

# A Faster $V$ -order String Comparison Algorithm

Ali Alatabbi<sup>1</sup>, Jacqueline W. Daykin<sup>1,2,3</sup>, Neerja Mhaskar<sup>4</sup>, M. Sohel Rahman<sup>5</sup>, and William F. Smyth<sup>1,4,6</sup>

<sup>1</sup> Department of Informatics, King's College London, UK  
ali.alatabbi@kcl.ac.uk, jackie.daykin@kcl.ac.uk

<sup>2</sup> Department of Computer Science  
Aberystwyth University, Wales & Mauritius

<sup>3</sup> Department of Information Science  
Stellenbosch University, South Africa

<sup>4</sup> Algorithms Research Group, Department of Computing & Software  
McMaster University, Canada  
pophlin@mcmaster.ca, smyth@mcmaster.ca

<sup>5</sup> Department of Computer Science & Engineering  
Bangladesh University of Engineering & Science  
msrahman@cse.buet.ac.bd

<sup>6</sup> School of Engineering & Information Technology  
Murdoch University, Western Australia

**Abstract.**  $V$ -order is a total order on strings that determines an instance of Unique Maximal Factorization Families (UMFFs) [7–10], a generalization of Lyndon words [12].  $V$ -order has also recently 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) [11]. The central problem of efficient  $V$ -ordering of strings was considered in [2–4, 9, 10], culminating in a remarkably simple, linear time, constant space comparison algorithm [1]. In this paper we improve on this result to achieve significant speed-up in almost all cases of interest.

**Keywords:** combinatorics, experiments, lexorder, linear, on-line algorithm, probability, string comparison,  $V$ -comparison,  $V$ -order

## 1 Introduction

This paper extends current knowledge on the non-lexicographic ordering technique known as  $V$ -order [6]. New combinatorial insights are obtained which are linked to computational settings. The first of these leads to an improvement on the linear-time  $V$ -order string comparison algorithm described in [1].

## 2 Preliminaries

We are given a finite totally ordered set of cardinality  $\sigma = |\Sigma|$ , called the *alphabet*, whose elements are *characters* (equivalently *letters*). A *string* is a sequence of zero or more characters over  $\Sigma$ . A string  $\mathbf{x} = x_1x_2 \cdots x_n$  of length  $|\mathbf{x}| = n$  is represented by  $\mathbf{x}[1..n]$ , where  $\mathbf{x}[i] \in \Sigma$  for  $1 \leq i \leq n$ . The set of all non-empty strings over the alphabet  $\Sigma$  is denoted by  $\Sigma^+$ . The *empty string* of zero length is denoted by  $\varepsilon$ , with  $\Sigma^* = \Sigma^+ \cup \varepsilon$ . If  $\mathbf{x} = \mathbf{uvw}$  for strings  $\mathbf{u}, \mathbf{w}, \mathbf{v} \in \Sigma^*$ , then  $\mathbf{u}$  is a *prefix*,  $\mathbf{w}$  is a *substring* or *factor*, and  $\mathbf{v}$  is a *suffix* of  $\mathbf{x}$ . We denote by  $\mathbf{s}[i \dots j]$ , or  $s_i \cdots s_j$ , the substring of  $\mathbf{s}$  that starts at position  $i$  and ends at position  $j$ . Notably, if  $i > j$ ,  $\mathbf{s}[i \dots j] = \varepsilon$ ; that is,  $\mathbf{s}[i \dots j]$  is the empty string. 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*. For further stringological definitions, theory and algorithmics see [5, 13].

We define an order that is not lexorder (dictionary order), called V-order, as well as some of its notable properties, both combinatorial and algorithmic. V-order was explored in [6] and subsequently studied extensively in the literature both combinatorially and algorithmically [1–4, 7, 9–11].

Let  $\mathbf{x} = x_1x_2 \cdots x_n$  be a string over  $\Sigma$ . 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 the letter  $x_h$ . We write  $\mathbf{x}^{s*}$  for  $(\dots(\mathbf{x}^*)^* \dots)^*$  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 in  $g^k, \dots, g^1, g^0 = \varepsilon$ . From all strings  $\mathbf{x}$  over  $\Sigma$ , we use this process to 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  $\varepsilon$  as the root.

**Definition 1.** [6, 7, 9, 10] We define V-order  $\prec$  between distinct strings  $\mathbf{x}, \mathbf{y}$ . 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 not, then 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  $\Sigma$  then  $\mathbf{x} \prec \mathbf{y}$ ; otherwise,  $\mathbf{y} \prec \mathbf{x}$ . Clearly  $\prec$  is a total order on all strings in  $\Sigma^*$ .

For instance, using the natural ordering of integers, if  $\mathbf{x} = 2631$ , then  $\mathbf{x}^* = 263$ ,  $\mathbf{x}^{2*} = 26$ ,  $\mathbf{x}^{3*} = 6$ ,  $\mathbf{x}^{4*} = \varepsilon$ , and so  $26 \prec 2631$  while  $2631 \prec 94$ .

**Definition 2.** [6, 7, 9, 10] 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 (possibly empty) 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. For clarity, when more than one string is involved, we use the notation  $g = \mathcal{L}\mathbf{x}$ ,  $k = \mathcal{C}\mathbf{x}$ .

**Lemma 3.** [6, 7, 9, 10] 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 hold:

- (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 substring pair  $\mathbf{x}_h, \mathbf{y}_h$  can likewise be decomposed into V-forms.

**Lemma 4.** [9, 10] For given strings  $\mathbf{x}$  and  $\mathbf{v}$ , if  $\mathbf{v}$  is a proper subsequence of  $\mathbf{x}$ , then  $\mathbf{v} \prec \mathbf{x}$ .

**Theorem 5.** [1, 4] For any strings  $\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}$ :  $\mathbf{x} \prec \mathbf{y} \Leftrightarrow \mathbf{uxv} \prec \mathbf{uyv}$ .

Note that, according to this result, a comparison of two strings can ignore equal prefixes (or suffixes) — see Algorithm COMPARE.

**Lemma 6.** [1] For any two strings  $\mathbf{x}, \mathbf{y}$  with  $\mathbf{x} \prec \mathbf{y}$  and any two letters  $\lambda, \mu$  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  $\mathbf{y}\mu \prec \mathbf{x}\lambda$ .

Lemma 6 led to Algorithm COMPARE (Figure 1), described in [1]<sup>1</sup>. Note that, in contrast to comments made in the Conclusion of [10], this tells us that in fact  $V$ -order comparison can be conducted in a positional manner as in lexorder comparison — in fact, COMPARE is an on-line algorithm [13] that requires only a one-character window. Notably, the first  $V$ -order comparison algorithm was presented in [9], and this was followed up by a couple of other interesting algorithms [2, 3, 10] before algorithm COMPARE was presented in [1].

```

procedure COMPARE( $\mathbf{x}, \mathbf{y}, \delta$ )
   $i \leftarrow 1$ ;  $j \leftarrow 1$ ;  $\delta \leftarrow 0$ 
  while  $j \leq |\mathbf{y}|$  and  $i \leq |\mathbf{x}|$  and  $\mathbf{x}[i] = \mathbf{y}[j]$  do
     $i \leftarrow i + 1$ ;  $j \leftarrow j + 1$ 
  if  $i > |\mathbf{x}|$  or  $j > |\mathbf{y}|$  then
    if  $|\mathbf{x}| = |\mathbf{y}|$  then
       $\delta \leftarrow 0$ ; return
    if  $i > |\mathbf{x}|$  then
       $\delta \leftarrow -1$ ; return
     $\delta \leftarrow 1$ ; return
  while true do
    while  $j \leq |\mathbf{y}|$  and  $\mathbf{x}[i] > \mathbf{y}[j]$  do
       $j \leftarrow j + 1$ 
    if  $j > |\mathbf{y}|$  then
       $\delta \leftarrow 1$ ; return
    else
       $i \leftarrow i + 1$ 
    while  $i \leq |\mathbf{x}|$  and  $\mathbf{x}[i] < \mathbf{y}[j]$  do
       $i \leftarrow i + 1$ 
    if  $i > |\mathbf{x}|$  then
       $\delta \leftarrow -1$ ; return
    else
       $j \leftarrow j + 1$ 

```

**Figure 1.** Comparing  $\mathbf{x}$  of length  $m$  and  $\mathbf{y}$  of length  $n$  in  $V$ -order:  $\mathbf{x} \prec, =, > \mathbf{y}$  according as  $\delta = -1, 0, 1$ .

**Lemma 7.** [1–3, 9, 10]  $V$ -comparison requires linear time and constant space.

<sup>1</sup> This version corrects the comparisons in lines 5 & 8 of the original presentation.

### 3 $V$ -order String Comparison

In this section, we present a new algorithm to compare strings in  $V$ -order, which is sensitive to the structure of the input strings – the **COMPARE-Sensitive** algorithm (Figure 2). We show that probabilistically **COMPARE-Sensitive** runs faster than **COMPARE** in almost all cases of interest. We also show experimental results comparing the two algorithms. From these experiments we see that **COMPARE-Sensitive** runs at least twice as fast as **COMPARE**, in almost all cases of interest. Moreover, the new algorithm can easily be converted to work in an on-line setting, just like **COMPARE**. In addition to the new algorithm, we give a minor correction to the **COMPARE** algorithm (Figure 1).

#### 3.1 **COMPARE-Sensitive** Algorithm

Suppose we are asked to compare two strings  $\mathbf{x}$ ,  $\mathbf{y}$  with  $V$ -forms

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

In this representation we can “most of the time” determine the order simply by applying conditions (C1) and (C2) of Lemma 3 to  $\mathcal{L}\mathbf{x}$  and  $\mathcal{L}\mathbf{y}$ : if one maximum is greater than the other, or if the maxima are equal but have different frequencies of occurrences, then we are done. The **COMPARE-Sensitive** Algorithm (Figure 2) uses this observation to compare strings  $\mathbf{x}$  and  $\mathbf{y}$ . The first few lines of the algorithm show the very simple preprocessing required for each string  $\mathbf{x}, \mathbf{y}$  to check for conditions (C1) and (C2).

```

procedure COMPARE-Sensitive( $\mathbf{x}, \mathbf{y}, \delta$ )
  SCAN( $\mathbf{x}, |\mathbf{x}|; \mathcal{L}\mathbf{x}, occx$ )
  SCAN( $\mathbf{y}, |\mathbf{y}|; \mathcal{L}\mathbf{y}, occy$ )
  if  $\mathcal{L}\mathbf{x} \neq \mathcal{L}\mathbf{y}$  then
     $\delta \leftarrow \text{SIGN}(\mathcal{L}\mathbf{x} - \mathcal{L}\mathbf{y})$ 
  else
    if  $occx \neq occy$  then
       $\delta \leftarrow \text{SIGN}(occx - occy)$ 
    else
      COMPARE-C3( $\mathbf{x}, \mathbf{y}, \mathcal{L}\mathbf{x}, \delta$ )

```

**Figure 2.** Apply conditions (C1) and (C2) to  $\mathbf{x}$  and  $\mathbf{y}$ : for  $n \gg \sigma$ , almost always  $\mathcal{L}\mathbf{x} = \mathcal{L}\mathbf{y}$ , while at the same time usually  $occx \neq occy$ , so that the routine **COMPARE-C3** does not need to be executed at all.

```

procedure SCAN( $\mathbf{x}, |\mathbf{x}|; \mathcal{L}\mathbf{x}, occx$ )
   $\mathcal{L}\mathbf{x} \leftarrow \mathbf{x}[1]; occx \leftarrow 1$ 
  for  $i \leftarrow 2$  to  $|\mathbf{x}|$  do
    if  $\mathbf{x}[i] \notin \mathcal{L}\mathbf{x}$  then
      if  $\mathbf{x}[i] = \mathcal{L}\mathbf{x}$  then
         $occx \leftarrow occx + 1$ 
      else
         $\mathcal{L}\mathbf{x} \leftarrow \mathbf{x}[i]; occx \leftarrow 1$ 

```

**Figure 3.** Scan traverses the string  $\mathbf{x}$  to compute the maximum  $\mathcal{L}\mathbf{x}$  and its frequency  $occx$ .

```

procedure COMPARE-C3( $x, y, \mathcal{L}x, \delta$ )
   $i \leftarrow 1; j \leftarrow 1; \delta \leftarrow 0$ 
  while  $j \leq |y|$  and  $i \leq |x|$  and  $x[i] = y[j]$  do
     $i \leftarrow i + 1; j \leftarrow j + 1$ 
  if  $i > |x|$  or  $j > |y|$  then
    if  $|x| = |y|$  then
       $\delta \leftarrow 0$ ; return
    if  $i > |x|$  then
       $\delta \leftarrow -1$ ; return
     $\delta \leftarrow 1$ ; return
  if  $x[i] = \mathcal{L}x$  then  $\delta \leftarrow -1$ ; return  $\triangleright s(\mathbf{x}_h) = \varepsilon$ 
  if  $y[j] = \mathcal{L}y$  then  $\delta \leftarrow 1$ ; return  $\triangleright s(\mathbf{y}_h) = \varepsilon$ 
  while true do
    while  $j \leq |y|$  and  $x[i] > y[j]$  and  $y[j] < \mathcal{L}x$  do
       $j \leftarrow j + 1$ 
    if  $y[j] = \mathcal{L}y$  or  $j > |y|$  then  $\triangleright \mathcal{L}x = \mathcal{L}y$  and  $y[j] \not\prec \mathcal{L}x$ 
       $\delta \leftarrow 1$ ; return
    else
       $i \leftarrow i + 1$ 
    while  $i \leq |x|$  and  $x[i] < y[j]$  and  $x[i] < \mathcal{L}x$  do
       $i \leftarrow i + 1$ 
    if  $x[i] = \mathcal{L}x$  or  $i > |x|$  then  $\triangleright x[i] \not\prec \mathcal{L}x$ 
       $\delta \leftarrow -1$ ; return
    else
       $j \leftarrow j + 1$ 

```

**Figure 4.** Comparing  $x$  of length  $m$  and  $y$  of length  $n$  in  $V$ -order, when  $\mathcal{L}x = \mathcal{L}y$  and  $occx = occy$ :  $x \prec, =, \succ y$  according as  $\delta = -1, 0, 1$ .

However if conditions (C1) and (C2) of Lemma 3 fail, we need to invoke the COMPARE algorithm to check for condition (C3). A potential disadvantage in invoking COMPARE is that it must rescan at least one of the strings in its full length. We propose a modification to the COMPARE algorithm, the COMPARE-C3 algorithm, which takes advantage of the values  $\mathcal{L}x = \mathcal{L}y$  computed by SCAN to avoid this rescanning. COMPARE-C3 is motivated by the following observations:

**Observation 1** Recall that the condition (C3) of Lemma 3 first eliminates a common prefix (up to the beginning of substrings  $\mathbf{x}_h$  and  $\mathbf{y}_h$ ) in strings  $x$  and  $y$ , and then compares only the substrings  $\mathbf{x}_h$  and  $\mathbf{y}_h$  containing the first mismatching letter in  $x$  and  $y$  while scanning the strings from left to right. Therefore, instead of comparing the entire strings  $x$  and  $y$  as in COMPARE, it suffices to first identify the substrings  $\mathbf{x}_h$  and  $\mathbf{y}_h$  and only compare them.

**Observation 2** Observe that the substring  $\mathbf{x}_h$  ( $\mathbf{y}_h$ ) is either followed by the empty string (when  $\mathbf{x}_h$  ( $\mathbf{y}_h$ ) is a suffix of  $x$  ( $y$ )), or  $\mathcal{L}x = \mathcal{L}y$  (when  $\mathbf{x}_h$  ( $\mathbf{y}_h$ ) is not a suffix of  $x$  ( $y$ )). We use this observation to identify the end of substring  $\mathbf{x}_h$  ( $\mathbf{y}_h$ ) in  $x$  ( $y$ ).

Similar to COMPARE, the first **while** loop in COMPARE-C3 identifies a common prefix (alternatively the first mismatching position in strings  $x$  and  $y$  from the left). Note that this prefix might include a common prefix of substrings  $\mathbf{x}_h$  and  $\mathbf{y}_h$  identified under condition (C3) of Lemma 3. We denote the suffixes of  $\mathbf{x}_h$  and  $\mathbf{y}_h$  without a common prefix by  $s(\mathbf{x}_h)$  and  $s(\mathbf{y}_h)$ , respectively. Then the second **while** loop in

COMPARE-C3, unlike COMPARE, only compares the suffixes  $s(\mathbf{x}_h)$  and  $s(\mathbf{y}_h)$ , to compare strings  $\mathbf{x}$  and  $\mathbf{y}$ . Therefore, it terminates when it encounters  $\mathcal{L}\mathbf{x}$  or  $\varepsilon$ , which marks the end of substrings  $\mathbf{x}_h$  and  $\mathbf{y}_h$  as seen in Observation 2. The two **if** statements before the second **while** loop ensure that COMPARE-C3 returns correct results when either  $s(\mathbf{x}_h) = \varepsilon$  or  $s(\mathbf{y}_h) = \varepsilon$ . For completeness we give the pseudocode for SCAN<sup>2</sup> (Figure 3), which is used in COMPARE-Sensitive to scan the string  $\mathbf{x}$  ( $\mathbf{y}$ ) to compute the maximum  $\mathcal{L}\mathbf{x}$  ( $\mathcal{L}\mathbf{y}$ ) and its frequency  $occx$  ( $occy$ ).

If the conditions (C1) and (C2) of Lemma 3 do not hold, then we need to execute COMPARE-C3 to compare strings  $\mathbf{x}$  and  $\mathbf{y}$ . For such strings, scanning them and computing  $\mathcal{L}\mathbf{x}$  ( $\mathcal{L}\mathbf{y}$ ) and  $occx$  ( $occy$ ) might result in an overhead. Therefore, such strings are possibly one of the worst case input strings for which the time required by COMPARE-Sensitive is the maximum. Therefore, we refer to them as “bad strings”.

**Lemma 3** *The probability of choosing a pair of bad strings  $\mathbf{x}$  and  $\mathbf{y}$  of length  $n$  from an alphabet of size  $\sigma$  is*

$$\frac{\sum_{\mathcal{L}\mathbf{x}=1}^{\sigma} \sum_{k=1}^n \left( \binom{n}{k} (\mathcal{L}\mathbf{x} - 1)^{n-k} \right)^2}{\sigma^{2n}}, \quad (1)$$

where  $\mathcal{L}\mathbf{x}$  is the maximum letter in  $\mathbf{x}$ , and  $k = occx$  is the number of times  $\mathcal{L}\mathbf{x}$  occurs in  $\mathbf{x}$ .

*Proof.*

For a given  $\mathcal{L}\mathbf{x}$  and  $k$  the number of strings of length  $n$  is computed as follows:

1. the  $k$  positions where  $\mathcal{L}\mathbf{x}$  occurs can be chosen in  $\binom{n}{k}$  ways;
2. the remaining  $n - k$  positions can be chosen in  $(\mathcal{L}\mathbf{x} - 1)^{n-k}$  ways.

Then for fixed values of  $\mathcal{L}\mathbf{x}$  and  $k$  the total number of ways to choose  $\mathbf{x}$  is

$$\binom{n}{k} (\mathcal{L}\mathbf{x} - 1)^{n-k}.$$

Since we choose  $\mathbf{y}$  independently of  $\mathbf{x}$ , and  $\mathcal{L}\mathbf{x} = \mathcal{L}\mathbf{y}$ , the number of ways in which we can choose  $\mathbf{y}$  for fixed values of  $\mathcal{L}\mathbf{x}$  and  $k$  is also  $\binom{n}{k} (\mathcal{L}\mathbf{x} - 1)^{n-k}$ .

Therefore the total number of ways to choose a pair of bad strings  $\mathbf{x}$  and  $\mathbf{y}$  for fixed values of  $\mathcal{L}\mathbf{x}$  and  $k$  is

$$\left( \binom{n}{k} (\mathcal{L}\mathbf{x} - 1)^{n-k} \right)^2.$$

Since  $\mathcal{L}\mathbf{x} \in [1..\sigma]$  and  $k \in [1..n]$ , the total number of ways to choose a pair of bad strings  $\mathbf{x}$  and  $\mathbf{y}$  is

$$\sum_{\mathcal{L}\mathbf{x}=1}^{\sigma} \sum_{k=1}^n \left( \binom{n}{k} (\mathcal{L}\mathbf{x} - 1)^{n-k} \right)^2.$$

The total number of ways in which we can choose two strings over an alphabet of size  $\sigma$  is  $\sigma^{2n}$ . Therefore the probability of choosing a bad pair of strings  $\mathbf{x}$  and  $\mathbf{y}$  is

<sup>2</sup> For reasons of efficiency, the implementation of SCAN does not necessarily conform to the pseudocode given in (Figure 3).

$$\frac{\sum_{\mathcal{L}\mathbf{x}=1}^{\sigma} \sum_{k=1}^n \left( \binom{n}{k} (\mathcal{L}\mathbf{x} - 1)^{n-k} \right)^2}{\sigma^{2n}}.$$

□

To simplify the computation in Lemma 3, we assume that  $|x| = |y| = n$ ; that is, the strings to be compared are of the same lengths. However, Lemma 3 can easily be extended where  $|x| \neq |y|$ .

We computed the probabilities for the length of strings ( $n$ ) ranging from 1 to 100 over alphabets of size  $\sigma = 2, 4, 20$ , and plotted graphs (see Appendix A) for the same. From the graphs we see that the probability of choosing a pair of bad strings  $\mathbf{x}$  and  $\mathbf{y}$  reaches a maximum value, and then drops significantly and stabilizes to a constant value as  $n$  approaches 100. Since for all practical purposes the strings of interest are very large ( $\gg 100$ ), we conclude that **COMPARE-Sensitive** will run faster than **COMPARE** in all cases of interest.

### 3.2 Experimental results

In this section, we show the results of the experiments conducted to compare the performance of **COMPARE** and **COMPARE-Sensitive** on different classes of bad strings. As expected from the Lemma 3, **COMPARE-Sensitive** runs much faster than **COMPARE**. In fact from the experimental results we see that **COMPARE-Sensitive** is twice as fast as **COMPARE**.

The experiments were conducted on a Windows 10 64-bit Operating System, with Intel(R) Xeon(R) CPU E31245 v3 @ 3.40GHz x64-based processor, having an installed memory (RAM) of 32.0 GB. The code was implemented in the C++ language using Visual Studio 2017.

Strings  $\mathbf{x}$  in the sample pair of strings  $(\mathbf{x}, \mathbf{y})$  were chosen from the sample strings found at:

<http://www.cas.mcmaster.ca/~bill/strings/>,

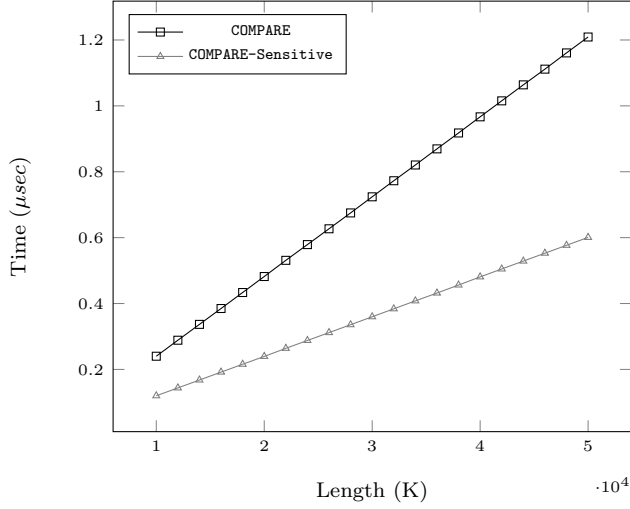
and the string  $\mathbf{y}$  was generated as a permutation of  $\mathbf{x}$ . Since  $\mathbf{y}$  is a permutation of  $\mathbf{x}$ , the strings  $\mathbf{x}$  and  $\mathbf{y}$  fail conditions (C1) and (C2) of Lemma 3. Therefore each pair  $(\mathbf{x}, \mathbf{y})$  is a pair of bad strings. In particular, strings  $\mathbf{x}$  were chosen from the sample DNA, protein, random ( $\sigma = 2, \sigma = 21$ ) and highly periodic strings available at the above URL. The lengths of the strings in a pair range from  $n = 10K$  to  $50K$ . To minimize the effects of external factors (for example delays caused due to interrupts etc.) on the experiments, we executed the algorithms **COMPARE** and **COMPARE-Sensitive** on the same pair of strings ten times, and used the *minimum* time taken by each of them. In addition to this, we take the average of the time taken for 100 different pairs, to get the final data point for comparison in the graphs.

Let  $T_c$  and  $T_{cs}$  be the slopes of the **COMPARE** and **COMPARE Sensitive** lines seen in the graphs. Let

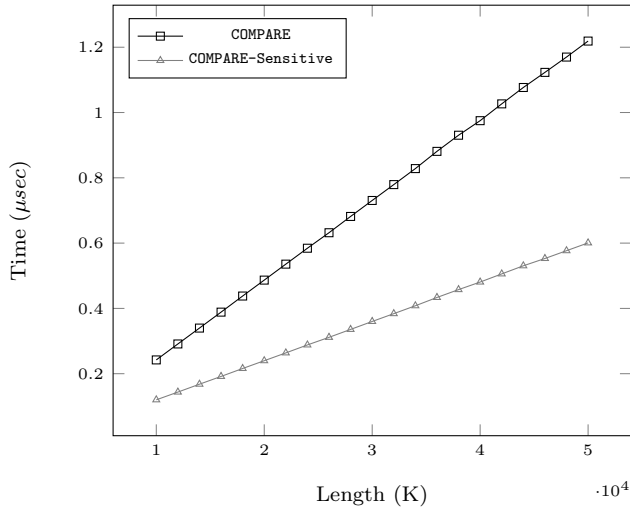
$$\alpha = \frac{T_c}{T_{cs}}.$$

Then the  $\alpha$  values computed from the graphs for DNA strings (see Figure 5), random strings over alphabets 2 and twenty one (see Figures 6, 7), protein strings

(see Figure 8), and highly periodic strings (see Figure 9) are:  $\alpha_{DNA} = 2.023$ ,  $\alpha_{rand2} = 2.031$ ,  $\alpha_{rand21} = 1.984$ ,  $\alpha_{protein} = 1.954$ , and  $\alpha_{hp} = 2.009$ , respectively. These  $\alpha$  values suggest that **COMPARE-Sensitive** is twice as fast as **COMPARE**. Extrapolating for a value of  $n = 1,000,000,000$  (1 billion) for a DNA string, the time taken by **COMPARE** is  $24.44 \cdot 10^{-3}$  secs while that for **COMPARE-Sensitive** is  $12.06 \cdot 10^{-3}$  secs. Moreover, it is seen that the size  $\sigma$  of the alphabet does not have any discernible effect on  $\alpha$ .

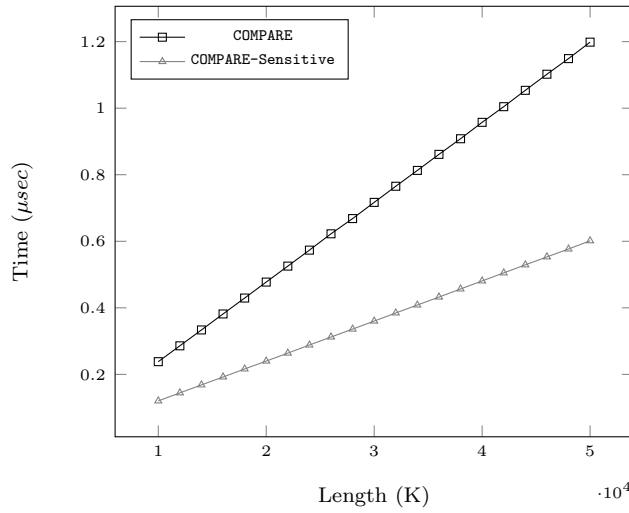


**Figure 5.** Runtime comparison of **COMPARE** and **COMPARE-Sensitive** on pairs of string  $(x, y)$ , where  $x$  is randomly selected from DNA strings of length  $2K$  to  $50K$  and  $\sigma = 4$ , and  $y$  is a permutation of  $x$ .

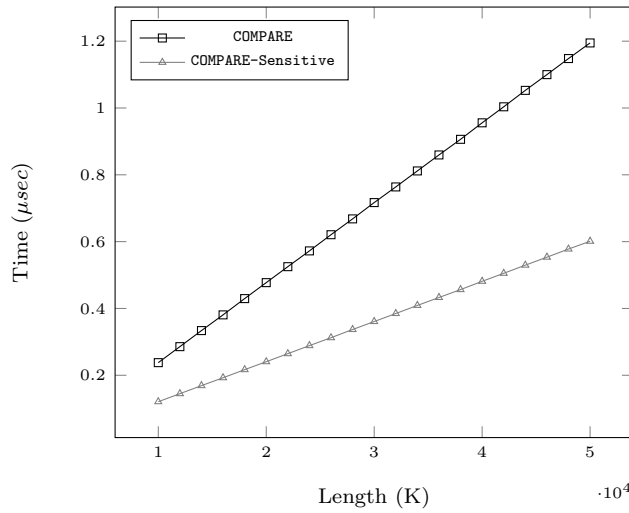


**Figure 6.** Runtime comparison of **COMPARE** and **COMPARE-Sensitive** on pairs of string  $(x, y)$ , where  $x$  is randomly selected from random strings of length  $2K$  to  $50K$  and  $\sigma = 2$ , and  $y$  is a permutation of  $x$ .

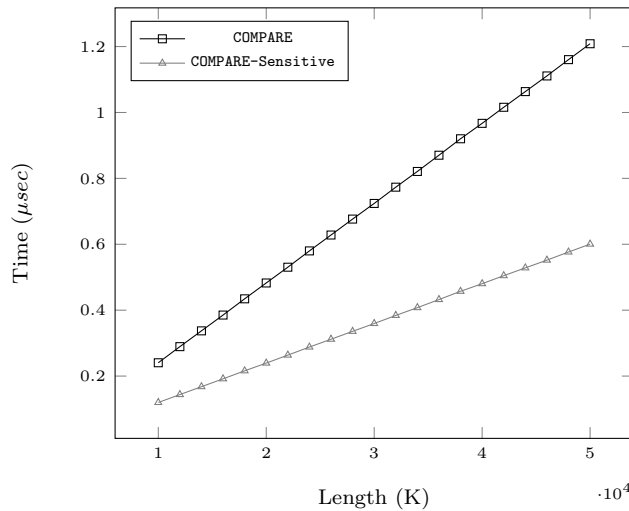




**Figure 7.** Runtime comparison of COMPARE and COMPARE-Sensitive on pairs of string  $(x, y)$ , where  $x$  is randomly selected from random strings of length  $2K$  to  $50K$  and  $\sigma = 21$ , and  $y$  is a permutation of  $x$ .



**Figure 8.** Runtime comparison of COMPARE and COMPARE-Sensitive on pairs of string  $(x, y)$ , where  $x$  is randomly selected from protein strings of length  $2K$  to  $50K$  and  $\sigma = 20$ , and  $y$  is a permutation of  $x$ .



**Figure 9.** Runtime comparison of COMPARE and COMPARE-Sensitive on pairs of string  $(x, y)$  where  $x$  is randomly selected from highly periodic strings of length  $2K$  to  $50K$  and  $\sigma = 2$ , and  $y$  is a permutation of  $x$ .

## Acknowledgments

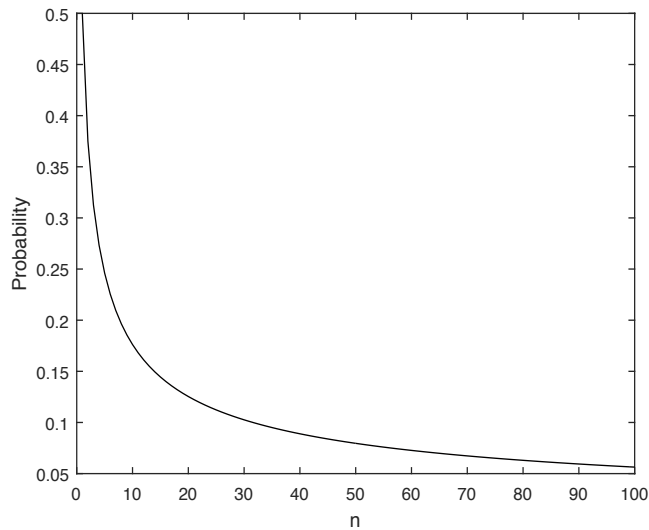


The third and fifth authors were funded by the NSERC Grant Number: 10536797. The second author was part-funded by the European Regional Development Fund through the Welsh Government.

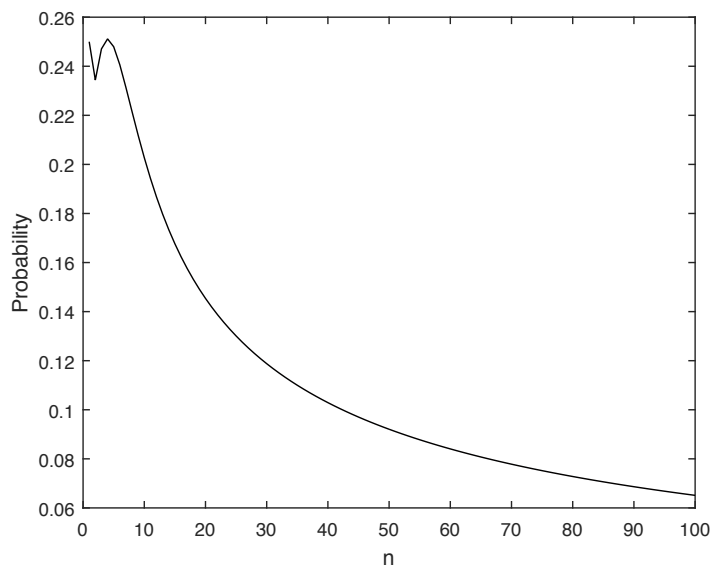
## References

1. A. ALATABBI, J. W. DAYKIN, J. KÄRKKÄINEN, M. S. RAHMAN, AND W. F. SMYTH: *V-order: New combinatorial properties & a simple comparison algorithm*. Discrete Appl. Math., 215 2016, pp. 41–46.
2. A. ALATABBI, J. W. DAYKIN, M. S. RAHMAN, AND 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.
3. A. ALATABBI, J. W. DAYKIN, M. S. RAHMAN, AND W. F. SMYTH: *Simple linear comparison of strings in V-order*. Fundam. Inform., 139(2) 2015, pp. 115–126.
4. A. ALATABBI, J. W. DAYKIN, M. S. RAHMAN, AND W. F. SMYTH: *String comparison in V-order: New lexicographic properties & on-line applications*. arXiv:1507.07038, 2015.
5. M. CROCHEMORE, C. HANCART, AND T. LECROQ: *Algorithms on Strings*, Cambridge University Press, New York, NY, USA, 2007.
6. T. N. DANH AND 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, pp. 43–53.
7. D. E. DAYKIN AND J. W. DAYKIN: *Lyndon-like and V-order factorizations of strings*. J. Discrete Algorithms, 1(3–4) 2003, pp. 357–365.
8. D. E. DAYKIN, J. W. DAYKIN, AND 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 (eds.), 2009, pp. 295–309.
9. D. E. DAYKIN, J. W. DAYKIN, AND 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.
10. D. E. DAYKIN, J. W. DAYKIN, AND W. F. SMYTH: *A linear partitioning algorithm for hybrid Lyndons using V-order*. Theoret. Comput. Sci., 483 2013, pp. 149–161.
11. J. W. DAYKIN AND W. F. SMYTH: *A bijective variant of the Burrows–Wheeler transform using V-order*. Theoret. Comput. Sci., 531 2014, pp. 77–89.
12. M. LOTHAIRE: *Combinatorics on Words*, Reading, MA (1983); 2nd Edition, Cambridge University Press, Cambridge (1997)., Addison–Wesley, 1983.
13. B. SMYTH: *Computing Patterns in Strings*, Pearson/Addison–Wesley, 2003.

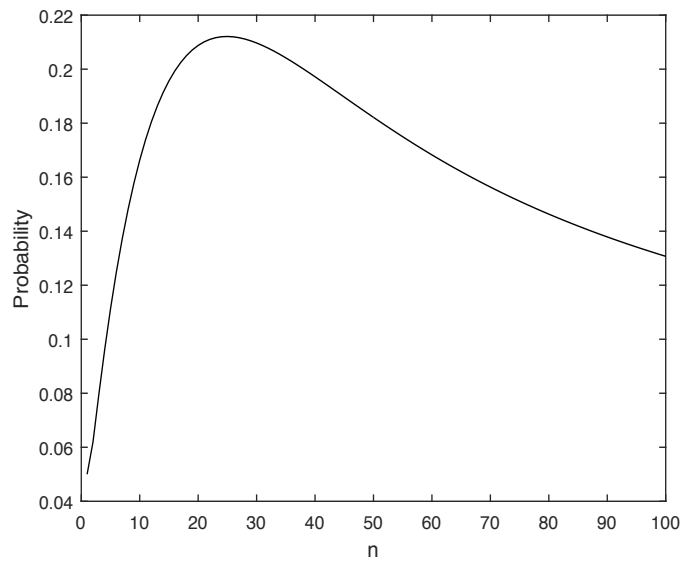
**A** Figures for experiments conducted to compute probabilities of choosing a pair of bad strings of length  $n$  from an alphabet of size  $\sigma$ .



**Figure 10.** Probabilities computed for strings over alphabet of size  $\sigma = 2$ , and length  $n \in [1..100]$  using Equation (1).



**Figure 11.** Probabilities computed for strings over alphabet of size  $\sigma = 4$ , and length  $n \in [1..100]$  using Equation (1).



**Figure 12.** Probabilities computed for strings over alphabet of size  $\sigma = 20$ , and length  $n \in [1..100]$  using Equation (1).