than letters, which we denote as VL . Specifically, VL consists of the last suffix $g\mathbf{x}_k$ of the V -form of each row in the V -BWT matrix. Likewise, we get the first column VF of prefix substrings in the V -BWT matrix, where the first k entries will be εg .

The substrings to be searched for (successfully) in a given string \mathbf{x} take the form of a substring \mathbf{X} of the \mathbf{x}_i substrings in the V -form of \mathbf{x} ; thus $\mathbf{X} = g\mathbf{x}_i \cdots g\mathbf{x}_j$ for $1 \leq i \leq j \leq k$. So the greatest letters g are basically for demarcation only. This demarcation ensures that string comparison is restricted to the \mathbf{x}_i substrings and avoids the complex intermingling illustrated in Figure 5. Furthermore, if useful for particular required demarcation, such as length of substrings, an artificial greatest letter can be inserted in specified positions of the string \mathbf{x} .

We now show a property for substrings in V -order analogous to the lexorder BWT Last First mapping, in which the i -th occurrence of a letter in L is the same as the i -th occurrence of the letter in F (see Lemma 24).

Lemma 34. *The Last First Mapping holds for substrings in V -order.*

Proof. The proof follows immediately from the proof of Lemma 24 by replacing the letter λ with a substring and applying Lemma 8. \square

In order to work with just the first k rows in the V -BWT matrix, we apply this lemma to modify the *FM-Index* [19] functions in the obvious way (see definitions of these functions for letters in Section 2):

- $C_{\mathbf{x}}(g\ell)$ gives, for each substring $g\ell \in \Sigma^+$ that occurs in \mathbf{x} , the number of substrings $g\mathbf{x}_i$ in \mathbf{x} smaller in V -order than $g\ell$;
- $\text{RANK}_{\mathbf{x}}(p, g\ell)$ gives the number of occurrences of substring $g\ell$ in the prefix of \mathbf{x} of length p .

9	21	9	12	9	21	9	<u>1314</u>	
9	21	9	1314	9	21	<u>9</u>	<u>12</u>	
9	12	9	21	9	1314	<u>9</u>	<u>21</u>	
9	1314	9	21	9	12	<u>9</u>	<u>21</u>	
<u>1</u>	9	12	9	21	9	1314	<u>9</u>	<u>2</u>
<u>1</u>	9	1314	9	21	9	12	<u>9</u>	<u>2</u>
<u>2</u>	9	21	9	1314	9	21	<u>9</u>	<u>1</u>
<u>21</u>	9	12	9	21	9	1314	<u>9</u>	
<u>21</u>	9	1314	9	21	9	12	<u>9</u>	
<u>12</u>	9	21	9	1314	9	21	<u>9</u>	
<u>4</u>	9	21	9	12	9	21	<u>9</u>	<u>131</u>
<u>14</u>	9	21	9	12	9	21	<u>9</u>	<u>13</u>
<u>314</u>	9	21	9	12	9	21	<u>9</u>	<u>1</u>
<u>1314</u>	9	21	9	12	9	21	<u>9</u>	

Figure 5: V -BWT $_{\mathbf{x}}$, where the columns VF and VL are underlined.

Example 35 (Substring matching). For $\mathbf{x} = 49219129219131$, the V -BWT matrix $M_{\mathbf{x}}^V$ is shown in Figure 5. Consider the substring 921 in the fourth row of the column VL . To calculate its index in VF we have $C_{\mathbf{x}}(921) = 0$ and $\text{RANK}_{VL}(4, 921) = 2$ (alternatively, $\text{RANK}_{VL}(3, 921) + 1 = 1 + 1 = 2$ depending on the formulation of the definition). Thus, the Last First mapping technique maps 921 to index 0+2, namely the second position in VF .

The substrings 921 in Example 35 are clustered into a run of length 2. Working with substrings potentially may allow for more effective preprocessing

for compression as the substrings can be encoded, such as $A = 921$, and we propose this avenue to be further examined.

Note that to compute the V -BWT matrix for substrings — that is, the first k rows — we encode each substring as above to generate a new alphabet and apply a linear sorting algorithm such as radix sort along with linear V -order string comparison to order the alphabet. Then we apply the techniques in [10] to compute the V -BWT matrix with respect to substrings in linear time overall. Generally this will result in space reduction for the columns VF and VL .

7. Future Research

As mentioned in Section 5, an immediate research problem is the application of VC and VRANK to BWT-type pattern matching in V -order, analogous to the lexorder results in [12].

Computing the V -BWT for substrings invites thorough experimental investigation in terms of the numbers of runs of similar substrings (rather than the usual letters) obtained. Furthermore, lex-extension order opens avenues for various techniques for substring clustering with appropriately defined BWTs: as discussed in [4], this hybrid ordering method can combine lexorder with any other total order on the factors induced in the factorization — here we combined lexorder with applying V -order to the factors, and propose considering other combinations of ordering methods.

The natural ordering of suffixes in V -order by their length suggests that further simple suffix computations and applications should be achievable. For instance, V -letters [4] are defined as V -words where the first letter in a string is strictly greater than the subsequent letters; for example, the V -word 812321. Then applying the trivial suffix-sorting to V -letters, the BWT V -transform is obtained simply by reversing the V -letter, and a V -letter can be similarly recovered from the V -transform [10]. We seek further efficient applications of these ideas. The computation of parallel V -ordering is also open.

Finally we propose the following problem. Suppose that $\mathbf{x}, \mathbf{y} \in \Sigma^+$ with

$\mathbf{x} \prec \mathbf{y}$. Under what permutations π — that is, $\mathbf{x} \rightarrow \pi(\mathbf{x})$ and $\mathbf{y} \rightarrow \pi(\mathbf{y})$ — does $\pi(\mathbf{x}) \prec \pi(\mathbf{y})$ hold? For instance, for $21 \prec 12$ no permutation works, whereas interchanging the first and last letters does work for $142 \prec 243$, since $241 \prec 342$. This example generalizes to requiring that the rightmost substrings of the V -forms are in V -order.

Acknowledgements

We thank the reviewers for carefully reading the manuscript and providing insightful comments which led to substantial improvements.

Funding: The first author was part-funded by the European Regional Development Fund through the Welsh Government [Grant Number 80761-AU-137 (West)]:



The second and third authors were funded by the Natural Sciences & Engineering Research Council of Canada [Grant Number 10536797].

References

- [1] A. Alatabbi, J. W. Daykin, N. Mhaskar, M. S. Rahman, W. F. Smyth, Applications of V -order: Suffix arrays, the Burrows-Wheeler transform & the FM-index, in: G. K. Das, P. S. Mandal, K. Mukhopadhyaya, S. ichi Nakano (Eds.), International Workshop on Algorithms & Computation (WALCOM), Vol. 11355 of Lecture Notes in Computer Science, Springer, 2019, pp. 329–338.
- [2] T.-N. Danh, D. E. Daykin, The structure of V -Order for integer vectors, Congr. Numer. Ed. A.J.W. Hilton. Utilas Mat. Pub. Inc., Winnipeg, Canada, 113 (1996) (1996) 43–53.

- [3] D. E. Daykin, J. W. Daykin, W. F. Smyth, String comparison and Lyndon-like factorization using V-Order in linear time, in: *Symp. on Combinatorial Pattern Matching*, Vol. 6661, 2011, pp. 65–76.
- [4] D. E. Daykin, J. W. Daykin, W. F. Smyth, A linear partitioning algorithm for hybrid Lyndons using V-Order, *Theoret. Comput. Sci.* 483 (2013) 149–161.
- [5] A. Alatabbi, J. W. Daykin, M. S. Rahman, W. F. Smyth, Simple linear comparison of strings in V-Order – (extended abstract), in: *International Workshop on Algorithms & Computation (WALCOM)*, Vol. 8344 of *Lecture Notes in Computer Science*, Springer, 2014, pp. 80–89.
- [6] A. Alatabbi, J. W. Daykin, M. S. Rahman, W. F. Smyth, String comparison in V-Order: New lexicographic properties & on-line applications, *arXiv:1507.07038* (2015).
- [7] A. Alatabbi, J. W. Daykin, M. S. Rahman, W. F. Smyth, Simple linear comparison of strings in V-Order, *Fundam. Inform.* 139 (2) (2015) 115–126. doi:10.3233/FI-2015-1228.
URL <https://doi.org/10.3233/FI-2015-1228>
- [8] A. Alatabbi, J. W. Daykin, J. Kärkkäinen, M. S. Rahman, W. F. Smyth, V-Order: New combinatorial properties & a simple comparison algorithm, *Discrete Appl. Math.* 215 (2016) 41–46.
- [9] A. Alatabbi, J. W. Daykin, N. Mhaskar, M. S. Rahman, W. F. Smyth, A faster V-Order string comparison algorithm, in: *Proc. Prague Stringology Conference*, 2018, pp. 38–49.
- [10] J. W. Daykin, W. F. Smyth, A bijective variant of the Burrows–Wheeler transform using V-Order, *Theoret. Comput. Sci.* 531 (2014) 77–89.
- [11] P. Ko, S. Aluru, Space efficient linear time construction of suffix arrays, in: R. A. Baeza-Yates, E. Chávez, M. Crochemore (Eds.), *Symp. on Combina-*

- torial Pattern Matching, Vol. 2676 of Lecture Notes in Computer Science, Springer, 2003, pp. 200–210.
- [12] S. Mantaci, A. Restivo, G. Rosone, M. Sciortino, Suffix array and Lyndon factorization of a text, *J. Discrete Algorithms* 28 (2014) 2–8. doi:10.1016/j.jda.2014.06.001.
URL <http://dx.doi.org/10.1016/j.jda.2014.06.001>
- [13] U. Baier, Linear-time suffix sorting — a new approach for suffix array construction, in: *Proc. 27th Annual Symposium on Combinatorial Pattern Matching (CPM 2016)*, 2016, pp. 23:1–23:12.
- [14] H. Bannai, T. I. S. Inenaga, Y. Nakashima, M. Takeda, K. Tsuruta, The “runs” theorem, *SIAM J. Comput.* 46 (5) (2017) 1501–1514.
- [15] J. W. Daykin, F. Franek, J. Holub, A. S. M. Sohidull Islam, W. F. Smyth, Reconstructing a string from its Lyndon arrays, *Theor. Comput. Sci.* 710 (2018) 44–51.
- [16] F. Franek, A. Paracha, W. F. Smyth, The linear equivalence of the suffix array and the partially sorted Lyndon array, in: *Proc. Prague Stringology Conf.*, 2017, pp. 85–92.
- [17] F. Franek, A. S. M. Sohidull Islam, M. S. Rahman, W. F. Smyth, Algorithms to compute the Lyndon array, in: *Proc. Prague Stringology Conf.*, 2016, pp. 172–184.
- [18] D. Adjeroh, T. Bell, A. Mukherjee, *The Burrows–Wheeler Transform: Data Compression, Suffix Arrays, and Pattern Matching*, Springer Publishing Company, 2008.
- [19] P. Ferragina, G. Manzini, Opportunistic data structures with applications, in: *Proc. 41st Annual Symposium on Foundations of Computer Science, (FOCS 2000)*, 2000, pp. 390–398.

- [20] K. T. Chen, R. H. Fox, R. C. Lyndon, Free differential calculus, iv – the quotient groups of the lower central series, *Ann. Math.* (1958) 68:81–95.
- [21] M. Crochemore, C. Hancart, T. Lecroq, *Algorithms on Strings*, Cambridge University Press, New York, NY, USA, 2007.
- [22] B. Smyth, *Computing Patterns in Strings*, Pearson/Addison–Wesley, 2003.
- [23] D. E. Daykin, J. W. Daykin, Lyndon–like and V-Order factorizations of strings, *J. Discrete Algorithms* 1 (3–4) (2003) 357–365.
- [24] D. E. Daykin, J. W. Daykin, W. F. Smyth, Combinatorics of unique maximal factorization families (UMFFs), *Fund. Inform.* 97–3, Special Issue on Stringology, R. Janicki, S. J. Puglisi and M. S. Rahman, editors (2009) 295–309.