

String Range Matching

Juha Kärkkäinen, Dominik Kempa, and Simon J. Puglisi

Department of Computer Science, University of Helsinki Helsinki, Finland
`firstname.lastname@cs.helsinki.fi`

Abstract. Given strings X and Y the *exact string matching* problem is to find the occurrences of Y as a substring of X . An alternative formulation asks for the lexicographically consecutive set of suffixes of X that begin with Y . We introduce a generalization called *string range matching* where we want to find the suffixes of X that are in an arbitrary lexicographical range bounded by two strings Y and Z . The problem has applications in distributed suffix sorting, where Y and Z are themselves suffixes of X . Exact string matching can be solved in linear time and constant extra space under the standard comparison model. Our conjecture is that string range matching is a harder problem and cannot be solved within the same time–space complexity. In this paper, we trace the upper bound on the complexity of string range matching by describing algorithms that are within a logarithmic factor of the time–space complexity of exact string matching, as well as variants of the problem and the model that can be solved in linear time and constant extra space.

1 Introduction

Exact string matching, the problem of finding all the occurrences of a string $Y[0..m)$ (the pattern) in a larger string $X[0..n)$ (the text) is a foundational problem in computer science, and has applications throughout modern computer software. Among the numerous algorithms for exact string matching [5], there are several that are optimal in space ($O(1)$ extra space) as well as in time ($O(n+m)$) under the comparison model [9, 8, 3, 2, 4]. We have recently considered a generalization called *longest prefix matching*, where we want to find the occurrences of the longest prefix of the pattern that occurs in the text, and have shown that (at least) one of these algorithms can be generalized to longest prefix matching within the same optimal time–space complexity [11].

In this paper we introduce a further generalization based on an alternative view of the above problems. If we regard the text as a collection of its suffixes, the above problems can be stated as reporting the starting positions of all suffixes that begin with the pattern or with the longest prefix of the pattern producing a non-empty result. The resulting subset of suffixes is lexicographically consecutive and the query could be expressed as a lexicographical range query on the set of suffixes. A natural generalization, which we call *string range matching*, then asks for suffixes in an arbitrary lexicographical range: Given strings $X[0..n)$, $Y[0..m_1)$ and $Z[0..m_2)$ report all i such that $Y \leq X[i..n) < Z$.

We describe two basic algorithms for string range matching. One is based Crochemore’s string matching algorithm [2] (or its simplification in [11]) and solves the problem in $O(n \log(m_1 + m_2))$ time using constant extra space. For certain important special cases the algorithm can be made to run in linear time, in particular when the strings Y and Z share a prefix of length $\epsilon(m_1 + m_2)$ for any constant $\epsilon > 0$. The second algorithm is based on the string matching algorithm by Galil and Seiferas [9] and solves the *counting* version of the problem in linear time using $O(\log(m_1 + m_2))$ extra space. In both cases, the modification to solve string range matching is non-trivial.

Furthermore, we show that the problem can be solved in linear time using just constant extra space by slightly cheating about the extra space aspect. We describe two distinct “cheats”: (i) The algorithm can overwrite the strings Y and Z with other data but must restore the strings to their original state at the end. The overwritten memory does not count as extra space. (ii) The algorithm has read access to its own output, which does not count as extra space. The algorithm works for most reasonable output formats but only for the reporting version of the problem, not for the counting version.

We conjecture that string range matching is a harder problem than exact string matching and cannot be solved in linear time and constant extra space in the comparison model without cheating, i.e., when the only access to the input is by character comparisons and the output is write-only.

Applications to Suffix Sorting. The suffix array [13] and the Burrows–Wheeler transform [1] are central to modern string processing and the key task in their construction is sorting the suffixes of a string. For long strings we may want to split the task into smaller subtasks in order to distribute the work load or reduce the memory requirements. One approach partitions the lexicographical space [10]. For each partition, we then need to collect all suffixes within the given lexicographical range. The collecting problem is exactly string range matching.

Another approach to suffix sorting for large strings is to split the string into smaller blocks, sort the suffixes of each block separately, and then merge [6]. Assume the string is ABC , where B is the current block. When sorting the suffixes of B , what we really want is the lexicographical ordering of the $|B|$ longest suffixes of BC , which can differ significantly from the ordering of the suffixes of B as a standalone string. The correct ordering can be obtained without accessing C if we know which of the $|B|$ longest suffixes of BC are smaller than C . This information can be obtained using a one-sided string range matching query.

In both of the papers mentioned above [10, 6], the string range matching problem is solved using an algorithm introduced in [10], which is based on the Knuth–Morris–Pratt [12] exact string matching algorithm, and is the only prior work on string range matching that we are aware of. However, the $O(m_1 + m_2)$ extra space needed by the algorithm is a problem since Y and Z are potentially very long suffixes of X . The techniques used in those papers to reduce the space requirement (a global data structure [10] and sequential processing of the blocks [6]) are not well suited for distributed computation. Thus our small space algorithms open up new possibilities for distributed suffix sorting.

2 General Framework

Consider a string $X = X[0..n] = X[0]X[1] \dots X[n-1]$ of $|X| = n$ symbols drawn from an ordered alphabet. Here and elsewhere we use $[i..j)$ as a shorthand for $[i..j-1]$. For $i \in [0..n)$ we write X_i to denote the *suffix* of X of length $n-i$, that is $X_i = X[i..n) = X[i]X[i+1] \dots X[n-1]$. We also generalize the notation to sets of suffixes: for any $S \subseteq [0..n)$, $X_S = \{X_i \mid i \in S\}$. The *empty string* is denoted by ε . For strings X and Y , we use $\text{lcp}(X, Y)$ to denote the length of the *longest common prefix* of X and Y . If $X \leq Y$, we write $[X, Y)$ to denote the *lexicographical range* of strings between X and Y , i.e., the set $\{Z \mid X \leq Z < Y\}$. For $X \leq Y \leq Z$, we have $[X, Z) = [X, Y) \cup [Y, Z)$, where \cup denotes the *disjoint union*. A positive integer p is a *period* of X if $X[i] = X[i+p]$ for all $i \in [0..n-p)$. If p and q are periods of X and $p+q \leq |X|$, then $\text{gcd}(p, q)$ is a period of X too (Weak Periodicity Lemma [7]). The smallest period of X is denoted $\text{per}(X)$. A string X is called *primitive* if it cannot be written as $X = Y^k$ for an integer $k > 1$.

We focus on the one-sided string range matching problem of computing the set $X_{[0..n)} \cap [\varepsilon, Y)$, to which the two-sided version can be reduced since $[Y, Z) = [\varepsilon, Z) \setminus [\varepsilon, Y)$. All our algorithms use a similar basic approach. In a generic step, they compute $\ell = \text{lcp}(X_i, Y)$ and then i is incremented by h , where either $h = p = \text{per}(Y[0..\ell))$ or $h < p$ and $h = \Theta(\ell)$. The efficient computation of ℓ and h is based on well-known exact string matching algorithms. For exact string matching, nothing else is needed as none of the skipped suffixes X_{i+j} , $j \in [1..h)$, can have Y as a prefix. Our contribution is to show how to add order comparisons between Y and the suffixes. Checking whether $X_i < Y$ needs just one symbol comparison. Comparing the skipped suffixes to Y is based on the following lemma.

Lemma 1. *Let $\ell = \text{lcp}(X_i, Y)$ and $p = \text{per}(Y[0..\ell))$. Then $X_{i+j} < Y$ iff $Y_j < Y$ for any $j \in [1..p)$.*

Proof. Fix $j \in [1..p)$. If we had $\text{lcp}(X_{i+j}, Y) \geq \ell - j$, this would imply that j is a period of $Y[0..\ell)$, which violates the assumption that $p = \text{per}(Y[0..\ell))$. Thus $\text{lcp}(X_{i+j}, Y) < \ell - j$ and the claim follows since $\text{lcp}(X_{i+j}, Y_j) = \ell - j$. \square

In other words, the status of the suffixes in the skipped segment depends only on Y enabling the use of precomputed information.

3 Linear Time and Logarithmic Extra Space

Our first algorithm solves the one-sided counting variant of the string range matching problem, that is, computes the value of $|X_{[0..n)} \cap [\varepsilon, Y)|$. The two-sided problem needs two calls to the algorithm. The algorithm is built on the exact string matching algorithm of Galil and Seiferas that uses $O(\log m)$ extra space [9] (or more precisely on the cleaner formulation due to Crochemore and Rytter [4]).

Assume that $k \geq 3$ is an integer constant. A prefix P of Y is called a *k-highly repeating prefix* (*k-hrp*) of Y if P is primitive and P^k is a prefix of Y . With each *k-hrp* P of Y we associate the interval of positions $[2|P|, \text{lcp}(Y, Y_{|P|}))$ called the *scope* of P . Scopes have a number of useful properties:

Lemma 2 ([4]). *A string Y of length m has $O(\log m)$ scopes and they can be computed in $O(m)$ time and $O(\log m)$ extra space. Two distinct scopes do not overlap and for any prefix $Y[0..\ell]$ of Y it holds that:*

$$\begin{aligned} \text{per}(Y[0..\ell]) &= b/2 && \text{if } \ell \in [b, e) \text{ and } [b, e) \text{ is a scope of some } k\text{-hrp of } Y; \text{ and} \\ \text{per}(Y[0..\ell]) &> \ell/k && \text{otherwise.} \end{aligned}$$

Using the scopes, we can compute the skip value h in the generic matching step as either $b/2$ or $\lfloor \ell/k \rfloor + 1$, which is sufficient for exact string matching [4]. For string range matching, we need additional precomputed information. First, for each k -hrp P of Y we precompute, in addition to its scope $[b, e)$, the value $c = |Y_{[0..|P|]} \cap [\varepsilon, Y]|$. We store all triples (b, e, c) in a sorted list \mathcal{S}_p , which allows us to count all suffixes that are skipped-over whenever $h = \text{per}(Y[0..\ell])$. If the match length ℓ does not belong to a scope, the Galil–Seiferas algorithm would use a skip value $h = \lfloor \ell/k \rfloor + 1$. Our strategy in this case is to do a shorter (by a constant factor) skip h for which we know the value $|Y_{[0..h]} \cap [\varepsilon, Y]|$. To be able to do that, we will precompute a second list $\mathcal{S}_n = ((b_1, c_1), \dots, (b_t, c_t))$ of length $O(\log m)$ that satisfies: $c_i = |Y_{[0..b_i]} \cap [\varepsilon, Y]|$ for $i \in [1..t]$ and $2b_i \leq b_{i+1} < 4b_i$ for $i \in [1..t-1]$. Moreover $4b_t > m$.

Pseudo-code for the algorithm is given in Fig. 1. The algorithm precomputing lists \mathcal{S}_p and \mathcal{S}_n is given in Fig. 2. It is essentially the same as in Fig. 1, except we are now matching the pattern against itself.

Algorithm OneSidedStringRangeCounting(X, Y)

Input: strings $X[0..n), Y[0..m)$.

Output: $|X_{[0..n)} \cap [\varepsilon, Y]|$.

```

1:  $(\mathcal{S}_p, \mathcal{S}_n) \leftarrow \text{Precompute}(Y)$ 
2:  $count \leftarrow i \leftarrow \ell \leftarrow 0$ 
3: while  $i < n$  do // Invariant:  $count = |X_{[0..i)} \cap [\varepsilon, Y]|$ 
4:   while  $i + \ell < n$  and  $\ell < m$  and  $X[i + \ell] = Y[\ell]$  do  $\ell \leftarrow \ell + 1$ 
5:    $(b, e, c) \leftarrow \text{contains}(\mathcal{S}_p, \ell)$ 
6:   if  $\ell < m$  and  $(i + \ell = n$  or  $X[i + \ell] < Y[\ell])$  then
7:      $count \leftarrow count + 1$ 
8:   if  $b \neq 0$  then //  $\text{per}(Y[0..\ell]) = b/2$ 
9:     //  $c = |Y_{[1..b/2)} \cap [\varepsilon, Y]| = |X_{[i+1..i+b/2)} \cap [\varepsilon, Y]|$ 
10:     $count \leftarrow count + c$ 
11:     $i \leftarrow i + b/2; \ell \leftarrow \ell - b/2$ 
12:   else //  $\text{per}(Y[0..\ell]) > \ell/k$ 
13:      $(b, c) \leftarrow \text{pred}(\mathcal{S}_n, \lfloor \ell/k \rfloor + 1)$  //  $(\lfloor \ell/k \rfloor + 1)/4 < b$ 
14:     //  $c = |Y_{[1..b)} \cap [\varepsilon, Y]| = |X_{[i+1..i+b)} \cap [\varepsilon, Y]|$ 
15:      $count \leftarrow count + c$ 
16:      $i \leftarrow i + b; \ell \leftarrow 0$ 
17: return  $count$ 

```

Fig. 1. String range counting in linear time and logarithmic space.

Algorithm Precompute(Y)

Input: a string $Y[0..m)$.
Output: lists \mathcal{S}_p and \mathcal{S}_n .

```

1:  $\mathcal{S}_p \leftarrow ()$ ;  $\mathcal{S}_n \leftarrow ((1, 0))$ 
2:  $i \leftarrow last \leftarrow 1$ ;  $\ell \leftarrow count \leftarrow 0$ 
3: while  $i < m$  do // Invariant:  $count = |Y_{[0..i)} \cap [\varepsilon, Y)|$ 
4:   while  $i + \ell < m$  and  $Y[i + \ell] = Y[\ell]$  do  $\ell \leftarrow \ell + 1$ 
5:    $(b, e, c) \leftarrow \text{contains}(\mathcal{S}_p, \ell)$ 
6:   if  $k \cdot i \leq i + \ell$  and  $b = 0$  then
7:      $(b, e, c) \leftarrow (2i, i + \ell, count)$ 
8:      $\text{add}(\mathcal{S}_p, (b, e, c))$  //  $Y[0..i)$  is a new  $k$ -hrp
9:   if  $2 \cdot last \leq i$  then
10:     $\text{add}(\mathcal{S}_n, (i, count))$ ;  $last \leftarrow i$ 
11:   if  $i + \ell = m$  or  $Y[i + \ell] < Y[\ell]$  then  $count \leftarrow count + 1$ 
12:   if  $b \neq 0$  then
13:      $count \leftarrow count + c$ 
14:      $i \leftarrow i + b/2$ ;  $\ell \leftarrow \ell - b/2$ 
15:   else
16:      $(b, c) \leftarrow \text{pred}(\mathcal{S}_n, \lfloor \ell/k \rfloor + 1)$ 
17:      $count \leftarrow count + c$ 
18:      $i \leftarrow i + b$ ;  $\ell \leftarrow 0$ 
19: return  $(\mathcal{S}_p, \mathcal{S}_n)$ 

```

Fig. 2. Computation of lists \mathcal{S}_p and \mathcal{S}_n .

The algorithms use the following additional procedures: $\text{add}(L, x)$ adds an element x at the end of a list L , $\text{pred}(\mathcal{S}_n, x)$ returns $(b, c) \in \mathcal{S}_n$, $b \leq x$ with the largest b , and $\text{contains}(\mathcal{S}_p, x)$ returns $(b, e, c) \in \mathcal{S}_p$ such that $b \leq x < e$ or $(0, \cdot, \cdot)$ if no such triple belongs to \mathcal{S}_p .

Theorem 1. *The algorithm in Figure 1 solves the counting version of the string range matching problem in linear time and $O(\log(m_1 + m_2))$ extra space.*

Proof. The correctness of the algorithm follows from Lemmas 1 and 2 as explained above and detailed in the invariants included in the code. The time complexity is linear by the same arguments as in [4]. In particular, the shifts made by our algorithm may be shorter but only by a constant factor. \square

4 $O(n \log(m_1 + m_2))$ Time and Constant Extra Space

In this section, we first describe an algorithm solving the string range matching problem in linear time and constant extra space for the special case, where the query range is of the form $[Y[0..r), Y)$ with $\lfloor 2m/3 \rfloor \leq r < m = |Y|$. We will then show how the general string range matching problem can be solved using $O(\log m)$ calls to the special case algorithm.

The backbone of the algorithm is the exact string matching algorithm of Crochemore [2] (see [11] for a simplified version, which we use here). Similarly

to the Galil–Seiferas algorithm, it computes the skip value h in the generic step either as $h = \text{per}(\mathbf{Y}[0..l]) \leq \ell/3$ or $h = \lfloor \ell/3 \rfloor + 1 \leq \text{per}(\mathbf{Y}[0..l])$, but now using only $O(1)$ extra space. Our addition is two precomputed values $q = \text{per}(\mathbf{Y}[0..r])$ and $e = q + \text{lcp}(\mathbf{Y}, \mathbf{Y}_q)$, i.e., $\mathbf{Y}[0..e]$ is the longest prefix of \mathbf{Y} having period q . Both values can be computed in linear time using constant extra space [2]. Determining of the status of the skipped suffixes is based on Lemma 1 and the following result. The pseudo-code is given in Figure 3.

Lemma 3. *Let $h \leq \lfloor m/3 \rfloor + 1$ be a positive integer and $g = \min(\lceil \frac{h-q}{q} \rceil, \lfloor \frac{e-r}{q} \rfloor)$. Then:*

$$\{j \in [0..h] \mid \mathbf{Y}[0..r] \leq \mathbf{Y}_j < \mathbf{Y}\} = \begin{cases} \{q, 2q, \dots, gq\} & \text{if } \mathbf{Y}_q < \mathbf{Y} \\ \emptyset & \text{otherwise} \end{cases}$$

Proof. First note that $\mathbf{Y}[0..r] \leq \mathbf{Y}_j < \mathbf{Y}$ iff $\mathbf{Y}_j < \mathbf{Y}$ and $\text{lcp}(\mathbf{Y}_j, \mathbf{Y}) \geq r$. If $\text{lcp}(\mathbf{Y}_j, \mathbf{Y}) \geq r$, then j is a period of $\mathbf{Y}[0..r]$. Since $q = \text{per}(\mathbf{Y}[0..r])$ and $j + q \leq 2j \leq 2\lfloor m/3 \rfloor \leq r$, the Weak Periodicity Lemma implies that j must be a multiple of q . Thus let $j = kq$ for some integer k such that $kq < h \leq e$. Then $\text{lcp}(\mathbf{Y}_j, \mathbf{Y}) = e - j$ and $\mathbf{Y}_j < \mathbf{Y}$ iff $\mathbf{Y}_q < \mathbf{Y}$. Thus, if $\mathbf{Y}_q \geq \mathbf{Y}$, the set in the lemma is empty. Otherwise, the set includes all $j = kq$ such that $kq < h$ and $\text{lcp}(\mathbf{Y}_j, \mathbf{Y}) \geq r$, which give the limits $k \leq \lceil (h - q)/q \rceil$ and $k \leq \lfloor (e - r)/q \rfloor$, respectively. \square

Algorithm RestrictedOneSidedStringRangeReporting(X, Y, r)

Input: strings $X[0..n]$, $Y[0..m]$ and an integer $r \geq \lfloor 2m/3 \rfloor$.

Output: the set $\mathcal{R} = \{j \in [0..n] \mid \mathbf{Y}[0..r] \leq X_j < \mathbf{Y}\}$.

```

1:  $(q, e) \leftarrow \text{Precompute}(\mathbf{Y}, r)$ 
2:  $\mathcal{R} \leftarrow \emptyset$ 
3:  $i \leftarrow \ell \leftarrow s \leftarrow p \leftarrow 0$ 
4: while  $i < n$  do // Invariant:  $\mathcal{R} = \{j \in [0..i] \mid \mathbf{Y}[0..r] \leq X_j < \mathbf{Y}\}$ 
5:   while  $i + \ell < n$  and  $\ell < m$  and  $X[i + \ell] = Y[\ell]$  do
6:      $\text{UpdateMS}(\mathbf{Y}, \ell, s, p)$  // Increments  $\ell$  by one
7:     if  $r \leq \ell < m$  and  $(i + \ell = n$  or  $X[i + \ell] < Y[\ell])$  then  $\mathcal{R} \leftarrow \mathcal{R} \cup \{i\}$ 
8:     if  $0 < p \leq \ell/3$  and  $\mathbf{Y}[0..s] = \mathbf{Y}[p..p + s]$  then //  $\text{per}(\mathbf{Y}[0..l]) = p$ 
9:        $h \leftarrow p; \ell \leftarrow \ell - p$ 
10:    else //  $\text{per}(\mathbf{Y}[0..l]) > \ell/3$ 
11:       $h \leftarrow \lfloor \ell/3 \rfloor + 1; (\ell, s, p) \leftarrow (0, 0, 0)$ 
12:    if  $e < m$  and  $Y[e] < Y[e \bmod q]$  then //  $\mathbf{Y}_q < \mathbf{Y}$ 
13:       $g \leftarrow \min(\lceil \frac{h-q}{q} \rceil, \lfloor \frac{e-r}{q} \rfloor)$ 
14:       $\mathcal{R} \leftarrow \mathcal{R} \cup \{i + q, i + 2q, \dots, i + gq\}$ 
15:     $i \leftarrow i + h$ 
16: return  $\mathcal{R}$ 

```

Fig. 3. Reporting suffixes $X_j \in [Y[0..r], Y)$ in linear time and constant extra space. Due to lack of space, we omit the code for $\text{UpdateMS}(\mathbf{Y}, \ell, s, p)$, which can be found in [11], and the precomputation of q and e .

Let us now consider more general forms of queries. First, consider the query $[Y[0..\ell), Y)$ with $\ell < r = \lfloor 2|Y|/3 \rfloor$. We can break the query into two distinct sub-queries: $[Y[0..\ell), Y) = [Y[0..\ell), Y[0..r)) \cup [Y[0..r), Y)$. This leads to a tail-recursive algorithm that makes $O(\log(|Y|/\ell))$ calls to the algorithm in Fig. 3 and thus runs in $O(n \log(|Y|/\ell))$ time still using only constant extra space. An alternative formulation for the query $[Y[0..\ell), Y)$ is to compute the set $\{X_j \mid X_j < Y \text{ and } \text{lcp}(X_j, Y) \geq \ell\}$. By a symmetrical procedure, we can also compute the set $\{X_j \mid Y \leq X_j \text{ and } \text{lcp}(X_j, Y) \geq \ell\}$ in $O(n \log(|Y|/\ell))$ time and constant extra space.

Consider now the fully general query $[Y, Z)$ and let $\ell = \text{lcp}(Y, Z)$. We can partition the result set $X_{[0..n)} \cap [Y, Z)$ into three disjoint sets $\mathcal{R}_1 = \{X_j \mid X_j < Z \text{ and } \text{lcp}(X_j, Z) > \ell\}$, $\mathcal{R}_2 = \{X_j \mid Y \leq X_j \text{ and } \text{lcp}(X_j, Y) > \ell\}$ and $\mathcal{R}_3 = \{X_j \mid Y \leq X_j < Z \text{ and } \text{lcp}(X_j, Y) = \text{lcp}(X_j, Z) = \ell\}$. The sets \mathcal{R}_1 and \mathcal{R}_2 can be computed as described above. \mathcal{R}_3 can be found using any time-space optimal exact string matching algorithm (e.g. [2]). We have proven the following result.

Theorem 2. *It is possible to solve the string range matching problem using constant extra space in $O(n \log((m_1 + m_2)/(1 + \text{lcp}(Y, Z))))$ time.*

The worst case time complexity is $O(n \log(m_1 + m_2))$. If $\text{lcp}(Y, Z) \geq \epsilon(m_1 + m_2)$ for some constant $\epsilon > 0$, the algorithm runs in linear time.

5 Linear Time and Constant Extra Space

In this section, we describe two algorithms for string range matching that run in linear time and constant extra space by taking advantage of the space reserved for input or output, which we do not count as extra space. Both algorithms are based on Crochemore's string matching algorithm.

We will describe the algorithms for one-sided queries $[\epsilon, Y)$. Since these algorithms produce their output in a sequential order, we can answer two-sided queries by running two one-sided queries in parallel: $[Y, Z) = [\epsilon, Z) \setminus [\epsilon, Y)$. To ease the description of the algorithms, we assume that the output is produced as a bitvector $B[0..n)$ satisfying, for all $i \in [0..n)$, $B[i] = 1$ iff $X_i < Y$. We will later describe how to produce output in other formats.

Copying Output. Consider the actions of Crochemore's algorithm as it compares the pattern Y against the text suffix X_i . It computes $\ell = \text{lcp}(X_i, Y)$ and then shifts the pattern forward by h positions. With one character comparison we can determine $B[i]$, so the problem now is computing the output for the skipped-over positions $B[i + 1..i + h)$.

Let ℓ_{\max} be the longest previously found match, which occurred at position j , i.e., $\ell_{\max} = \text{lcp}(X_j, Y)$. If $\ell \leq \ell_{\max}$, then by Lemma 1 we have $B[i + 1..i + h) = B[j + 1..j + h)$. Thus, if we have access to the previous output, we can simply copy the missing output. Note that this works even if the two intervals $[i + 1..i + h)$ and $[j + 1..j + h)$ overlap.

If $\ell > \ell_{\max}$, we cannot obtain all the missing values by copying. In this case, we set ℓ to ℓ_{\max} (after computing $B[i]$) and continue as if the match had been a shorter one. This shorter match length leads to a shorter shift and allows copying the output for all the skipped-over positions. Shortening the match length does not violate the correctness of the algorithm as long as the state of algorithm is adjusted accordingly. The full match length becomes the new value of ℓ_{\max} , i.e., we are actually swapping the values of ℓ and ℓ_{\max} .

Algorithm OneSidedStringRangeReporting(X, Y)

Input: strings $X[0..n)$ and $Y[0..m)$.

Output: bitvector $B[0..n)$, where $B[i] = 1$ iff $X[i..n) < Y$.

```

1:  $B \leftarrow (0, 0, \dots, 0)$ 
2:  $i \leftarrow \ell \leftarrow p \leftarrow s \leftarrow 0$ 
3:  $i_{\max} \leftarrow \ell_{\max} \leftarrow p_{\max} \leftarrow s_{\max} \leftarrow 0$ 
4: while  $i < n$  do
5:   while  $i + \ell < n$  and  $\ell < m$  and  $X[i + \ell] = Y[\ell]$  do
6:      $(\ell, s, p) \leftarrow \text{UpdateMS}(Y, \ell, s, p)$  // Increments  $\ell$  by one
     //  $\ell = \text{lcp}(X_i, Y)$ 
7:   if  $\ell < m$  and ( $i + \ell = n$  or  $X[i + \ell] < Y[\ell]$ ) then  $B[i] = 1$ 
8:    $j \leftarrow i_{\max}$ 
9:   if  $\ell > \ell_{\max}$  then
10:    swap  $(\ell, s, p)$  and  $(\ell_{\max}, s_{\max}, p_{\max})$ 
11:     $i_{\max} \leftarrow i$ 
12:   if  $0 < p \leq \ell/3$  and  $Y[0..s) = Y[p..p + s)$  then //  $\text{per}(Y[0..\ell]) = p$ 
13:      $B[i + 1..i + p) \leftarrow B[j + 1..j + p)$ 
14:      $i \leftarrow i + p$ 
15:      $\ell \leftarrow \ell - p$ 
16:   else //  $\text{per}(Y[0..\ell]) > \ell/3$ 
17:      $h \leftarrow \lfloor \ell/3 \rfloor + 1$ 
18:      $B[i + 1..i + h) \leftarrow B[j + 1..j + h)$ 
19:      $i \leftarrow i + h$ 
20:      $(\ell, s, p) \leftarrow (0, 0, 0)$ 
21: return  $B$ 

```

Fig. 4. Linear time, constant extra space algorithm.

The pseudocode for the algorithm is given in Fig. 4. The key changes to the original algorithm in [11, Figure 2] are the computation of B on lines 7, 13 and 18, and the swap of ℓ and ℓ_{\max} on line 10. The two other variables, s and p , swapped together with ℓ represent the rest of the state of the algorithm that needs to change when ℓ changes. We refer to [11] for their explanation. Here it suffices to know that they are functions of ℓ (and Y).

Theorem 3. *The algorithm in Fig. 4 solves the string range reporting problem in linear time using constant extra space.*

Proof. The correctness of the algorithm has been established above, and the algorithm clearly uses only constant extra space. In [11], it is shown that the time spent in any round of the original algorithm is proportional to the increase of the value $6i + \ell + s$ during that round. Because of the swap on line 11, this might not hold for the modified algorithm, but this can be corrected by using the value $6i + \ell + s + \ell_{\max} + s_{\max}$ instead. The copying of the output does not add more than $O(n)$ time. Thus the time complexity is still $O(n + m)$. \square

The algorithm can be easily modified to handle other output formats such as a sequence of starting positions. On the other hand, the algorithm cannot be used if the output is streamed or compressed in such a way that it cannot be quickly decompressed from an arbitrary position. Also, the counting version of the problem cannot be solved this way without extra space.

Overwriting Input. Suppose we have separate storage space of m bits and consider the following variant of the algorithm in Fig. 4. Whenever i_{\max} changes we write $B[i_{\max} + 1..i_{\max} + m]$ to our separate storage at the same time as it is being written to the output. Whenever we need to copy a part of the output, we can copy from our separate storage instead of from the output itself. Now the algorithm never needs to access the output, and so can stream the output, encode it in any way, or just count the number of occurrences without producing any output until the end.

Let $\hat{\ell}$ be the length of the longest prefix of Y that also occurs elsewhere in Y . We will use the space occupied by $Y[0..\hat{\ell}]$ as a separate storage of at least $\hat{\ell}$ bits. Any access to the characters of $Y[0..\hat{\ell}]$ are redirected to that other occurrence. Clearly this can be made to work even if that other occurrence overlaps $Y[0..\hat{\ell}]$. The value $\hat{\ell}$ is easily computed in $O(m)$ time using $O(1)$ extra space.

The largest number of bits the algorithm may need to copy from the separate storage at a given stage is at most $\ell_{\max}/3$, since the longest possible shift is $\lfloor \ell_{\max}/3 \rfloor + 1$. Thus as long as $\ell_{\max} \leq 3\hat{\ell}$, the algorithm can operate as described.

Now consider a round i of the algorithm with $\ell = \text{lcp}(X_i, Y) > 3\hat{\ell}$. We call this a *long match*. In this case, the algorithm behaves as it would for a shorter match except the value of ℓ_{\max} (and the associated s_{\max} and p_{\max}) is not modified. Thus ℓ_{\max} never grows too large for the separate storage. The extra time spent in round i because of this is at most $O(\ell)$, and the total extra time is bounded by the total length of the long matches. Thus the following lemma shows that the time complexity of the algorithm is still $O(n)$.

Lemma 4. *The total length of all long matches is at most $1.5n$.*

Proof. Consider a long match of length ℓ at position i , i.e., $X[i..i + \ell) = Y[0..\ell)$, and assume that the next long match starts at $j > i$ and has length ℓ' . Then $\min\{\ell, \ell'\} > 3\hat{\ell}$. First observe that we must have $j + \ell' > i + \ell$. Otherwise, $Y[0..\ell')$ has another occurrence as $Y[j - i..j - i + \ell')$ and $\hat{\ell} \geq \ell'$. Furthermore, we must have $j > i + \ell - \lceil \ell/3 \rceil$. Otherwise, $Y[0..\lceil \ell/3 \rceil)$ has another occurrence as $Y[j - i..j - i + \lceil \ell/3 \rceil)$ and $3\hat{\ell} \geq \ell$. Thus at most one third of a long match can be overlapped by a later long match, which proves the lemma. \square

We have proven the following.

Theorem 4. *The string range matching problem can be solved in linear time and constant extra space if each character of the range boundary strings occupies at least one bit and that bit can be overwritten by the algorithm.*

Note that the algorithm can restore the overwritten part of Y at the end as long as equal characters are always represented by identical bit patterns. The restoration might not be possible in some cases, for example, when upper and lower case letters are considered equal in comparisons.

6 Concluding Remarks

This article is the first to directly consider the string range matching problem. Many interesting avenues for future work remain, but perhaps the most challenging open problem is to establish if the algorithms described in this paper for the general setting (read-only input, inaccessible output) are optimal, or if linear-time constant extra-space methods indeed exist.

References

1. M. Burrows and D.J. Wheeler. A block sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, Palo Alto, California, 1994.
2. M. Crochemore. String-matching on ordered alphabets. *Theor. Comp. Sci.*, 92:33–47, 1992.
3. M. Crochemore and D. Perrin. Two-way string matching. *J. ACM*, 38(3):651–675, 1991.
4. M. Crochemore and W. Rytter. Squares, cubes, and time-space efficient string searching. *Algorithmica*, 13(5):405–425, 1995.
5. S. Faro and T. Lecroq. The exact online string matching problem: A review of the most recent results. *ACM Comp. Surv.*, 45(2):13, 2013.
6. P. Ferragina, T. Gagie, and G. Manzini. Lightweight data indexing and compression in external memory. *Algorithmica*, 63(3):707–730, 2012.
7. N. J. Fine and H. S. Wilf. Uniqueness theorems for periodic functions. *Proc. Amer. Math. Soc.*, 16(1):109–114, 1965.
8. Z. Galil and J. Seiferas. Time-space optimal string matching. *J. Comp. Sys. Sci.*, 26:280–294, 1983.
9. Z. Galil and J. I. Seiferas. Saving space in fast string-matching. *SIAM J. Comp.*, 9(2):417–438, 1980.
10. J. Kärkkäinen. Fast BWT in small space by blockwise suffix sorting. *Theor. Comp. Sci.*, 387(3):249–257, 2007.
11. J. Kärkkäinen, D. Kempa, and S. J. Puglisi. Crochemore’s string matching algorithm: Simplification, extensions, applications. In *Proc. PSC 2013*, pages 168–175, Czech Technical University, 2013.
12. D. Knuth, J. H. Morris, and V. Pratt. Fast pattern matching in strings. *SIAM J. Comp.*, 6(2):323–350, 1977.
13. U. Manber and G. W. Myers. Suffix arrays: a new method for on-line string searches. *SIAM J. Comp.*, 22(5):935–948, 1993.