

Sorting suffixes of two-pattern strings

Frantisek Franek and W. F. Smyth*

Algorithms Research Group
Department of Computing & Software
McMaster University
Hamilton, Ontario
Canada L8S 4L7

e-mail: {franek, smyth}@mcmaster.ca

Abstract. Recently, several authors presented linear recursive algorithms for sorting suffixes of a string. All these algorithms employ a similar three-step approach, based on an initial division of the suffixes of x into two sets: in step 1 sort the first set using recursive reduction of the problem, in step 2 determine the order of the suffixes in the second set based on the order of the suffixes in the first set, and in step 3 merge the two sets together. To optimize such an algorithm either for space or time, it may not be sufficient to optimize one of the three steps, since in doing so, one might increase the resources required for the others to an unacceptable extent.

Franek, Lu, and Smyth introduced two-pattern strings as a generalization of Sturmian strings. Like Sturmian strings, two-pattern strings are generated by iterated morphisms, but they exhibit a much richer structure.

In this paper we show that the suffixes of two-pattern strings can be sorted in linear time using a variant of the three step approach outlined above. It turns out that, given the order of the suffixes in a two-pattern string, one can almost directly list in linear time all the suffixes of its expansion under a two-pattern morphism.

1 Introduction

Ever since Manber and Myers in [MM93] introduced suffix arrays as data structures comparable to suffix trees for most pattern matching tasks in strings, yet requiring significantly less memory, the search was on for a linear time algorithm for their construction. Such an algorithm for suffix tree construction had been known since 1997 [F97]. In 2003 to our knowledge three different groups of researchers independently proposed linear recursive algorithms to sort string suffixes: [KA03, KSPP03, KS03]. Though different, all three algorithms employ three steps, based on a separation of the suffixes into two sets. In step 1 the first set is ordered using recursive reduction of the problem, in step 2 the suffixes of the second set are sorted based on the order of the suffixes in the first set, and in step 3 both ordered sets are merged together. The fact that all three algorithms follow this basic approach, yet use a completely different

*also Department of Computing, Curtin University, Perth WA 6845, Australia.

separation into sets, a different way of ordering the second set based on the first set, and a different merge technique, points to some common fundamental aspect of these algorithms. To optimize such an algorithm either for space or time, it may not be sufficient to optimize one of the three steps, since in doing so, one might increase the resources required for the others to an unacceptable extent.

Two-pattern strings were introduced in [FLS03] as a generalization of Sturmian strings. Like Sturmian strings, two-pattern strings are generated by iterated morphisms, but they exhibit a much richer structure. It was shown in [FLS04] that the iterated construction of these strings could be used to compute all the repetitions and near-repetitions in time linear in string length.

This paper was motivated by our investigation of the three different linear suffix sorting algorithms discussed above and our desire to fully understand the underlying phenomena. Thus, we investigated whether the recursive nature of two-pattern strings could be used in sorting of the suffixes in the approach of the three algorithms mentioned. As it turned out, the “natural” recursive reduction of two-pattern strings can be used for step 1, and then steps 2 and 3 can be simplified into a single step: from having the suffixes of the reduced string ordered, one can almost directly list the suffixes of the two-pattern string in the right order.

For the sake of completeness, let us recall the definition of a two-pattern string (see [FLS03]), including all supporting definitions. Throughout this paper, a **binary string** means a string over the alphabet $\{a, b\}$.

Definition 1.1 *A binary string \mathbf{q} is said to be **p-regular** if and only if $\mathbf{q} = \mathbf{u}\mathbf{p}\mathbf{v}\mathbf{u}$ for some choice of (possibly empty) substrings \mathbf{u} and \mathbf{v} .*

Definition 1.2 *An ordered pair (\mathbf{p}, \mathbf{q}) of nonempty binary strings is said to be **suitable** if and only if*

- \mathbf{p} is **primitive** (that is, \mathbf{p} has no nonempty border);
- \mathbf{p} is not a suffix of \mathbf{q} ;
- \mathbf{q} is neither a prefix nor a suffix of \mathbf{p} ;
- \mathbf{q} is not **p-regular**.

Note: Since a two-pattern string is a concatenation of blocks $\mathbf{p}^i\mathbf{q}$ and $\mathbf{p}^j\mathbf{q}$, the above two definitions make sure that \mathbf{p} and \mathbf{q} are dissimilar enough to be recognized efficiently.

Definition 1.3 $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ is an **expansion of scope λ** , if (\mathbf{p}, \mathbf{q}) is suitable, $|\mathbf{p}| \leq \lambda$, $|\mathbf{q}| \leq \lambda$, $1 \leq i, j$, $i \neq j$ are integers, and λ is an integer ≥ 1 .

Note: An expansion $\sigma = [\mathbf{p}, \mathbf{q}, i, j]$ is applied to a binary string \mathbf{x} in the following fashion: each occurrence of a in \mathbf{x} is replaced by $\mathbf{p}^i\mathbf{q}$ and each occurrence of b by $\mathbf{p}^j\mathbf{q}$. The resulting string is denoted as $\sigma(\mathbf{x})$. We define $\sigma(\varepsilon) = \varepsilon$. The composition of two expansions σ_1 and σ_2 , $(\sigma_1 \circ \sigma_2)(\mathbf{x})$ is defined by $(\sigma_1 \circ \sigma_2)(\mathbf{x}) = \sigma_1(\sigma_2(\mathbf{x}))$. The role of the scope λ is to limit the size of \mathbf{p} and \mathbf{q} that can be used in the following definition.

Definition 1.4 A binary string \mathbf{x} is a **two-pattern string of scope λ** if there exists a sequence $\{\sigma_1, \sigma_2, \dots, \sigma_m\}$ of expansions of scope λ so that $\mathbf{x} = \sigma_1 \circ \dots \circ \sigma_m(a)$.

It was mentioned at the end of [FLS03] that if the definition of \mathbf{p} -regularity were made more restrictive, a larger class of complete two-pattern strings could be obtained. The more restrictive definition, sufficient to give two-pattern strings all their desired properties, contained a few typographical errors as it was given in [FLS03], and so we provide a corrected definition here:

Definition 1.5 A binary string \mathbf{q} is said to be **\mathbf{p} -regular** (\mathbf{p} a binary string) if and only if there exist (possibly empty) strings \mathbf{u}, \mathbf{v} together with nonnegative integers $n_1, n_2, \dots, n_k, k \geq 1, r \geq 0$, such that

- the integers n_i assume at most two distinct values — that is,

$$|\{n_i : i \in 1..k\}| \leq 2;$$

- $\mathbf{q} = (\mathbf{u}\mathbf{p}^r\mathbf{v}\mathbf{p}^{n_1})(\mathbf{u}\mathbf{p}^r\mathbf{v}\mathbf{p}^{n_2}) \dots (\mathbf{u}\mathbf{p}^r\mathbf{v}\mathbf{p}^{n_k})\mathbf{u}$ for some $\mathbf{u}, \mathbf{v}, r \geq 0$, where $\mathbf{v} = \varepsilon$ if $r = 0$.

Note: the definition 1.5 can be used to replace the definition 1.1. In fact, all the proofs accompanying this paper are compliant with the more restrictive definition 1.5.

Certain finite fragments of the well-known infinite Fibonacci string and the equally well-known infinite Sturmian strings are in fact complete two-pattern strings of scope $\lambda = 1$ (see [FLS03]).

Here are a few simple examples of two-pattern strings:

1. a , now apply $\sigma_2 = [ab, ba, 2, 3]$ to it, we get
2. $\sigma_2(a) = ababba$, now apply $\sigma_1 = [abb, aa, 1, 4]$ to it, we get
3. $\sigma_1(\sigma_2(ababba)) = abbaa(abb)^4aaabbaa(abb)^4aa(abb)^4aaabbaa$.

Strings 1, 2, and 3 are all two-pattern strings of scope 3 (string 2 is in fact of scope 2, and string 1 is in fact of scope 1).

It was shown in [FLS03] that complete two-pattern strings can be recognized in linear time: the recognition algorithm outputs an essentially unique sequence of expansions to construct the string from a . So in the following we can assume that not only do we have a complete two-pattern string, but also the sequence of expansions that iteratively generates the string.

In the next section we describe the principles underlying the algorithm for sorting suffixes of a two-pattern string. In Section 3 we provide an overview of the algorithm itself, while Section 5 we list some of the main lemmas on which the algorithm is based. We conclude with Section 6.

2 The Principles Underlying the Algorithm

For the sake of clarity and brevity, we introduce several symbols: we use the symbol $\mathbf{u} < \mathbf{v}$ for strings \mathbf{u}, \mathbf{v} to express that \mathbf{u} is lexicographically smaller than \mathbf{v} . We use the symbol \prec in $\mathbf{u} \prec \mathbf{v}$ (or \succ in $\mathbf{u} \succ \mathbf{v}$) to express the fact that $\mathbf{u} < \mathbf{v}$ yet \mathbf{u} is not

a prefix of \mathbf{v} (or $\mathbf{v} < \mathbf{u}$ yet \mathbf{v} is not a prefix of \mathbf{u}). Note that $\mathbf{u} < \mathbf{v}$ iff ($\mathbf{u} \prec \mathbf{v}$ or \mathbf{u} is a prefix of \mathbf{v}). We use the symbol $\mathbf{u} \asymp \mathbf{v}$ to indicate that either $\mathbf{u} \prec \mathbf{v}$ or $\mathbf{u} \succ \mathbf{v}$.

For a binary string \mathbf{u} , we will use $\bar{\mathbf{u}}$ to denote its ones-complement; that is, the string formed by interchanging a 's and b 's in \mathbf{u} .

In accordance with [FLS03], if \mathbf{x} , \mathbf{y} are complete two-pattern strings, σ an expansion, and $\mathbf{y} = \sigma(\mathbf{x})$, then the occurrences of copies of \mathbf{p} and copies of \mathbf{q} in the concatenation of blocks $\mathbf{p}^i\mathbf{q}$ and $\mathbf{p}^j\mathbf{q}$ as defined by $\sigma(\mathbf{x})$ are called **restrained** copies. Any other occurrence of \mathbf{p} or \mathbf{q} is referred to as **free**. A consecutive sequence of restrained copies of \mathbf{p} 's and/or \mathbf{q} 's will also be referred to as a **restrained configuration** or a **restrained substring** of \mathbf{y} .

Throughout the following discussion we assume that the scope λ is fixed and that $\mathbf{y} = \sigma(\mathbf{x})$, where \mathbf{x} is a complete two-pattern string of scope λ and $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ an expansion of scope λ . Moreover we assume that all suffixes of \mathbf{x} are lexicographically sorted: $\rho_1 < \dots < \rho_{|\mathbf{x}|}$. We then describe how to order the suffixes of \mathbf{y} . We may assume further that $\mathbf{q} < \mathbf{p}$. If it were not the case, according to Lemma 5.2 (see section 5), $\bar{\mathbf{q}} < \bar{\mathbf{p}}$, we sort all of the suffixes of $\bar{\mathbf{y}} = \bar{\sigma}(\mathbf{x})$, where $\bar{\sigma} = [\bar{\mathbf{p}}, \bar{\mathbf{q}}, i, j]_\lambda$, and reversing the order, we get all suffixes of \mathbf{y} ordered properly.

Since we are assuming $\mathbf{q} < \mathbf{p}$, according to Lemma 5.1 (see section 5), for any suffixes ρ_1, ρ_2 of \mathbf{x} , if $\rho_1 < \rho_2$, then $\sigma(\rho_1) < \sigma(\rho_2)$. In simple terms, the assumption $\mathbf{q} < \mathbf{p}$ makes all expansions to preserve the order of suffixes.

We put all the suffixes of \mathbf{y} into disjoint buckets of five types **A–E**. Their definitions follow (*note that the expansion $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ is fixed*):

- For every nontrivial suffix δ of \mathbf{p} and for every integer k , $0 < k < i$,
 $\mathbf{A}_{\delta,k} = \{\delta\mathbf{p}^k\mathbf{q}\sigma(\rho) : \rho \text{ is a proper suffix of } \mathbf{x} \text{ or } \rho = \varepsilon\}$;
- for every nontrivial suffix δ of \mathbf{p} that is also a suffix of \mathbf{q} ,
 $\mathbf{A}_{\delta,i} = \{\delta\mathbf{p}^i\mathbf{q}\sigma(\rho) : \rho \text{ is a proper suffix of } \mathbf{x} \text{ or } \rho = \varepsilon\}$;
- for every nontrivial suffix δ of \mathbf{p} that is not a suffix of \mathbf{q} ,
 $\mathbf{A}_{\delta,i} = \{\delta\mathbf{p}^i\mathbf{q}\sigma(\rho) : b\rho \text{ is a proper suffix of } \mathbf{x}, \rho \text{ can be empty}\}$;
- for every nontrivial suffix δ of \mathbf{p} and for every integer k , $i < k < j$,
 $\mathbf{A}_{\delta,k} = \{\delta\mathbf{p}^k\mathbf{q}\sigma(\rho) : b\rho \text{ is a proper suffix of } \mathbf{x}, \rho \text{ can be empty}\}$.
- for every nontrivial suffix δ of \mathbf{p} ,
 $\mathbf{B}_\delta = \{\delta\mathbf{q}\sigma(\rho) : \rho \text{ is a proper nontrivial suffix of } \mathbf{x}\}$;
- for every nontrivial suffix δ of \mathbf{q} that is not a suffix of \mathbf{p} ,
 $\mathbf{C}_\delta = \{\delta\mathbf{p}^i\mathbf{q}\sigma(\rho) : a\rho \text{ is a proper suffix of } \mathbf{x}, \rho \text{ can be empty}\}$;
- for every nontrivial suffix δ of \mathbf{q} ,
 $\mathbf{D}_\delta = \{\delta\mathbf{p}^j\mathbf{q}\sigma(\rho) : b\rho \text{ is a proper suffix of } \mathbf{x}, \rho \text{ can be empty}\}$;
- $\mathbf{E} = \{\delta\mathbf{q} : \delta \text{ is a nontrivial suffix of } \mathbf{p}\} \cup \{\delta : \delta \text{ is a nontrivial suffix of } \mathbf{q}\}$.

(where the term *proper suffix* refers to a suffix that is not equal to the whole string and the term *trivial suffix* refers to the empty suffix).

It is straightforward to check that any suffix of \mathbf{y} belongs to one of the buckets $\mathbf{A-E}$ (for proof see the supplement, see below). We are going to order the suffixes in buckets $\mathbf{A-D}$ based on the ordering of the suffixes for \mathbf{x} (Step 1), then merge in the suffixes from \mathbf{E} (Steps 2 & 3); since $|\mathbf{E}| \leq 2\lambda$, this will not destroy the linearity of the algorithm. Note that the order within each bucket is determined by the order of suffixes of \mathbf{x} :

- in the bucket $\mathbf{A}_{\delta,k}$: $\delta \mathbf{p}^k \mathbf{q} \sigma(\rho_1) < \delta \mathbf{p}^k \mathbf{q} \sigma(\rho_2)$ if $\rho_1 < \rho_2$;
- in the bucket \mathbf{B}_{δ} : $\delta \mathbf{q} \sigma(\rho_1) < \delta \mathbf{q} \sigma(\rho_2)$ if $\rho_1 < \rho_2$;
- in the bucket \mathbf{C}_{δ} : $\delta \mathbf{p}^i \mathbf{q} \sigma(\rho_1) < \delta \mathbf{p}^i \mathbf{q} \sigma(\rho_2)$ if $\rho_1 < \rho_2$;
- and in the bucket \mathbf{D}_{δ} : $\delta \mathbf{p}^j \mathbf{q} \sigma(\rho_1) < \delta \mathbf{p}^j \mathbf{q} \sigma(\rho_2)$ if $\rho_1 < \rho_2$.

Thus, it is straightforward to list the suffixes in each bucket in the correct order, given the order of the suffixes of \mathbf{x} .

We make use of the following notation: if X, Y are sets of suffixes of \mathbf{y} , we write $X \ll Y$ iff $(\forall x \in X)(\forall y \in Y)(x < y)$. The major observation our algorithm is based on is that the buckets are linearly ordered by \ll ; that is, pairwise orderings can be made between bucket pairs of types

$$\mathbf{AA}, \mathbf{AB}, \mathbf{AC}, \mathbf{AD}, \mathbf{BB}, \mathbf{BC}, \mathbf{BD}, \mathbf{CC}, \mathbf{CD}, \mathbf{DD}, \quad (1)$$

based on five mutually exclusive (and exhaustive) conditions on any pair δ_1, δ_2 of suffixes of \mathbf{p} and/or \mathbf{q} :

- (C1) $\delta_1 \prec \delta_2$;
- (C2) $\delta_1 \succ \delta_2$;
- (C3) δ_1 is a proper prefix of δ_2 ;
- (C4) δ_2 is a proper prefix of δ_1 ;
- (C5) $\delta_1 = \delta_2 = \delta$.

Observe that, given δ_1 and δ_2 , to determine which of these conditions holds requires at most λ letter comparisons (since $|\delta_1| \leq \lambda, |\delta_2| \leq \lambda$).

Thus, for example, two \mathbf{A} buckets can be compared as follows:

- (C1) $\mathbf{A}_{\delta_1, k_1} \ll \mathbf{A}_{\delta_2, k_2}$.
- (C2) $\mathbf{A}_{\delta_2, k_2} \ll \mathbf{A}_{\delta_1, k_1}$.
- (C3) Let $\delta_2 = \delta_1 \delta'_1$ for some nonempty δ'_1 :
 - (a) if $\delta'_1 \prec \mathbf{p}$, then $\mathbf{A}_{\delta_2, k_2} \ll \mathbf{A}_{\delta_1, k_1}$;
 - (b) otherwise, $\mathbf{A}_{\delta_1, k_1} \ll \mathbf{A}_{\delta_2, k_2}$.
- (C4) Let $\delta_1 = \delta_2 \delta'_2$ for some nonempty δ'_2 :
 - (a) If $\delta'_2 \prec \mathbf{p}$, then $\mathbf{A}_{\delta_1, k_1} \ll \mathbf{A}_{\delta_2, k_2}$;
 - (b) otherwise, $\mathbf{A}_{\delta_2, k_2} \ll \mathbf{A}_{\delta_1, k_1}$.

- (C5) (a) If $k_1 < k_2$, then $\mathbf{A}_{\delta,k_1} \ll \mathbf{A}_{\delta,k_2}$;
 (b) if $k_1 = k_2$, then $\mathbf{A}_{\delta,k_1} = \mathbf{A}_{\delta,k_2}$;
 (c) if $k_1 > k_2$, then $\mathbf{A}_{\delta,k_2} \ll \mathbf{A}_{\delta,k_1}$.

It is not very hard to prove that this ordering is correct. The demonstration for cases (C1), (C2) and (C5) is straightforward. For (C3), observe that we are comparing $\delta_1 \mathbf{p}^{k_1} \mathbf{q} \dots$ with $\delta_2 \mathbf{p}^{k_2} \mathbf{q} \dots$, hence $\mathbf{p}^{k_1} \mathbf{q} \dots$ with $\delta'_1 \mathbf{p}^{k_2} \mathbf{q} \dots$. Since δ'_1 is a suffix of δ_2 , it is also a suffix of \mathbf{p} and so cannot be a prefix of \mathbf{p} . It follows that either $\delta'_1 \prec \mathbf{p}$ or $\delta'_1 \succ \mathbf{p}$, and the result follows. The proof for (C4) is exactly analogous.

Furthermore the \mathbf{AA} ordering is efficient, since the cases (a) and (b) in (C3) and (C4) can be processed in at most λ constant-time steps in addition to the λ steps that may be required to identify which condition holds: thus a total of at most 2λ steps altogether.

The results for the other pairs listed in (1) are similar: the details vary slightly from one case to another. The main result is that any of the pairs can be processed in at most 3λ steps, a constant. To avoid distracting the reader with unnecessary and uninteresting detail, we do not include the other cases here. For those details, please access the web supplement of this paper at

<http://www.cas.mcmaster.ca/~franek/web-publications.html>

This supplement will be soon available as a technical report of Department of Computing and Software, McMaster University, Hamilton, Ontario, L8S 4K1 Canada.

3 The High-Level Logic of the Algorithm

We describe only the recursive step (Step 1) that takes us from \mathbf{x} and its sorted suffixes to the corresponding sorted suffixes of $\mathbf{y} = \sigma(\mathbf{x})$, where $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$. Recall that we assume $\mathbf{q} < \mathbf{p}$.

1. Create names (A, δ) for every suffix δ of \mathbf{p} . (*This requires at most λ steps. Each name will be eventually replaced by a sequence of buckets, see below.*)
2. Sort the names according to the order described in the previous section for mutual comparison of the four \mathbf{A} buckets (of course, according to (C1)-(C4) only). (*This requires at most $2\lambda^3$ steps as we are sorting λ names and each comparison requires $\leq 2\lambda$ steps.*)
3. Replace every name (A, δ) by a sequence of names (A, δ, k) , $1 \leq k < j$. Let us call the resulting sequence BUCKETS. (*Now we have the names of \mathbf{A} buckets in the proper order. This requires at most $|\mathbf{y}|$ steps as the size of BUCKETS is $\leq |\mathbf{y}|$. Each name (A, δ, k) will eventually be replaced by a corresponding bucket $\mathbf{A}_{\delta,k}$, see below.*)
4. Create names (B, δ) for every suffix δ of \mathbf{p} . (*This requires at most λ steps. Each name (B, δ) will eventually be replaced by a corresponding bucket \mathbf{B}_δ , see below.*)

5. Merge into BUCKETS all names (B, δ) according to comparisons as described in comparing \mathbf{A} buckets to \mathbf{B} buckets. (*This requires at most $|\text{BUCKETS}|3\lambda^2$ steps, as we are merging in λ names and each comparison requires $\leq 3\lambda$ steps, hence at most $|\mathbf{y}|3\lambda^2$ steps.*)
6. Create names (C, δ) for every suffix δ of \mathbf{q} that is not a suffix of \mathbf{p} . (*This requires at most λ^2 steps. Each name (C, δ) will eventually be replaced by a bucket \mathbf{C}_δ , see below.*)
7. Merge into BUCKETS all names (C, δ) according to comparisons as described in comparing \mathbf{A} buckets to \mathbf{C} buckets and \mathbf{B} buckets to \mathbf{C} buckets. (*This requires at most $|\text{BUCKETS}|3\lambda^2$ steps, hence at most $|\mathbf{y}|3\lambda^2$ steps.*)
8. Create names (D, δ) for every suffix δ of \mathbf{q} . (*This requires at most λ steps. Each name (D, δ) will eventually be replaced by a bucket \mathbf{D}_δ , see below.*)
9. Merge into BUCKETS all names (D, δ) according to comparisons as described in comparing \mathbf{A} buckets to \mathbf{D} buckets, \mathbf{B} buckets to \mathbf{D} buckets, \mathbf{C} buckets to \mathbf{D} buckets. (*Now we have all required bucket names, except \mathbf{E} , in proper order. This requires at most $|\text{BUCKETS}|3\lambda^2$ steps, hence at most $|\mathbf{y}|3\lambda^2$ steps.*)
10. Traverse BUCKETS and replace each name by a sequence of suffixes according to the sequence of suffixes of \mathbf{x} . Let us call this sequence SUFFIXES. (*We turned the names into proper buckets and merged them all together in a single list. Now we have all suffixes from buckets \mathbf{A} – \mathbf{D} in proper order. This requires at most $|\mathbf{y}|$ steps as the size of SUFFIXES is $\leq |\mathbf{y}|$.)*)
11. Merge into SUFFIXES the suffixes from the bucket \mathbf{E} . (*This requires at most $|\text{SUFFIXES}|4\lambda^2$ steps, as we are merging in 2λ suffixes, each of length $\leq 2\lambda$, hence at most $|\mathbf{y}|4\lambda^2$ steps.*)

SUFFIXES is now a sorted list of all suffixes of \mathbf{y} and it took less than $\alpha|\mathbf{y}|$ steps, where we set $\alpha = 2\lambda^3 + 14\lambda^2 + 3\lambda + 2$. Since every reduction of a complete two-pattern string at least halves its length, altogether the algorithm with all iterative steps included took less than $\alpha n + \alpha \frac{n}{2} + \alpha \frac{n}{4} + \dots < 2\alpha n$ steps, where n is the size of the input string.

4 An example

Let $\mathbf{x} = ababa$, and let $\sigma = [ba, ab, 1, 2]$. (Thus $\mathbf{q} = ab < \mathbf{p} = ba$.) Hence $\mathbf{y} = \sigma(\mathbf{x}) = baabbabaabbaabbabaabbaab$.

All nontrivial proper suffixes of \mathbf{x} are a , aba , b , and $baba$. All nontrivial suffixes of \mathbf{p} are ba and a , and all nontrivial suffixes of \mathbf{q} are ab and b . Let see the buckets:

$$\begin{aligned} \mathbf{A}_{ba,1} &= \{babaab\sigma(a), babaab\sigma(aba)\} = \\ &\quad \{babaabbaab, babaabbaabbabaab\} = \{\mathbf{y}[15..24], \mathbf{y}[5..24]\}. \\ \mathbf{A}_{a,1} &= \{abaab\sigma(a), abaab\sigma(aba)\} = \{abaabbaab, abaabbaabbabaab\} = \\ &\quad \{\mathbf{y}[16..24], \mathbf{y}[6..24]\}. \\ \mathbf{B}_{ba} &= \{baab\sigma(a), baab\sigma(aba), baab\sigma(ba), baab\sigma(baba)\} = \\ &\quad \{baabbaab, baabbaabbabaab, baabbabaab, \end{aligned}$$

$$\begin{aligned}
 & \{baabbabaabbaabbabaab\} = \{\mathbf{y}[17..24], \mathbf{y}[7..24], \mathbf{y}[11..24], \mathbf{y}[1..24]\} \\
 \mathbf{B}_a &= \{aab\sigma(a), aab\sigma(aba), aab\sigma(ba), aab\sigma(baba)\} = \\
 & \{aabbaab, abbaabbabaabbaab, aabbabaabbaab, \\
 & aabbabaabbaabbabaab\} = \{\mathbf{y}[18..24], \mathbf{y}[8..24], \mathbf{y}[12..24], \mathbf{y}[2..24]\} \\
 \mathbf{C}_{ab} &= \{abbaab\sigma(\varepsilon), abbaab\sigma(a)\} = \{abbaab, abbaabbabaabbaab\} = \\
 & \{\mathbf{y}[19..24], \mathbf{y}[9..24]\} \\
 \mathbf{C}_b &= \{bbaab\sigma(\varepsilon), bbaab\sigma(ba)\} = \{bbaab, bbaabbabaabbaab\} = \\
 & \{\mathbf{y}[20..24], \mathbf{y}[10..24]\} \\
 \mathbf{D}_{ab} &= \{abbabaab\sigma(a), abbabaab\sigma(aba)\} = \\
 & \{abbabaabbaab, abbabaabbaabbabaab\} = \{\mathbf{y}[13..24], \mathbf{y}[3..24]\} \\
 \mathbf{D}_b &= \{bbabaab\sigma(a), bbabaab\sigma(aba)\} = \\
 & \{bbabaabbaab, bbabaabbaabbabaab\} = \{\mathbf{y}[14..24], \mathbf{y}[4..24]\} \\
 \mathbf{E} &= \{baab, aab, ab, b\} = \{\mathbf{y}[21..24], \mathbf{y}[22..24], \mathbf{y}[23..24], \mathbf{y}[24..24]\}
 \end{aligned}$$

First note that we really got all nontrivial suffixes of \mathbf{y} . Also note that the suffixes in the buckets are in proper order. Let us see the mutual relationship of buckets:

$$\begin{aligned}
 & \mathbf{A}_{ba,1} \gg \mathbf{A}_{a,1} \text{ (by (C1))}, \mathbf{A}_{ba,1} \gg \mathbf{B}_{ba} \text{ (by (C5))}, \mathbf{A}_{ba,1} \gg \mathbf{B}_a \text{ (by (C2))}, \\
 & \mathbf{A}_{ba,1} \gg \mathbf{C}_{ab} \text{ (by (C2))}, \mathbf{A}_{ba,1} \ll \mathbf{C}_b \text{ (by (C4a))}, \mathbf{A}_{ba,1} \gg \mathbf{D}_{ab} \text{ (by (C2))}, \\
 & \mathbf{A}_{ba,1} \ll \mathbf{D}_b \text{ (by (C4a))}, \mathbf{A}_{a,1} \ll \mathbf{B}_{ba} \text{ (by (C1))}, \mathbf{A}_{a,1} \gg \mathbf{B}_a \text{ (by (C5))}, \\
 & \mathbf{A}_{a,1} \ll \mathbf{C}_{ab} \text{ (by (C3b))}, \mathbf{A}_{a,1} \ll \mathbf{C}_b \text{ (by (C1))}, \mathbf{A}_{a,1} \ll \mathbf{D}_{ab} \text{ (by (C3b))}, \\
 & \mathbf{A}_{a,1} \ll \mathbf{D}_b \text{ (by (C1))}, \mathbf{B}_{ba} \gg \mathbf{B}_a \text{ (by (C2))}, \mathbf{B}_{ba} \gg \mathbf{C}_{ab} \text{ (by (C2))}, \\
 & \mathbf{B}_{ba} \ll \mathbf{C}_b \text{ (by (C4b))}, \mathbf{B}_{ba} \gg \mathbf{D}_{ab} \text{ (by (C2))}, \mathbf{B}_{ba} \ll \mathbf{D}_b \text{ (by (C4a))}, \\
 & \mathbf{B}_a \ll \mathbf{C}_{ab} \text{ (by (C3b))}, \mathbf{B}_a \ll \mathbf{C}_b \text{ (by (C1))}, \mathbf{B}_a \ll \mathbf{D}_{ab} \text{ (by (C3b))}, \\
 & \mathbf{B}_a \ll \mathbf{D}_b \text{ (by (C1))}, \mathbf{C}_{ab} \ll \mathbf{C}_b \text{ (by (C1))}, \mathbf{C}_{ab} \ll \mathbf{D}_{ab} \text{ (by (C5))}, \\
 & \mathbf{C}_{ab} \ll \mathbf{D}_b \text{ (by (C1))}, \mathbf{C}_b \gg \mathbf{D}_{ab} \text{ (by (C2))}, \mathbf{C}_b \ll \mathbf{D}_b \text{ (by (C5))}, \\
 & \mathbf{D}_{ab} \ll \mathbf{D}_b \text{ (by (C1))}.
 \end{aligned}$$

Now follow the 11 steps.

1. create names $(A, ba), (A, a)$
2. sort them: $(A, a), (A, ba)$ (according to (C1))
3. “refine” the names to BUCKETS= $(A, a, 1), (A, ba, 1)$
4. create names to $(B, ba), (B, a)$
5. merge them into BUCKETS= $(B, a), (A, a, 1), (B, ba), (A, ba, 1)$
6. create names to $(C, ab), (C, b)$
7. merge them into BUCKETS= $(B, a), (A, a, 1), (C, ab), (B, ba), (A, ba, 1), (C, b)$
8. create names to $(D, ba), (D, a)$
9. merge them into BUCKETS= $(B, a), (A, a, 1), (C, ab), (D, ab), (B, ba), (A, ba, 1), (C, b), (D, b)$
10. replace the names by buckets: SUFFIXES=
 $(\mathbf{y}[18..24], \mathbf{y}[8..24], \mathbf{y}[12..24], \mathbf{y}[2..24]), (\mathbf{y}[16..24], \mathbf{y}[6..24]),$
 $(\mathbf{y}[19..24], \mathbf{y}[9..24]), (\mathbf{y}[13..24], \mathbf{y}[3..24]),$

$(\mathbf{y}[17..24], \mathbf{y}[7..24], \mathbf{y}[11..24], \mathbf{y}[1..24]), (\mathbf{y}[15..24], \mathbf{y}[5..24]),$
 $(\mathbf{y}[20..24], \mathbf{y}[10..24]), (\mathbf{y}[14..24], \mathbf{y}[4..24])$

11. merge in \mathbf{E} bucket: SUFFIXES= $\mathbf{y}[22..24], \mathbf{y}[18..24], \mathbf{y}[8..24],$
 $\mathbf{y}[12..24], \mathbf{y}[2..24], \mathbf{y}[23..24], \mathbf{y}[16..24], \mathbf{y}[6..24], \mathbf{y}[19..24], \mathbf{y}[9..24],$
 $\mathbf{y}[13..24], \mathbf{y}[3..24], \mathbf{y}[24..24], \mathbf{y}[21..24], \mathbf{y}[17..24], \mathbf{y}[7..24], \mathbf{y}[11..24],$
 $\mathbf{y}[1..24], \mathbf{y}[15..24], \mathbf{y}[5..24], \mathbf{y}[20..24], \mathbf{y}[10..24], \mathbf{y}[14..24], \mathbf{y}[4..24]$

5 The Supporting Lemmas

For the proofs, see the supplement as mention before.

The first lemma establishes that the ordering of suffixes is invariant under an expansion with $\mathbf{q} < \mathbf{p}$.

Lemma 5.1 *Let $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ be an expansion and $\mathbf{q} < \mathbf{p}$. Let \mathbf{x} and \mathbf{y} be two-pattern strings of scope λ and let $\mathbf{y} = \sigma(\mathbf{x})$. Let ρ_1, ρ_2 be suffixes of \mathbf{x} so that $\rho_1 < \rho_2$. Then $\sigma(\rho_1) < \sigma(\rho_2)$.*

The next lemma tells us that interchanging a and b in a binary string reverses the order of the suffixes.

Lemma 5.2 *Let $\rho_1 < \dots < \rho_n$ be the sequence of all suffixes of a binary string \mathbf{u} in an ascending lexicographic order. Then $\overline{\rho_1} > \dots > \overline{\rho_n}$ is the sequence of all suffixes of $\overline{\mathbf{u}}$ in a descending lexicographic order.*

The next three lemmas are technical lemmas required for some of the proofs (see website referenced above) that the pairs (1) can be processed correctly in $O(3\lambda)$ time. Essentially these lemmas tell us that the ordering of restrained suffixes of \mathbf{y} can be accomplished in at most 2λ constant-time algorithmic steps.

Lemma 5.3 *Let \mathbf{x}, \mathbf{y} be two-pattern strings of scope λ , $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ an expansion, and $\mathbf{y} = \sigma(\mathbf{x})$. Let \mathbf{u} be a non-empty binary string and let \mathbf{uqp} be a suffix of a restrained configuration \mathbf{pqp} of \mathbf{y} and let \mathbf{qp} be a restrained configuration of \mathbf{y} . Then $\mathbf{uqp} \succ \mathbf{qp}$ and whether $\mathbf{uqp} \prec \mathbf{qp}$ or $\mathbf{uqp} \succ \mathbf{qp}$ can be determined in $\leq 2\lambda$ steps.*

Lemma 5.4 *Let \mathbf{x}, \mathbf{y} be two-pattern strings of scope λ , $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ an expansion, and $\mathbf{y} = \sigma(\mathbf{x})$. Let \mathbf{u} be a non-empty binary string and let \mathbf{up} be a suffix of a restrained configuration \mathbf{qp} of \mathbf{y} . Let $1 \leq k$, and let $\mathbf{p}^k \mathbf{q}$ be a restrained configuration of \mathbf{y} . Then $\mathbf{up} \succ \mathbf{p}^k \mathbf{q}$ and whether $\mathbf{up} \prec \mathbf{p}^k \mathbf{q}$ or $\mathbf{up} \succ \mathbf{p}^k \mathbf{q}$ can be determined in $\leq 2\lambda$ steps.*

Lemma 5.5 *Let \mathbf{x}, \mathbf{y} be two-pattern strings of scope λ , $\sigma = [\mathbf{p}, \mathbf{q}, i, j]_\lambda$ an expansion, and $\mathbf{y} = \sigma(\mathbf{x})$. Let \mathbf{u} be a non-empty binary string and let $\mathbf{up}^k \mathbf{q}$, $1 \leq k$, be a suffix of a restrained configuration $\mathbf{p}^{k+1} \mathbf{q}$ or $\mathbf{qp}^k \mathbf{q}$ of \mathbf{y} . Let \mathbf{qp} be a restrained configuration of \mathbf{y} . Then $\mathbf{up}^k \mathbf{q} \succ \mathbf{qp}$ and whether $\mathbf{up}^k \mathbf{q} \prec \mathbf{qp}$ or $\mathbf{up}^k \mathbf{q} \succ \mathbf{qp}$ can be determined in $\leq 2\lambda$ steps.*

6 Conclusion

Even though it is known that suffixes for all strings can be sorted in linear time using recursive algorithms, our research verified that for the class of complete two-pattern strings the sorting can be done iteratively, also in linear time. The analysis shows that the approach presented here is rather straightforward, thus providing additional evidence of how two-pattern strings are well-suited for computational processing, the main goal of this effort.

References

- [F97] M. Farach, **Optimal suffix tree construction with large alphabets**, in *Proc. 38th Annual Symposium on Foundations of Computer Science*, IEEE (1997) pp. 137–143.
- [FLS03] F. Franek, W. Lu, and W. F. Smyth, **Two-pattern strings I — a recognition algorithm**, *J. Discrete Algorithms 1–5/6* (2003) pp. 445–460.
- [FLS04] F. Franek, W. Lu, and W. F. Smyth, **Two-pattern strings II — computing all repetitions and near-repetitions**, submitted to *J. Discrete Algorithms*.
- [KA03] P. Ko and S. Aluru, **Space efficient linear time construction of suffix arrays**, *Proceedings of the 14th Annual Symposium CPM*, LNCS 2676, Springer (2003) pp. 200–210.
- [KSPP03] D. K. Kim, J. S. Sim, H. Park, and K. Park, **Linear-time construction of suffix arrays**, *Proceedings of the 14th Annual Symposium CPM*, LNCS 2676, Springer (2003) pp. 186–199.
- [KS03] J. Kärkkäinen and P. Sanders, **Simple linear work suffix array construction**, *Proceedings of the 30th International Colloquium on Automata, Languages and Programming*, LNCS 2719, Springer (2003) pp. 943–955.
- [MM93] U. Manber and G. Myers, **Suffix arrays: a new method for on-line string searches**, *SIAM Journal on Computing 22–5* (1993) pp. 935–948.

Acknowledgements

The first author would like to acknowledge the support and hospitality of the School of Computing, Curtin University, Perth, Australia during the research for this paper. The research of both authors was supported in part by their respective research grants from the Natural Sciences and Engineering Research Council of Canada.