

# Extended Parameterized Burrows–Wheeler Transform

Eric M. Osterkamp<sup>a</sup>, Dominik Köppl<sup>b</sup>

<sup>a</sup>*University of Münster, Department of Computer Science, Münster, Germany*

<sup>b</sup>*University of Yamanashi, Faculty of Engineering, Kofu, Japan*

---

## Abstract

The Burrows–Wheeler transform (BWT) lies at the heart of succinct and compressed full-text indexes for pattern matching queries. Notable variants are (a) the extended BWT (eBWT) capable to index multiple circular texts for pattern matching, or (b) the parameterized BWT (pBWT) for parameterized pattern matching. A natural extension is the combination of the virtues of both variants into a new data structure, whose name we coin with *extended parameterized BWT* (epBWT). We show that the epBWT supports pattern matching in context of parameterized pattern matching on multiple circular texts, within the same complexities as known solutions presented for the pBWT [Kim and Cho, IPL’21] for patterns not longer than the shortest indexed text. Additionally, we show how to compute the epBWT within the same complexities as [Iseri et al., ICALP’24], i.e., in compact space and quasi-linear time. As an application, we extend the matching statistics problem to the parameterized pattern matching setting on circular texts.

*Keywords:* extended Burrows–Wheeler transform, parameterized pattern matching, circular texts, matching statistics

---

## 1. Introduction

The seek for compact indexing data structures stems from the need to manage and analyze the ever-increasing amount of available data [2]. A mainstream line of research is devoted to fast yet space-economical variations of the FM-index [3], a full-text index built upon the Burrows–Wheeler transform (BWT) [4]. Given a pattern, the FM-index allows us to query for

---

\*This article is an extension of a contribution [1] to the Data Compression Conference 2024.

the number of its occurrences in the indexed text, called *count* query, or the starting positions of these occurrences, called *locate* query. For the former, it suffices to encode the BWT with a data structure supporting rank and select queries such as the wavelet tree [5]. For the latter, we add suffix array samples to map found positions in the BWT into text-range. The few number of tools make the FM-index a predestined data structure for compact text indexing. However, the standard variant only supports exact pattern matching in the classic sense. On the one hand, a common variation in bioinformatics and computational geometry is to search in circular texts [6]. For instance, the genetic data like of the herpes simplex virus (HSV-1) [7] or circular RNA [8, 9] can be modeled as circular strings. With the aim for an index built on circular texts for pattern matching, Mantaci et al. [10] proposed the extended BWT (eBWT), which can additionally index multiple circular strings. On the other hand, structural analysis of RNA data involves a different kind of pattern matching, coined as *parameterized pattern matching* by Baker [11], and introduced to the context of RNA data by Shibuya [12]. BWT-based indexes such as the *parameterized BWTs* (pBWTs) by Ganguly et al. [13] and Kim and Cho [14] emerge as compact indexes for parameterized pattern matching. Recently, Iseri et al. [15] showed how to build the latter pBWT variant.

In the light that bioinformatic problems ask for RNA indexes that can cope with circular as well as parameterized pattern matching, a natural question to ask is whether an extension of the FM-index exists that allows a combination of both types of queries. A straight-forward answer to this question is to use the pBWT to index the concatenation of the text with itself and apply some post-processing after a query. However, this has certain drawbacks: we need twice the amount of space for the second duplicate of the text to index and cannot find conjugates of texts whose iterated concatenation contain the pattern as a prefix. For instance, we consider the pattern  $P = \text{abcabcab}$  to cyclic match at the beginning of the text  $T = \text{abc}$  by regarding the text as its infinite concatenation  $\text{abcabcabc}\dots$ . Such an extension of pattern matching is of interest when querying for repetitive patterns, so that we can support texts that store only short snippets of the repetitive pattern we want to find.

Without the need for text duplication and post-processing, we present the *extended parameterized BWT* (epBWT), which is based on the ideas of Mantaci et al. [10] for the eBWT and the tricks used by Iseri et al. [15] for the pBWT. The epBWT achieves the same complexities as one of the

pBWT versions of Ganguly et al. [13] or of Kim and Cho [14]. As an application, we generalize the *matching statistics* problem proposed by Chang and Lawler [16] to a parameterized version on multiple circular texts, and show that the epBWT index can solve this problem.

### 1.1. Related Work

We briefly review work on circular texts and parameterized string matching. To avoid confusion, we remind the reader aware of the circular pattern matching problem that this problem asks for reporting all conjugates of the pattern, while we here study the problem to find a pattern in all conjugates of the circular texts subject to index.

*Circular Texts.* A folklore solution for a pattern matching index on a circular text is to index twice the linearized input by a standard indexing technique such as the suffix tree [17], and use a filter after a query to avoid reporting pattern occurrences completely within the second copy of the input. Let  $n$  be the total length of the  $d$  circular texts we want to index. A solution that does not need this duplication trick is the eBWT [10]. The eBWT takes as an input a multiset of primitive, cyclic strings and is defined as the last characters of all the rotations of all the cyclic strings in the multiset sorted in a lexicographic-like total order  $\leq_{\omega}$  [10, Definition 4], which does not coincide with the lexicographic order in general. For example, the eBWT of the multiset  $\{\text{ab}, \text{ab}, \text{aba}\}$  is  $\text{babbaaa}$  since  $\text{aab} \leq_{\omega} \text{aba} \leq_{\omega} \text{ab} \leq_{\omega} \text{ab} \leq_{\omega} \text{baa} \leq_{\omega} \text{ba} \leq_{\omega} \text{ba}$ .<sup>1</sup> The inverse of the eBWT, known as the *Gessel–Reutenauer transform* [18], has already been conceived more than ten years prior to the discovery of the eBWT. Given a string  $T$ , the Gessel–Reutenauer transform maps  $T$  to the multiset of strings whose eBWT is  $T$ . However, Mantaci et al. [10] did not cover algorithmic details like the construction — yet it is clear that the construction time is upper bounded by the total length of all conjugates [10, after Example 9], which is  $O(n^2)$ . Later, Hon et al. [19] proposed a circular pattern matching algorithm using eBWT. They use  $n \log \sigma(1 + o(1)) + O(n) + O(d \log n)$  bits and can answer a count query in  $O(m \lg \sigma)$  for a pattern of length  $m$ . Subsequently, with a larger number of authors [20], they gave a construction algorithm for the

---

<sup>1</sup>The relation  $\leq_{\omega}$  is not to be confused with the here later defined preorder  $\preceq_{\omega}$  since the former assigns a linear order even on periodic strings by measuring the lengths of their roots for breaking ties.

eBWT in  $O(n \lg n)$  time. Their idea is to construct a so-called *circular suffix array* CSA such that the  $i$ -th position of the eBWT is given by  $T[\text{CSA}[i] - 1]$ , where  $T$  is the concatenation of all input texts. For their algorithm to work, they make the assumption that all input strings are non-periodic, belong to different conjugacy classes, and have asymptotically the same lengths. Very recently, Cotumaccio [21] showed how to get rid of these assumptions while retaining the same time and space complexities. Bonomo et al. [22] presented an online algorithm building the eBWT in  $O(n \lg n / \lg \lg n)$  time; here, online means that a construction algorithm updates the eBWT on reading a new input string. An interesting aspect regarding the online construction is that the up so far only known online technique [23, 24] for computing the traditional BWT needs the text to be given in reversed order (starting with the last character). Recently, Bannai et al. [25] showed that the eBWT can be constructed in linear time. Subsequently, Boucher et al. [26] practically simplified this construction. One of the currently fastest eBWT construction algorithms in practice is due to Olbrich et al. [27].

*Parameterized Matching.* The parameterized suffix tree as an extension of suffix trees to the parameterized indexing problem by Baker [11, 28] pioneered a new research area that adapts classic indexes to parameterized pattern matching. We briefly mention some achievements in this field and refer to Mendivilso et al. [29] for a more comprehensive overview of this research area. In timeline, we are aware of parameterized versions of the suffix array [30, 31], the longest previous factor array [32], the position heap [33], the LCP array [34], the linear-sized suffix trie [35], the suffix tray [36], compact suffix trees [37], and DAWGs [38]. Finally, Ganguly et al. [13] proposed the pBWT as the first parameterized index that is compact, i.e., the number of required bits is linear in the length of the input text. Kim and Cho [14] gave a simplification of this technique, requiring however twice the amount of space. The construction of Kim and Cho's variant has been addressed by Hashimoto et al. [39] whose time, like other construction algorithm of preceding data structures, is linearly dependent on both the text length and the alphabet size of the parameterized symbols. Finally, a breakthrough in the time complexity was achieved by Iseri et al. [15], who improved the latter linear dependency to a logarithmic one, and thus gave for the very first time a construction algorithm for a parameterized index with quasilinear time for arbitrarily large alphabets.

### 1.2. Structure of this Article

We start with Section 2 to cover basic notations and tools. We introduce the pBWT and formalize the type of index that is subject to this article. Subsequently, we present the epBWT index in Section 3 and give a construction algorithm in Section 4. As an application, we solve a generalization of matching statistics to parameterized matching on circular texts in Section 5 and finally give a conclusion in Section 6 with open problems and future work. Our main result can be stated as follows.

**Theorem 1.1.** *Let  $T_1, \dots, T_d$  non-empty, parameterized strings over an alphabet  $\Sigma$  of size  $\sigma$  and  $n = |T_1 \cdots T_d|$ . The epBWT index of  $\{T_1, \dots, T_d\}$  is a data structure that takes  $O(n \lg \sigma)$  bits of space and can be computed in  $O(n \lg \sigma)$  bits of space and  $O(nt_{\text{query}})$  time, where  $t_{\text{query}} = \frac{\lg \sigma \lg n}{\lg \lg n}$ . Given a pattern  $P$  of length  $m$ , it counts all conjugates of  $T_1, \dots, T_d$  whose infinite concatenation has a prefix that  $p$ -matches with  $P$ , in  $O(mt_{\text{query}})$  time. It can restore the conjugacy classes of  $T_1, \dots, T_d$  up to  $p$ -match equivalence, in  $O(nt_{\text{query}})$  time. It can add and remove a text in  $O(n't_{\text{query}})$  time, where  $n'$  is the maximum length of all texts (including the one to insert). By switching to static data structures, we can improve  $t_{\text{query}}$  to  $\lg \sigma$  and the space to  $2n \lg \sigma + 2n + O(n \frac{\lg \lg n}{\lg n})$  bits of space.*

Compared to the conference version [1], we give details on the construction and matching statistics as an application. The presentation of the epBWT index has been simplified and didactically improved with examples and explanations. In particular, we could simplify the description of  $\omega$ -preorder of the conjugates in the epBWT by introducing a new encoding  $\langle .. \rangle_r$  (cf. Lemma 3.1) such that we no longer need to rely on the  $\llbracket .. \rrbracket$ -encoding like in the conference paper.

## 2. Preliminaries

Let  $\lg = \log_2$  and  $\mathbb{N}_{\geq 1}$  denote, respectively, the logarithm to base two and the set of natural numbers starting at 1. An interval  $\{i, i+1, \dots, j-1, j\}$  of integers is denoted by  $[i..j]$ , where  $[i..j] = \emptyset$  if  $j < i$ . Appendix A gives an overview of the variables used in this article. We start with basic notations for strings, queries we want to answer on strings, data structures that can answer those queries, and then move to the parameterized matching. In the end of this section, we formalize the problem we tackle in this article.

## 2.1. Strings

Let  $\Sigma$  denote an *alphabet*, i.e., a non-empty, finite and totally ordered set. An element of  $\Sigma$  is called a *symbol*, a sequence of symbols from  $\Sigma$  a *string over*  $\Sigma$  and a subsequence  $U$  of a consecutive range of elements from a string  $V$  over  $\Sigma$  a *substring* of  $V$ . A substring  $U$  of  $V$  is *proper* if  $U \neq V$ . Let  $\Sigma^*$  denote the set of all finite sequences of symbols from  $\Sigma$ ,  $\Sigma^\omega$  the set of all countably infinite sequences of symbols from  $\Sigma$ , and  $\Sigma^\infty = \Sigma^* \cup \Sigma^\omega$ . Then an element of  $\Sigma^*$  is called a *finite* string over  $\Sigma$  and an element of  $\Sigma^\omega$  an *infinite* string over  $\Sigma$ . Whether we are talking about a finite or an infinite string should be clear from the context at all times. The concatenation of any  $U \in \Sigma^*$  and  $V \in \Sigma^\infty$  is denoted by  $U \cdot V$  or  $UV$ . We write  $U^k$  if we concatenate  $k \in \mathbb{N}$  instances of  $U \in \Sigma^*$ , and we denote by  $U^\omega$  the infinite string obtained by infinitely iterating  $U$ , i.e.,  $U^\omega = UUU\cdots \in \Sigma^\omega$ . We call a string  $V \in \Sigma^*$  *primitive* if  $V = U^k$  for  $U \in \Sigma^*$  and  $k \in \mathbb{N}$  already implies  $V = U$  and  $k = 1$ . It is known that for every  $V \in \Sigma^*$  there exists a unique primitive  $U \in \Sigma^*$  and a unique  $k \in \mathbb{N}$  such that  $V = U^k$ , denoted by  $\text{root}(V)$  and  $\exp(V)$ , respectively.

If  $U = VW$  for  $V \in \Sigma^*$  and  $W \in \Sigma^\infty$ , then the substrings  $V$  and  $W$  are called *prefix* and *suffix* of  $U$ , respectively. Assume  $W \in \Sigma^*$ . Then  $WV$  is called a *conjugate* of  $U = VW$ . For  $U \in \Sigma^\infty$ , let  $|U|$  denote the *length* of  $U$ , i.e.,  $|U| = \infty$  if  $U \in \Sigma^\omega$ , or the number of symbols in  $U$  otherwise. Let  $\varepsilon$  denote the unique string of length 0 and denote the length of the longest common prefix of  $U, V \in \Sigma^\infty$  by  $\text{lcp}(U, V)$ . Also, let  $\Sigma^+ = \Sigma^* - \{\varepsilon\}$ .

Fix some  $V \in \Sigma^\infty$  and  $i, j \in [1..|V|]$ . Then  $V[i]$  denotes the  $i$ -th symbol of  $V$ . A substring  $V[i] \cdots V[j]$  starting at position  $i$  and ending at position  $j$  of  $V$  is denoted by  $V[i..j]$ , where we set  $V[i..j] = \varepsilon$  if  $j < i$ . Let  $V[..i]$  denote the prefix  $V[1..i]$  of  $V$  and  $V[..]$  the suffix  $W$  of  $V$  such that  $V = V[..i-1]W$ . For notational convenience, if  $k \leq 0$ , then  $V[..k] = \varepsilon$ , and if  $k \geq |V| + 1$ , then  $V[k..] = \varepsilon$ .

We assume  $V \in \Sigma^+$  for the remainder of the section. Let  $\text{Rot}(V, 0) = V$  and  $\text{Rot}(V, k+1) = \text{Rot}(V, k)[2..]\text{Rot}(V, k)[1]$  for  $k \in \mathbb{N}$ , i.e.,  $\text{Rot}(V, k)$  denotes the  $k$ -th *left rotation* of  $V$ . Note that  $\text{Rot}(V, k) = \text{Rot}(V, k \bmod |V|)$  for every  $k \in \mathbb{N}$ . We call  $k \in [1..|V|]$  satisfying  $V[i] = V[i+k]$  for every  $i \in [1..|V|-k]$  a *period* of  $V$ .

Let  $<$  denote the lexicographical order on  $\Sigma$ . For  $U, W \in \Sigma^*$ , we write  $U < W$  if and only if  $U$  is a proper prefix of  $W$  or  $U[\text{lcp}(U, W)+1] < W[\text{lcp}(U, W)+1]$ . For  $U, W \in \Sigma^\omega$ , we write  $U < W$  if and only if there exists  $k \in \mathbb{N}$  such that  $U[..k-1] = W[..k-1]$  and  $U[k] < W[k]$ .

## 2.2. Queries on Strings

Let  $f \leq g \in \Sigma$  and  $i, j \in [1..|V|]$ . The query  $\text{rank}_g(V, i)$  returns the number of occurrences of  $g$  in  $V[..i]$ , the query  $\text{select}_g(V, i)$  returns the index of the  $i$ -th occurrence of  $g$  in  $V$  if it exists, and 0 otherwise and the query  $\text{rnkent}_V(i, j, f, g)$  returns the number of entries in  $V[i..j]$  having a value in  $[f..g]$ . Moreover, the operation  $\text{insert}_V(i, g)$  inserts the symbol  $g$  at position  $i$  of  $V$  and the operation  $\text{delete}_V(i)$  deletes the  $i$ -th entry of  $V$ , e.g.,  $\text{insert}_{312}(2, 4)$  transforms 312 into 3412 and  $\text{delete}_{3412}(3)$  transforms 3412 into 342. For  $i \leq j$ , the query  $\text{RmQ}_V(i, j)$  returns the index of a minimal value in  $V[i..j]$ , the query  $\text{RMQ}_V(i, j)$  returns the index of a maximal value in  $V[i..j]$ , and the query  $\text{RNQ}_V(i, j, g)$  returns the smallest value in  $V[i..j]$  larger than  $g$  if it exists, and  $g$  otherwise.

We borrow notation from Iseri et al. [15, Section 2.1], who introduced *find previous queries* (FPQ) and *find next queries* (FNQ) as syntactic elements for range searches in a string. Both queries are generalizations of predecessor and successor queries by additionally requiring the output to satisfy a predicate.<sup>2</sup> In detail, they are defined as follows. Let  $i \in [0..|V|]$ . The query  $\text{FNQ}_g(V, i)$  returns the smallest index  $k \in [\max\{i, 1\}..|V|]$  such that  $V[k] = g$  if it exists and  $|V|+1$  otherwise. Similarly, the query  $\text{FNQ}_{\geq g}(V, i)$  returns the smallest  $k \in [\max\{i, 1\}..|V|]$  such that  $V[k] \geq g$  if it exists, and  $|V|+1$  otherwise. The query  $\text{FPQ}_g(V, i)$  returns the largest  $k \in [1..i]$  such that  $V[k] = g$  if it exists and 0 otherwise, and the query  $\text{FPQ}_{\geq g}(V, i)$  returns the largest  $k \in [1..i]$  such that  $V[k] \geq g$  if it exists, and 0 otherwise. Last, the query  $\text{MI}_V(i, g)$  returns the maximal interval  $[\ell..r]$  such that  $0 \leq \ell \leq i \leq r \leq |V|$ ,  $V[k] \geq g$  for every  $k \in [\ell + 1..r]$ . Note that FPQ <sub>$g$</sub>  and FNQ <sub>$g$</sub>  queries can be answered by a constant number of rank and select queries, and an MI query can be answered by a constant number of FPQ <sub>$\geq g$</sub>  and FNQ <sub>$\geq g$</sub>  queries. See Figure 1 for examples.

Depending on the queries and operations required, we may represent a string by one of the following data structures. We assume a random access model with word size  $\Omega(\lg n)$ , where  $n$  denotes the length of an input string.

**Lemma 2.1** ([41, 42]). *A static string of length  $n$  over  $[0..\sigma]$  with  $\sigma \leq n^{O(1)}$  can be stored in a data structure occupying  $n \lg \sigma + O(n \frac{\lg \lg n}{\lg n})$  bits of space,*

---

<sup>2</sup>A cousin of FPQs and FNQs are *next smaller value queries* [40], which can express, together with successor queries, FPQ and FNQs.

$i$	1	2	3	4	5	6	7	8
$V[i]$	7	14	5	1	11	27	11	7

Figure 1: Example integer string to illustrate queries. Here,  $\text{rank}_{11}(V, 5) = 1$ ,  $\text{select}_{11}(V, 2) = 7$ ,  $\text{select}_{23}(V, 1) = 0$ ,  $\text{rnkcnt}_V(2, 7, 7, 14) = 4$ ,  $\text{RmQ}_V(2, 5) = 4$ ,  $\text{RMQ}_V(2, 5) = 2$ ,  $\text{RNQ}_V(1, 5, 8) = 11$ ,  $\text{RNQ}_V(3, 5, 13) = 13$ ,  $\text{FNQ}_1(V, 2) = 4$ ,  $\text{FNQ}_{\geq 23}(V, 7) = 9$ ,  $\text{FPQ}_1(V, 2) = 0$ ,  $\text{FPQ}_{\geq 1}(V, 2) = 2$ ,  $\text{MI}_V(5, 6) = [4..8]$ , and  $\text{MI}_V(2, 15) = [2..2]$ .

supporting the queries access, rank, select and rnkcnt in  $O(\lg \sigma)$  time. It can be constructed in  $O(n \lg \sigma)$  time with additional  $n \lg \sigma + O(1)$  bits of space.

**Lemma 2.2** ([43]). *A static string of length  $n$  over an alphabet of size  $\sigma \in O(n)$  can be stored in a data structure occupying  $2n + O(n \frac{\lg \lg n}{\lg n})$  bits of space, supporting either the query RmQ or the query RMQ in  $O(1)$  time. It can be constructed in  $O(n)$  time with additional  $O(n)$  bits of space.*

**Lemma 2.3** ([44]). *A dynamic string of length  $n$  over  $[0..\sigma]$  with  $\sigma \leq n^{O(1)}$  can be stored in a data structure occupying  $O(n \lg \sigma)$  bits of space, supporting insertion, deletion and the queries access, rank and select in  $O(\frac{\lg n}{\lg \lg n})$  time.*

**Lemma 2.4** ([15, Lemma 4]). *A dynamic string of length  $n$  over  $[0..\sigma]$  with  $\sigma \leq n^{O(1)}$  can be stored in a data structure occupying  $O(n \lg \sigma)$  bits of space, supporting insertion, deletion and the queries access, rank, select, RNV, FPQ $_{\geq}$ , FNQ $_{\geq}$  and rnkcnt in  $O(\frac{\lg \sigma \lg n}{\lg \lg n})$  time.*

### 2.3. Parameterized Matching

Let  $\Sigma_s$  and  $\Sigma_p$  denote two disjoint alphabets,  $\sigma_s = |\Sigma_s|$  and  $\sigma_p = |\Sigma_p|$ . A symbol from  $\Sigma_s$  is called a *static symbol (s-symbol)* and a symbol from  $\Sigma_p$  is called a *parameterized symbol (p-symbol)*. We set  $\Sigma = \Sigma_s \cup \Sigma_p$ ,  $\sigma = |\Sigma|$ , and call a string over  $\Sigma$  a *parameterized string (p-string)*. Finally, we introduce two special symbols that do not convey any information of the input we want to index, but are needed for technical reasons. The first is  $\$$  that denotes the smallest s-symbol. The second is the symbol  $\infty \notin \Sigma_s \cup \Sigma_p$ , which we stipulate to be larger than any integer.

We assume  $\sigma \leq n^{O(1)}$ , i.e., each symbol from  $\Sigma$  can be represented by  $O(\lg n)$  bits of space. Since we can define order preserving mappings from  $\Sigma_s$  and  $\Sigma_p$  to  $[-\sigma_s .. -1]$  and  $[1..\sigma_p]$ , respectively, we can assume that  $\mathbb{N} \cap \Sigma_s = \emptyset$ , that any s-symbol is distinguishable from both p-symbols and positive

integers in  $O(1)$  time and that the order of two s-symbols can be determined in  $O(1)$  time. For any examples throughout, we will use  $\Sigma_s = \{\$, a, b\}$  and  $\Sigma_p = \{A, B, C\}$ , where  $\$ < a < b$ .

We say that two p-strings  $U, V \in \Sigma^*$  of equal length *parameterized match* (*p-match*) if and only if there exists a bijection  $\psi : \Sigma_p \rightarrow \Sigma_p$  such that  $U[i] = V[i]$  if  $V[i] \in \Sigma_s$ , or  $U[i] = \psi(V[i])$  otherwise, for each  $i \in [1..|U|]$ , and write  $U \approx V$ . For example,  $U = BCAaCB \approx CABaAC = V$ , where  $\psi(A) = C$ ,  $\psi(B) = A$  and  $\psi(C) = B$ .

The goal of this subsection is to study *parameterized matching*, which we define as follows.

**Problem 1** (P-COUNT). *Given a text p-string  $T$  and a pattern p-string  $P$ , count the number of all occurrences of  $P$  in  $T$ .*

By counting occurrences we mean to count all substrings of  $T$  that p-match with  $P$ , or more formally, to compute the cardinality of the subset  $\mathcal{I} \subset [1..n]$  with  $i \in \mathcal{I}$  if and only if  $P \approx T[i..i + |P| - 1]$ . As an example, we consider the text p-string  $T = CABaAC$  and the pattern p-strings  $P_1 = bA$  and  $P_2 = CB$ . It is easy to see that  $P_1$  has no occurrences in  $T$  since the s-symbol  $b$  does not occur in  $T$ . The pattern p-string  $P_2$ , on the other hand, has exactly 3 occurrences since  $P_2 \approx T[1..2] \approx T[2..3] \approx T[5..6]$ ,  $P_2 \not\approx T[3..4]$  and  $P_2 \not\approx T[4..5]$ .

While Baker [11] gives a linear time algorithm to check if two p-strings p-match, they note that it does not generalize conveniently to solving P-COUNT. Instead, they present an encoding that reduces checking whether two strings p-match to the exact matching of their encodings. In what follows, we introduce two such convenient encodings, the first of which is a variant of Baker's.

Kim and Cho's [14] *prev-encoding*  $\langle V \rangle$  of  $V \in \Sigma^*$  is a string of length  $|V|$  over  $\Sigma_s \cup [1..|V| - 1] \cup \{\infty\}$  such that

$$\langle V \rangle[i] = \begin{cases} V[i] & \text{if } V[i] \in \Sigma_s, \\ \infty & \text{if } V[i] \in \Sigma_p \wedge \forall (1 \leq j < i) : V[i] \neq V[j], \text{ and} \\ \text{FPQ}_{V[i]}(V, i - 1) & \text{otherwise,} \end{cases}$$

for every  $i \in [1..|V|]$ , i.e., the leftmost occurrence of a p-symbol in  $V$  is replaced by  $\infty$  and each successive occurrence with the distance to its previous occurrence. Note that  $\langle V \rangle[.i] = \langle V[.i] \rangle$  for each  $V \in \Sigma^*$  and  $i \in [1..|V|]$ .

$i$	$\langle V[i..] \rangle$	$\langle \text{Rot}(V, i) \rangle$	$\langle V[..i] \rangle$	$\llbracket V[i..] \rrbracket$	$\llbracket \text{Rot}(V, i) \rrbracket$	$\llbracket V[..i] \rrbracket$
1	$\infty\infty\infty a35$	$\infty\infty a3\infty 1$	$\infty$	323a21	23a213	1
2	$\infty\infty a3\infty$	$\infty a\infty\infty 13$	$\infty\infty$	23a23	3a2132	22
3	$\infty a\in\infty\infty$	$a\in\infty\infty 13\infty$	$\infty\in\infty\infty$	3a33	a21323	333
4	$a\in\infty\in\infty$	$\infty\in\infty 13\infty a$	$\infty\in\infty\in\infty a$	a22	21323a	333a
5	$\infty\in\infty$	$\infty 1\in\infty\infty a3$	$\infty\in\infty\in\infty a3$	22	1323a2	323a2
6	$\infty$	$\infty\in\infty\infty a35$	$\infty\in\infty\in\infty a35$	1	323a21	323a21

Figure 2: Encodings of substrings and conjugates of  $V = \text{CABaAC}$ .

For any  $U \in \Sigma^*$ , let  $|U|_p = \text{rank}_\infty(\langle U \rangle, |U|)$  denote the number of distinct p-symbols in  $U$ . The encoding  $\llbracket V \rrbracket$  of  $V \in \Sigma^*$  proposed by Hashimoto et al. [39] is given by a string of length  $|V|$  over  $\Sigma_s \cup [1..|\sigma_p|]$  such that  $\llbracket V \rrbracket[i] = V[i]$  if  $V[i] \in \Sigma_s$ , or  $\llbracket V \rrbracket[i] = |\text{Rot}(V, i)[.. \text{select}_{V[i]}(\text{Rot}(V, i), 1)]|_p$  otherwise, for each  $i \in [1..|V|]$ . In other words, we view  $V$  circularly and replace each occurrence of a p-symbol in  $V$  by the number of distinct p-symbols until its next occurrence. See Figure 2 for some examples.

**Lemma 2.5** ([11, 39, Propositions 1]). *Let  $U, V \in \Sigma^*$ . Then  $\langle U \rangle = \langle V \rangle \Leftrightarrow U \approx V \Leftrightarrow \llbracket U \rrbracket = \llbracket V \rrbracket$ .*

By Lemma 2.5, a naive solution for P-COUNT is to check either  $\langle P \rangle = \langle T[i..i + |P| - 1] \rangle$  or  $\llbracket P \rrbracket = \llbracket T[i..i + |P| - 1] \rrbracket$  for every  $i \in [1..n - |P| + 1]$ .

In what follows, we count the number of  $\infty$ 's in the longest common prefix of two p-strings  $U$  and  $V$ , which helps us to judge, given both are prefixed by symbols, the order between the prev-encoding of the prefixed versions, if the order of the prev-encoding of  $U$  and  $V$  is already known. We later use that fact for adapting the backward search, the centerpiece of the FM-index, to parameterized matching.

For  $V \in \Sigma^+$ , let  $\pi(V) = \llbracket V \rrbracket[1]$ . For  $U, V \in \Sigma^*$ , let  $\text{lcp}^\infty(U, V) = \text{rank}_\infty(\langle U \rangle, \text{lcp}(\langle U \rangle, \langle V \rangle))$ . Since  $\infty$  does not appear in  $\langle V^i \rangle[|V| + 1..]$  for any  $V \in \Sigma^*$  and integer  $i \geq 1$ ,  $\text{lcp}^\infty(U, V) = \text{lcp}^\infty(U^j, V^k)$  for each  $U, V \in \Sigma^*$  and integers  $j, k \geq 1$ . The following results are obtained by examination of the definitions.

**Lemma 2.6** ([15, Corollary 7]). *Let  $U, V \in \Sigma^+$ ,  $e = \text{lcp}^\infty(U[2..], V[2..])$  and  $\langle U[2..] \rangle < \langle V[2..] \rangle$ . Then  $\langle V \rangle < \langle U \rangle$  if and only if  $\pi(V) < \pi(U) \wedge (\pi(V) \in \Sigma_s \vee (\pi(U), \pi(V)) \notin \Sigma_s \wedge \pi(V) \leq e)$ .*

**Lemma 2.7** ([15, Table 1]). *Let  $U, V \in \Sigma^+$  and  $e = \text{lcp}^\infty(U[2..], V[2..])$ . Then*

$$\text{lcp}^\infty(U, V) = \begin{cases} 0 & \text{if } \pi(U) \neq \pi(V) \wedge (\pi(U) \in \Sigma_s \vee \pi(V) \in \Sigma_s), \\ e & \text{if } \pi(U) = \pi(V) \leq e, \\ \pi(U) & \text{if } \pi(U) \notin \Sigma_s \wedge \pi(U) \leq e \wedge \pi(U) < \pi(V), \\ \pi(V) & \text{if } \pi(V) \notin \Sigma_s \wedge \pi(V) \leq e \wedge \pi(V) < \pi(U), \\ e + 1 & \text{otherwise, i.e., if } e < \min\{\pi(U), \pi(V)\}. \end{cases}$$

#### 2.4. The Parameterized Burrows–Wheeler Transform

Let  $T \in \Sigma^+$  denote a string of length  $n$  such that  $\text{rank}_\$(T, n) = 1$  and  $T[n] = \$$ , i.e.,  $\$$  is the unique delimiter of the string  $T$ . In what follows, we present a variant of Ganguly et al.’s [13] original, space-efficient index solving P-COUNT for  $T$ .

The (parameterized) *suffix array* [30]  $\text{SA}_T$  is a string of length  $n$  over  $[1..n]$  such that  $\text{SA}_T[i] = j$  if and only if there are  $i - 1$  suffixes of  $T$  whose prev-encodings are smaller than the prev-encoding of the suffix  $T[j..]$ , i.e.,

$$\begin{aligned} i - 1 &= |\{k \in [1..n] \mid \langle \text{Rot}(T, k) \rangle < \langle \text{Rot}(T, j - 1) \rangle\}| \\ &= |\{k \in [1..n] \mid \langle T[k..] \rangle < \langle T[j..] \rangle\}|. \end{aligned}$$

Then the *parameterized Burrows–Wheeler Transform (pBWT)*  $\text{L}_T$  of  $T$  is a string of length  $n$  over  $\Sigma_s \cup [1..\sigma_p]$  such that  $\text{L}_T[i] = [\![\text{Rot}(T, \text{SA}_T[i] - 1)]\!]_n$  for each  $i \in [1..n]$ . The pBWT is further augmented to act as an index. We let  $\text{F}_T$  denote a string of length  $n$  over  $\Sigma_s \cup [1..\sigma_p]$  such that  $\text{F}_T[i] = [\![\text{Rot}(T, \text{SA}_T[i] - 1)]\!]_1$  for each  $i \in [1..n]$ , and we let  $\text{LCP}_T^\infty$  denote a string of length  $n$  over  $[0..\sigma_p]$  such that  $\text{LCP}_T^\infty[1] = 0$  and  $\text{LCP}_T^\infty[i] = \text{lcp}^\infty(\text{Rot}(T, \text{SA}_T[i] - 1), \text{Rot}(T, \text{SA}_T[i - 1] - 1))$  for each  $i \in [2..n]$ .

The statement regarding construction of the following result is originally for an index of a suffix-based variant of our pBWT, which differs in the definition of  $\text{F}_T$  and  $\text{L}_T$ . The transformation into  $\text{F}_T$  and  $\text{L}_T$  as defined here is done after construction by computing  $\llbracket T \rrbracket$  and utilization of the so-called LF-mapping.

**Lemma 2.8** ([14, 15]). *Let  $T \in \Sigma^+$ ,  $n = |T|$ ,  $T[n] = \$$  and  $\text{rank}_\$(T, n) = 1$ . Then  $\text{F}_T$ ,  $\text{L}_T$  and  $\text{LCP}_T^\infty$  form an index that occupies  $O(n \lg \sigma)$  bits of space and solves P-COUNT in  $O(m \lg \sigma)$  time, where  $m$  is the length of the pattern. The index can be computed in  $O(n \lg \sigma)$  bits of space and  $O(n \frac{\lg \sigma \lg n}{\lg \lg n})$  time, and the original text  $T$  restored up to  $p$ -match equivalence in  $O(n \lg \sigma)$  time.*

$i$	$\text{SA}_T[i]$	$F_T[i]$	$L_T[i]$	$V_T[i]$	$i$	$\text{SA}_T[i]$	$F_T[i]$	$L_T[i]$	$V_T[i]$
1	25	\$	1	0	14	8	1	2	1
2	2	a	1	0	15	10	1	2	2
3	5	a	2	1	16	18	1	b	2
4	3	b	a	0	17	12	1	2	3
5	22	b	2	1	18	20	3	3	1
6	17	b	3	1	19	15	3	2	3
7	6	b	a	1	20	7	2	b	2
8	24	1	1	0	21	9	2	1	2
9	1	1	\$	1	22	11	2	1	2
10	4	2	b	1	23	14	2	3	2
11	21	2	3	1	24	19	3	1	2
12	16	3	3	2	25	13	3	1	3
13	23	1	b	1					

Figure 3: The pBWT index of  $T = \text{BabBabABBAABBACAbBBCAbBB\$}$ . Here,  $V_T = \text{LCP}_T^\infty$ .

We call the index of Lemma 2.8 the *pBWT index of  $T$* . See Figure 3 for an example. The pBWT index of a text  $T \in \Sigma^+$  can compute the *suffix range*  $\text{SR}_T(P)$  of a pattern  $P \in \Sigma^*$ , which is the unique maximal interval  $[\ell..r] \subseteq [1..n]$  such that  $P \approx T[\text{SA}_T[i]..\text{SA}_T[i]+|P|-1]$  if and only if  $i \in [\ell..r]$ . The length of this interval is consequently the number of occurrences of  $P$  in  $T$ , i.e., the solution to P-COUNT.

The computation of the suffix range is also called the *backward search*, which works similarly to its namesake in the BWT. At its core is the so-called *LF-mapping*  $\text{LF}_T$ , which maps  $i \in [1..n]$  to  $\text{SA}_T^{-1}[\text{SA}_T[i]-1]$  if  $\text{SA}_T[i] > 1$ , and otherwise to  $\text{SA}_T^{-1}[n] = 1$ , where  $\text{SA}_T^{-1}$  denotes the inverse of the suffix array. It can be computed efficiently by finitely many queries on  $F_T$  and  $L_T$ , which is due to Lemma 2.6 and basing the definition of  $F_T$  and  $L_T$  on the  $\llbracket .. \rrbracket$ -encoding.

**Corollary 2.9** ([14]). *Let  $T \in \Sigma$  and  $i \in [1..|T|]$ . Then*

$$\text{LF}_T(i) = \text{select}_{L_T[i]}(F_T, \text{rank}_{L_T[i]}(L_T, i)).$$

Given a pattern p-string  $P \in \Sigma^*$ , a symbol  $g \in \Sigma$  and the suffix range  $\text{SR}_T(P)$ , we can use the LF-mapping to compute  $\text{SR}_T(gP)$ .

**Corollary 2.10** ([14]). *Let  $T, P \in \Sigma^*$  and  $g \in \Sigma$ . If  $g$  is an s-symbol or a p-symbol that already occurred in  $P$ , then  $\text{LF}_T(i) \in \text{SR}_T(gP)$  if and only*

if  $i \in \text{SR}_T(P)$  and  $\text{L}_T[i] = \pi(gP)$ , and otherwise  $\text{LF}_T(i) \in \text{SR}_T(gP)$  if and only if  $i \in \text{SR}_T(P)$  and  $\text{L}_T[i] \geq \pi(gP)$ .

We give an exemplary, naive backward search for  $P = \text{bCC}$  with the pBWT index of  $T = \text{BabBabABBAABBACAbBBCAbBB\$}$ , which is presented in Figure 3. Initially,  $\text{SR}_T(P[4..]) = \text{SR}_T(\varepsilon) = [1..25]$ . Since the p-symbol  $P[3] = \text{C}$  did not yet occur in  $P[4..] = \varepsilon$ , we have to compute the LF-mapping for each  $i \in [1..25]$  satisfying  $\text{L}_T[i] \geq \pi(\text{C}) = 1$  to obtain  $\text{SR}_T(\text{C}) = [8..25]$  by Corollary 2.10. Next, we find that the p-symbol  $P[2] = \text{C}$  already occurred in  $P[3..] = \text{C}$ . Thus, we have to compute the LF-mapping of each  $i \in [8..25]$  satisfying  $\text{L}_T[i] = \pi(\text{CC}) = 1$  to obtain  $\text{SR}_T(\text{CC}) = [13..17]$ . Last, we find that  $P[1] = \text{b}$  is an s-symbol. Hence, we have to compute the LF-mapping of each  $i \in [13..17]$  such that  $\text{L}_T[i] = \pi(\text{bCC}) = \text{b}$ , which yields  $\text{SR}_T(P) = [5..6]$ . This concludes the backward search. We find that  $\text{bCC}$  occurs twice, p-matching  $T[22..24]$  and  $T[17..19]$ .

### 2.5. Circular, Parameterized Matching

Throughout, let  $\mathcal{T} = \{T_1, \dots, T_d\} \subset \Sigma^+$  with  $d \in \mathbb{N}_{\geq 1}$  and  $n_k = |T_k|$  for each  $k \in [1..d]$ . For our running example, let  $d = 3$ ,  $T_1 = \text{Bab}$ ,  $T_2 = \text{ABBA}$  and  $T_3 = \text{CAbBB}$ . We are interested in solving the following problem on  $\mathcal{T}$ .

**Problem 2** (EP-COUNT). *Given a set of p-strings  $\mathcal{T}$  and a pattern p-string  $P$ , count each conjugate of the p-strings in  $\mathcal{T}$  whose infinite iteration has a prefix p-matching  $P$ . More formally, return*

$$\sum_{x=1}^d |\{i \in [1..n_x] \mid \text{Rot}(T_x, i)^\omega[1..|P|] \approx P\}|.$$

We call the query described in Problem 2 a *count* query, and call a count query for a pattern no longer than the shortest text a *shortcount* query. As an example, consider the count query for  $P = \text{CCB}$  in our running example  $\mathcal{T}$ , which is a shortcount query. We find that only the conjugates  $\text{Rot}(T_2, 1)^\omega[.3] = \text{BBA}$ ,  $\text{Rot}(T_2, 3)^\omega[.3] = \text{AAB}$  and  $\text{Rot}(T_3, 3)^\omega[.3] = \text{BBC}$  of text p-strings in  $\mathcal{T}$  p-match  $\text{CCB}$ , i.e., the count query returns the value 3.

In the remainder of this section, we present a straight-forward solution that supports shortcount queries based on Lemma 2.8. We assume that the special character  $\$$  does not appear in the input, i.e.,  $\text{rank}_\$(T_k, n_k) = 0$  for each  $k \in [1..d]$ . In the preprocessing, we construct the pBWT index

	$T_1$	$T_1$	$T_2$	$T_2$	$T_3$	$T_3$	$T_3$	$\$$
$T[i]$	B	a	b	B	a	b	A	B
$A_T[i]$	1	1	1	0	0	0	1	1
$B_T[i]$	0	1	0	1	0	1	1	0
$i$	1	2	3	4	5	6	7	8
	9	10	11	12	13	14	15	16
	17	18	19	20	21	22	23	24
	25							

Figure 4: Bit strings  $A_T[i]$  and  $B_T[i]$  augmenting the pBWT index for answering EP-COUNT, instantiated with our running example whose pBWT index is given in Figure 3.

of  $T = T_1^2 \cdots T_d^2 \$$ . If we compute the suffix range  $SR_T(P)$  of a pattern p-string  $P \in \Sigma^*$  no longer than the shortest text, then the length of the suffix range  $SR_T(P)$  is at least the number of occurrences since each conjugate of  $T_1, \dots, T_d$  appears once as a prefix of distinct suffixes of  $T$ . However, an occurrence might be counted twice or cross text boundaries. In our running example, compare the occurrences of  $P = a$  in  $T_1, \dots, T_3$  totaling 1 with the occurrences in  $T = T_1^2 \cdot T_2^2 \cdot T_3^2 \cdot \$$ , which total 2.

To remedy this issue, let  $A_T = 1^{n_1} 0^{n_1} \cdots 1^{n_d} 0^{n_d} 0$  and let  $B_T$  denote a bit string of length  $|T|$  such that  $B_T[i] = 1$  if and only if  $A_T[SA_T[i]] = 1$ . Then the entries of  $B_T$  evaluating to 1 are one-to-one to the conjugates of  $T_1, \dots, T_d$  and consequently  $\text{rank}_1(B_T, r) - \text{rank}_1(B_T, \ell) + B_T[\ell]$  is the shortcount of  $P$ , where  $[\ell..r] = SR_T(P)$ . Thus, the pBWT index of  $T$  augmented with  $A_T$  and  $B_T$  solves EP-COUNT in the case of shortcount queries with the same time complexity as a pBWT index for P-COUNT. See Figure 3 and Figure 4 for the solution to our running example. Note that this index cannot solve count queries in general because a conjugate of a text is considered in EP-COUNT even if it is shorter than  $P$ , e.g., the count query for  $P = T_1^9$  should return 1, but the solution based on Lemma 2.8 for our running example returns 0.

In what follows, we apply Mantaci and colleagues' [10] techniques to the pBWT to obtain a space-efficient index capable of count queries on multiple texts that does without any artificial symbols or an increase in input size.

### 3. An Extension of the pBWT

Throughout, let  $T = T_1 \cdots T_d$ ,  $n = |T|$  and  $\text{LEN}_{\mathcal{T}} = 10^{n_1-1} \cdots 10^{n_k-1} 1$ . Then  $T_k = T[\text{select}_1(\text{LEN}_{\mathcal{T}}, k).. \text{select}_1(\text{LEN}_{\mathcal{T}}, k+1)-1]$  for each  $k \in [1..d]$ . Let  $\text{txt}_{\mathcal{T}}(i) = \text{rank}_1(\text{LEN}_{\mathcal{T}}, i)$ ,  $\text{pos}_{\mathcal{T}}(i) = i - \text{select}_1(\text{LEN}_{\mathcal{T}}, \text{txt}_{\mathcal{T}}(i))$  and  $C_{\mathcal{T}}(i) = \text{Rot}(T_{\text{txt}_{\mathcal{T}}(i)}, \text{pos}_{\mathcal{T}}(i))$  for each  $i \in [1..n]$ , i.e., each position  $i$  in  $T$  corresponds to exactly one of the conjugates  $C_{\mathcal{T}}(i)$  we are considering, which

in turn can be reported by means of this position in the event of a p-matching prefix. In our running example,  $C_{\mathcal{T}}(3) = \text{Rot}(T_{\text{txt}_{\mathcal{T}}(3)}, \text{pos}_{\mathcal{T}}(3)) = \text{Rot}(T_1, 3 - 1) = \mathbf{b}\mathbf{B}\mathbf{a}$  and  $C_{\mathcal{T}}(9) = \text{Rot}(T_{\text{txt}_{\mathcal{T}}(9)}, \text{pos}_{\mathcal{T}}(9)) = \text{Rot}(T_3, 9 - 8) = \mathbf{A}\mathbf{b}\mathbf{B}\mathbf{C}$ . See Figure 5 for more examples.

### 3.1. Conjugate Array

We extend the ideas of Mantaci et al. [10] to p-matching, and introduce a new preorder on  $\Sigma^*$  to compare conjugates across texts. For notational convenience, we define the *rotational prev-encoding*  $\langle V \rangle_r$  of  $V \in \Sigma^*$  by  $\langle V \rangle_r = \langle V^2 \rangle[|V| + 1..]$ .

**Lemma 3.1.** *Let  $V, U \in \Sigma^+$ . Then  $\langle V \rangle_r = \langle U \rangle_r$  if and only if  $V \approx U$ .*

*Proof.* ( $\Rightarrow$ ). Let  $U \not\approx V$ . If  $|U| \neq |V|$ , then  $\langle V \rangle_r \neq \langle U \rangle_r$ . Hence, assume  $|U| = |V|$ . Then  $\langle U \rangle \neq \langle V \rangle$  by Lemma 2.5. Let  $i \in [1..|U|]$  minimal such that  $\langle U \rangle[i] \neq \langle V \rangle[i]$ . If  $\max\{\langle U \rangle[i], \langle V \rangle[i]\} < \infty$ , then  $\langle U \rangle_r[i] = \langle U \rangle[i] \neq \langle V \rangle[i] = \langle V \rangle_r[i]$ , i.e.,  $\langle V \rangle_r \neq \langle U \rangle_r$ . Thus, assume  $\max\{\langle U \rangle[i], \langle V \rangle[i]\} = \infty$ . Without loss of generality,  $\langle U \rangle[i] = \infty$ . Then  $\langle V \rangle_r[i] = \langle V \rangle[i] < i \leq \langle U \rangle_r[i]$ , i.e.,  $\langle V \rangle_r \neq \langle U \rangle_r$ .

( $\Leftarrow$ ). Let  $U \approx V$ . Then  $\langle U \rangle = \langle V \rangle$  by Lemma 2.5. In particular,  $|\langle U \rangle_r| = |\langle V \rangle_r|$ . Assume  $\langle U \rangle_r \neq \langle V \rangle_r$ . Let  $i \in [1..|U|]$  minimal such that  $\langle U \rangle_r \neq \langle V \rangle_r$ . Since  $\langle W \rangle[k] < \infty$  implies  $\langle W \rangle[k] = \langle W \rangle_r[k]$  for each  $W \in \Sigma^+$  and  $k \in [1..|W|]$ ,  $\langle U \rangle[i] = \langle V \rangle[i] = \infty$ , i.e.,  $\min\{\langle V \rangle_r[i], \langle U \rangle_r[i]\} \geq i$ . Without loss of generality,  $\langle V \rangle_r[i] < \langle U \rangle_r[i] \leq |V|$ . Let  $j = i - \langle V \rangle_r[i] + |V|$ . Then  $\langle V \rangle[j] < \infty$ , i.e.,  $V[j]$  occurs at least twice in  $V$ . However,  $U[j] \neq U[i]$ , which contradicts  $U \approx V$  since  $V[i] = V[j]$  is mapped to two different p-symbols or an s-symbol and a p-symbol. Hence,  $\langle U \rangle_r = \langle V \rangle_r$ .  $\square$

For any  $U, V \in \Sigma^*$ , we write  $U \preceq_\omega V$  if and only if there exists an integer  $i \geq 1$  such that  $\langle U^\omega[..i] \rangle < \langle V^\omega[..i] \rangle$  or  $\text{root}(\langle V \rangle_r) = \text{root}(\langle U \rangle_r)$  holds.

**Corollary 3.2.** *The relation  $\preceq_\omega$  defines a total preorder on  $\Sigma^*$ , i.e., the relation is binary, reflexive, transitive and connected.*

We call the relation  $\preceq_\omega$  the  $\omega$ -preorder. We write  $U =_\omega V$  if and only if  $U \preceq_\omega V \wedge V \preceq_\omega U$ , and  $U \prec_\omega V$  if and only if  $U \preceq_\omega V \wedge U \neq_\omega V$ . Note that  $=_\omega$  is an equivalence relation. For example, BBAA  $\prec_\omega$  ABBA, aBBaB  $\prec_\omega$  aBB and CA  $=_\omega$  BCBC  $=_\omega$  ABABAB. The latter example also shows that the  $\omega$ -preorder is not a total order since CA  $\neq$  BCBC  $\wedge$  CA  $=_\omega$  BCBC violates antisymmetry. Adapting Hon et al. [20, Lemma 5], we give a convenient but not necessarily

optimal way to compute the  $\omega$ -preorder of two p-strings making use of the weak periodicity lemma (cf. Crochemore and Rytter's book [45]).

**Lemma 3.3** (Weak Periodicity Lemma). *Let  $j$  and  $k$  be two periods of a string  $V$ . If  $j + k \leq |V|$ , then  $\gcd(j, k)$  is also a period of  $V$ .*

**Lemma 3.4.** *Let  $V, U \in \Sigma^+$  and  $z = \max\{|V|, |U|\}$ . Then  $V =_{\omega} U$  if and only if  $\langle V^{\omega}[..3z] \rangle = \langle U^{\omega}[..3z] \rangle$ .*

*Proof.* Without loss of generality,  $|U| = z$ . Let  $|V| = i$ ,  $|\text{root}(\langle V \rangle_r)| = j$  and  $|\text{root}(\langle U \rangle_r)| = k$ .

( $\Leftarrow$ ). Assume  $\langle V^{\omega}[..3z] \rangle = \langle U^{\omega}[..3z] \rangle$ . Then, on the one hand,

$$\langle \text{Rot}(V, z) \rangle_r^{\omega}[..2z] = \langle V^{\omega}[..3z] \rangle[z+1..] = \langle U^{\omega}[..3z] \rangle[z+1..] = \langle U \rangle_r \cdot \langle U \rangle_r.$$

Since the rotational prev-encoding is commutative with left rotations,  $j$  is a period of  $\langle \text{Rot}(V, z) \rangle_r$ . Consequently, both  $z$  and  $j$  are periods of  $\langle U \rangle_r \cdot \langle U \rangle_r$ . Since  $j + z \leq 2z = |\langle U \rangle_r \cdot \langle U \rangle_r|$ , we can apply Lemma 3.3 and find that  $\gcd(j, z)$  is a period of  $\langle U \rangle_r \cdot \langle U \rangle_r$ . As  $\gcd(j, z)$  divides  $z = |\langle U \rangle_r|$ ,  $\langle U \rangle_r$  can be formed by repeating  $\langle U \rangle_r[.. \gcd(j, z)]$  an integral number of times, which implies  $\gcd(j, z) \geq k$ , i.e.,  $i \geq j \geq k$ . On the other hand,

$$\langle V \rangle_r \cdot \langle V \rangle_r = \langle V^{\omega}[..3z] \rangle[i+1..3i] = \langle U^{\omega}[..3z] \rangle[i+1..3i] = \langle \text{Rot}(U, i) \rangle_r^{\omega}[..2i].$$

Since the rotational prev-encoding is commutative with left rotations,  $k$  is a period of  $\langle \text{Rot}(U, i) \rangle_r$ . As  $k \leq i$ ,  $k$  is a period of  $\langle V \rangle_r \cdot \langle V \rangle_r$  in addition to  $i$ . Then  $k + i \leq 2i = |\langle V \rangle_r \cdot \langle V \rangle_r|$  and Lemma 3.3 imply that  $\gcd(k, i)$  is a period of  $\langle V \rangle_r \cdot \langle V \rangle_r$ . As  $\gcd(k, i)$  divides  $i = |\langle V \rangle_r|$ ,  $\langle V \rangle_r[.. \gcd(k, i)]$  can be repeated an integral number of times to form  $\langle V \rangle_r$ , which implies  $\gcd(k, i) \geq j$ . Consequently,  $k \geq j$ , and therefore  $|\text{root}(\langle V \rangle_r)| = j = k = |\text{root}(\langle U \rangle_r)|$ . In particular,  $\text{root}(\langle V \rangle_r) = \langle V^{\omega}[..3z] \rangle[3z-j+1..] = \langle U^{\omega}[..3z] \rangle[3z-k+1..] = \text{root}(\langle U \rangle_r)$ , i.e.,  $V =_{\omega} U$ .

( $\Rightarrow$ ). Let  $V =_{\omega} U$ . Assume  $\langle V^{\omega}[..3z] \rangle \neq \langle U^{\omega}[..3z] \rangle$  with  $x$  minimal such that  $\langle V^{\omega}[..3z] \rangle[x] \neq \langle U^{\omega}[..3z] \rangle[x]$ . If  $\max\{\langle V^{\omega}[..3z] \rangle[x], \langle U^{\omega}[..3z] \rangle[x]\} < \infty$ , then  $\text{Rot}(\text{root}(\langle V \rangle_r), x-1)[1] = \langle V^{\omega}[..3z] \rangle[x] \neq \langle U^{\omega}[..3z] \rangle[x] = \text{Rot}(\text{root}(\langle U \rangle_r), x-1)[1]$ , which contradicts  $V =_{\omega} U$ . Hence, and without loss of generality, assume  $\langle V^{\omega}[..3z] \rangle[x] = \infty$ . Then  $\langle U^{\omega}[..3z] \rangle[x] < \infty$ ,  $x \leq |\text{root}(\langle V \rangle_r)|$ , and consequently  $\text{root}(\langle U \rangle_r)[x] = \langle U^{\omega}[..3z] \rangle[x] < x \leq \text{root}(\langle V \rangle_r)[x]$ , a contradiction to  $V =_{\omega} U$ . Thus,  $\langle V^{\omega}[..3z] \rangle = \langle U^{\omega}[..3z] \rangle$ .  $\square$

$i$	$C_{\mathcal{T}}(i)$	$\langle C_{\mathcal{T}}(i)^\omega[..15] \rangle$	$CA_{\mathcal{T}}^{-1}[i]$	$CA_{\mathcal{T}}[i]$	$\langle C_{\mathcal{T}}(CA_{\mathcal{T}}[i])^\omega[..15] \rangle$
1	Bab	$\infty ab3ab3ab3ab3ab$	4	2	$ab\infty ab3ab3ab3ab3$
2	abB	$ab\infty ab3ab3ab3ab3$	1	3	$b\infty ab3ab3ab3ab3a$
3	bBa	$b\infty ab3ab3ab3ab3a$	2	10	$b\infty 1\infty b4155b4155$
4	ABBA	$\infty\infty 1313131313131$	10	1	$\infty ab3ab3ab3ab3ab$
5	BBAA	$\infty 1\infty 131313131313$	6	9	$\infty b\infty 1\infty 5b4155b415$
6	BAAB	$\infty\infty 1313131313131$	11	5	$\infty 1\infty 131313131313$
7	AABB	$\infty 1\infty 131313131313$	7	7	$\infty 1\infty 131313131313$
8	CAbBB	$\infty\infty b\infty 155b4155b41$	9	11	$\infty 1\infty b4155b4155b$
9	AbBBC	$\infty b\infty 1\infty 5b4155b415$	5	8	$\infty\infty b\infty 155b4155b41$
10	bBBCA	$b\infty 1\infty\infty b4155b4155$	3	4	$\infty\infty 1313131313131$
11	BBCAb	$\infty 1\infty\infty b4155b4155b$	8	6	$\infty\infty 1313131313131$
12	BCAbB	$\infty\infty\infty b4155b4155b4$	12	12	$\infty\infty\infty b4155b4155b4$

Figure 5: The conjugate array  $CA_{\mathcal{T}}$  of our running example  $\mathcal{T} = \{Bab, ABBA, CABBB\}$ .

**Corollary 3.5.** *Let  $V, U \in \Sigma^+$  and  $z = \max\{|V|, |U|\}$ . Then  $V \prec_{\omega} U$  if and only if  $\langle V^\omega[..3z] \rangle < \langle U^\omega[..3z] \rangle$ .*

We generalize the suffix array similarly to Boucher et al. [26]. The *conjugate array*  $CA_{\mathcal{T}}$  of  $\mathcal{T}$  is a string of length  $n$  over  $[1..n]$  such that  $CA_{\mathcal{T}}[i] = j$  if and only if  $i-1 = |\{k \in [1..n] \mid C_{\mathcal{T}}(k) \prec_{\omega} C_{\mathcal{T}}(j) \vee (C_{\mathcal{T}}(k) =_{\omega} C_{\mathcal{T}}(j) \wedge k < j)\}|$ , i.e.,  $i-1$  is the number of all conjugates smaller than  $C_{\mathcal{T}}(j)$  according to  $\omega$ -preorder, where we break ties first with respect to text index, and then with respect to text position. By resolving all ties this way, we enforce a linear order on the conjugates considered in  $\mathcal{T}$ , and conclude that  $CA_{\mathcal{T}}$  is well-defined. Since  $CA_{\mathcal{T}}$  is a permutation, its inverse, which we denote by  $CA_{\mathcal{T}}^{-1}$ , is also well-defined. For the conjugate array of our running example, see Figure 5.

The *conjugate range*  $CR_{\mathcal{T}}(P)$  of a pattern  $P \in \Sigma^*$  of length  $m$  and  $\mathcal{T}$  is a maximal interval  $[\ell..r] \subseteq [1..n]$  such that  $P \approx C_{\mathcal{T}}(CA_{\mathcal{T}}[i])^\omega[..m]$  for every  $i \in [\ell..r]$ . Note that  $CR_{\mathcal{T}}(\varepsilon) = [1..n]$ . The following result is a consequence of Lemma 2.5.

**Corollary 3.6.** *Let  $T_1, \dots, T_d \in \Sigma^+$ ,  $n = |T_1 \cdots T_d|$ ,  $P \in \Sigma^*$ ,  $m = |P|$  and  $[\ell..r] = CR_{\mathcal{T}}(P)$ . Then  $P \approx C_{\mathcal{T}}(CA_{\mathcal{T}}[i])^\omega[..m]$  if and only if  $i \in CR_{\mathcal{T}}(P)$ .*

We call the computation of  $CR_{\mathcal{T}}(P)$  for some pattern  $P \in \Sigma^*$  the *backward search* for  $P$  in  $\mathcal{T}$  and the length of  $CR_{\mathcal{T}}(P)$  is the return value to a count query for  $P$  in  $\mathcal{T}$  by Corollary 3.6, i.e., the solution to EP-COUNT.

$i$	$T[i]$	$\text{LEN}_{\mathcal{T}}[i]$	$\text{ENC}_{\mathcal{T}}[i]$	$\text{PRV}_{\mathcal{T}}[i]$	$\text{prev}_{\mathcal{T}}(i)$
1	B	1	1	1	3
2	a	0	a	0	1
3	b	0	b	0	2
4	A	1	2	1	5
5	B	0	1	0	4
6	B	0	2	1	7
7	A	0	1	0	6
8	C	1	3	1	12
9	A	0	3	0	8
10	b	0	b	0	9
11	B	0	1	0	10
12	B	0	3	0	11

Figure 6: The permutation  $\text{prev}_{\mathcal{T}}$  of our running example and its associated data structures. Here,  $\text{LEN}_{\mathcal{T}}[13] = \text{PRV}_{\mathcal{T}}[13] = 1$ .

For our running example, we consider the pattern p-strings  $P_1 = \text{ACAB}$ ,  $P_2 = \text{CCB}$  and  $P_3 = \text{a}$ . Then  $\text{CR}_{\mathcal{T}}(P_1) = \emptyset$ ,  $\text{CR}_{\mathcal{T}}(P_2) = [6..8]$  and  $\text{CR}_{\mathcal{T}}(P_3) = [1..1]$ , and the count queries for  $P_1$ ,  $P_2$  and  $P_3$  return 0, 3 and 1, respectively. In what follows, we define the necessary tools to allow for an efficient backward search.

### 3.2. LF-Mapping

Unlike Mantaci et al. [10], we did not assume our texts or their encodings to be primitive, e.g.,  $\llbracket T_2 \rrbracket = 2121$  is not primitive. This warrants the following. Let  $\text{ENC}_{\mathcal{T}} = \llbracket T_1 \rrbracket \cdots \llbracket T_d \rrbracket = \text{root}(\llbracket T_1 \rrbracket)^{\exp(\llbracket T_1 \rrbracket)} \cdots \text{root}(\llbracket T_d \rrbracket)^{\exp(\llbracket T_d \rrbracket)}$ . We define a bit string  $\text{PRV}_{\mathcal{T}}$  of length  $n + 1$  marking the positions of  $\text{ENC}_{\mathcal{T}}$  with a 1 if the position is a starting position of an encoded text root, and set  $\text{PRV}_{\mathcal{T}}[n + 1] = 1$  as a sentinel. Note that  $\text{LEN}_{\mathcal{T}} = \text{PRV}_{\mathcal{T}}$  if and only if  $\llbracket T_k \rrbracket$  is primitive for every  $k \in [1..d]$ . For every  $i \in [1..n]$ , let  $\text{prev}_{\mathcal{T}}(i) = \text{select}_1(\text{PRV}_{\mathcal{T}}, \text{rank}_1(\text{PRV}_{\mathcal{T}}, i) + 1) - 1$  if  $\text{PRV}_{\mathcal{T}}[i] = 1$ , or  $\text{prev}_{\mathcal{T}}(i) = i - 1$  otherwise, i.e., the permutation  $\text{prev}_{\mathcal{T}}$  of  $[1..n]$  moves circularly to the previous position inside a root of an encoded input text. See Figure 6 for an example. Note that  $\text{prev}_{\mathcal{T}}$  can be derived from the rotational  $\text{prev}$ -encoding instead by Lemma 2.5 and Lemma 3.1.

Then the *LF-mapping*  $\text{LF}_{\mathcal{T}}$  of  $\mathcal{T}$  is a string of length  $n$  over  $[1..n]$  such that  $\text{LF}_{\mathcal{T}}[i] = \text{CA}_{\mathcal{T}}^{-1}[\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])]$  for each  $i \in [1..n]$ . Since the LF-mapping is

a permutation, it has an inverse  $\text{LF}_{\mathcal{T}}^{-1}$ , which we call the *FL-mapping* of  $\mathcal{T}$  and denote by  $\text{FL}_{\mathcal{T}}$ . The LF-mapping will be at the core of the backward search. However, storing LF-mapping and FL-mapping in their plain form creates the need for two integer arrays of length  $n$  and entries of  $\lg n$  bits. To represent both the LF- and FL-mapping of  $\mathcal{T}$  in compact space, we need one more result on the  $\omega$ -preorder of two strings, which now follows.

**Lemma 3.7.** *Let  $U, V \in \Sigma^+$ ,  $e = \text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V, 1))$  and assume  $\text{Rot}(U, 1) \prec_\omega \text{Rot}(V, 1)$ . Then  $V \prec_\omega U$  if and only if  $\pi(V) < \pi(U) \wedge (\pi(V) \in \Sigma_s \vee (\pi(U), \pi(V) \notin \Sigma_s \wedge \pi(V) \leq e))$ .*

*Proof.* Application of Lemma 2.6 and Corollary 3.5.  $\square$

Finally, we introduce our representation of the LF- and FL-mapping, which consists of two strings  $L_{\mathcal{T}}$  and  $F_{\mathcal{T}}$ , and are defined as follows. First, the *extended parameterized Burrows–Wheeler Transform (epBWT)*  $L_{\mathcal{T}}$  of  $\mathcal{T}$  is defined as a string of length  $n$  over  $[1..n] \cup \Sigma_s$  such that  $L_{\mathcal{T}}[i] = \pi(C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[i]]))$  for each  $i \in [1..n]$ .

**Lemma 3.8** (Non-crossing Property). *Let  $\mathcal{T} \subset \Sigma^+$  and  $i, j \in [1..n]$  with  $i < j$ , where  $n$  is the accumulated length of all strings of  $\mathcal{T}$ . If  $L_{\mathcal{T}}[i] = L_{\mathcal{T}}[j]$ , then  $\text{LF}_{\mathcal{T}}[i] < \text{LF}_{\mathcal{T}}[j]$ .*

*Proof.* Let  $L_{\mathcal{T}}[i] = L_{\mathcal{T}}[j]$ . Since  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[k]) = \text{Rot}(C_{\mathcal{T}}(\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[k])), 1)$  for every  $k \in [1..n]$ , Lemma 3.7 implies

$$C_{\mathcal{T}}(\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])) \preceq_\omega C_{\mathcal{T}}(\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])) \Leftrightarrow C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i]) \preceq_\omega C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j]).$$

As  $i < j$ ,  $C_{\mathcal{T}}(\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])) \preceq_\omega C_{\mathcal{T}}(\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j]))$ . Then by definition of the conjugate array,  $\text{CA}_{\mathcal{T}}^{-1}[\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])] < \text{CA}_{\mathcal{T}}^{-1}[\text{prev}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])]$ , i.e.,  $\text{LF}_{\mathcal{T}}[i] < \text{LF}_{\mathcal{T}}[j]$ .  $\square$

Second, let  $F_{\mathcal{T}}$  denote a string of length  $n$  over  $[1..n] \cup \Sigma_s$  such that  $F_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[i]] = L_{\mathcal{T}}[i]$  for each  $i \in [1..n]$ . By the following result,  $L_{\mathcal{T}}$  and  $F_{\mathcal{T}}$  suffice to compute both LF- and FL-mapping of  $\mathcal{T}$ .

**Corollary 3.9.** *Let  $\mathcal{T} \subset \Sigma^+$  and  $n$  be the accumulated length of all strings of  $\mathcal{T}$ . Then*

$$\begin{aligned} \text{LF}_{\mathcal{T}}[i] &= \text{select}_{L_{\mathcal{T}}[i]}(F_{\mathcal{T}}, \text{rank}_{L_{\mathcal{T}}[i]}(L_{\mathcal{T}}, i)) \text{ and} \\ \text{FL}_{\mathcal{T}}[i] &= \text{select}_{F_{\mathcal{T}}[i]}(L_{\mathcal{T}}, \text{rank}_{F_{\mathcal{T}}[i]}(F_{\mathcal{T}}, i)) \text{ for each } i \in [1..n]. \end{aligned}$$

We also achieve reversibility to a degree.

**Corollary 3.10.** *Let  $\mathcal{T} = \{T_1, \dots, T_d\} \subset \Sigma^+$ . Given  $F_{\mathcal{T}}$  and  $L_{\mathcal{T}}$ , we can restore the conjugacy classes of  $\mathcal{T}$  up to  $=_{\omega}$ -equivalence. In particular, we can restore the conjugacy classes of  $\mathcal{T}$  up to p-match equivalence if  $\llbracket T \rrbracket_1, \dots, \llbracket T \rrbracket_d$  are primitive.*

*Proof.* We obtain the text-roots of the  $\llbracket .. \rrbracket$ -encoded input texts from the cycles of the LF-mapping and either  $F_{\mathcal{T}}$  or  $L_{\mathcal{T}}$ . Since the LF-mapping can be computed from  $F_{\mathcal{T}}$  and  $L_{\mathcal{T}}$  by Corollary 3.9, the claims follow from Lemma 2.5 and Lemma 3.1.  $\square$

The LF-mapping of our running example can be found in Figure 7. To restore the input texts, we examine its cycles. We start with  $F_{\mathcal{T}}[1] = a$  and compute  $LF_{\mathcal{T}}[1] = 4$ , which yields  $F_{\mathcal{T}}[4] \cdot F_{\mathcal{T}}[1] = 1a$ . Next, we proceed to compute  $LF_{\mathcal{T}}[4] = 2$  resulting in  $F_{\mathcal{T}}[2] \cdot F_{\mathcal{T}}[4] \cdot F_{\mathcal{T}}[1] = b1a$ . Since  $LF_{\mathcal{T}}[2] = 1$ , the cycle ends; we have found a root of a  $\llbracket .. \rrbracket$ -encoded conjugate of an input text string, namely  $\llbracket \text{Rot}(T_1, 2) \rrbracket = b1a$ . If we do this for all the other cycles of the LF-mapping, we obtain  $F_{\mathcal{T}}[8] \cdot F_{\mathcal{T}}[12] \cdot F_{\mathcal{T}}[9] \cdot F_{\mathcal{T}}[5] \cdot F_{\mathcal{T}}[3] = 1333b = \llbracket \text{Rot}(T_3, 3) \rrbracket$  and  $F_{\mathcal{T}}[10] \cdot F_{\mathcal{T}}[6] = F_{\mathcal{T}}[11] \cdot F_{\mathcal{T}}[7] = 21 = \text{root}(\llbracket T_2 \rrbracket)$ .

To enable restoration up to p-match and conjugacy class without the assumption that the  $\llbracket .. \rrbracket$ -encoded input texts are primitive, we can proceed as follows. For each  $k \in [1..d]$ , our goal is to ignore all but one cycle in the LF-mapping corresponding to  $T_k$  and store  $\exp(\llbracket T_k \rrbracket)$  in an easily accessible way. Let  $D_{\mathcal{T}}$  denote a string of length  $n$  over  $[1..n]$  such that  $D_{\mathcal{T}}[i] = \exp(\llbracket T_{\text{txt}_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])} \rrbracket)$  if  $\text{select}_1(\text{LEN}_{\mathcal{T}}, \text{rank}_1(\text{LEN}_{\mathcal{T}}, \text{CA}_{\mathcal{T}}[i])) = \text{select}_1(\text{PRV}_{\mathcal{T}}, \text{rank}_1(\text{PRV}_{\mathcal{T}}, \text{CA}_{\mathcal{T}}[i]))$ , and  $D_{\mathcal{T}}[i] = 0$  otherwise, i.e., non-zero entries of  $D_{\mathcal{T}}$  correspond to entries in the first iteration of a root of a  $\llbracket .. \rrbracket$ -encoded text and yield the exponent of the corresponding  $\llbracket .. \rrbracket$ -encoded text. Restoration of up to p-match and conjugacy class is then straight-forward.

Note that we have to construct  $D_{\mathcal{T}}$  during the construction of  $F_{\mathcal{T}}$  and  $L_{\mathcal{T}}$  since we are unable to discern the cycles of any two texts that have  $\omega$ -equal conjugates. Moreover, since  $\sum_{i=1}^n D_{\mathcal{T}}[i] = n$ , we can encode  $D_{\mathcal{T}}$  as a bit string  $D'_{\mathcal{T}}$  of length  $2n$ , where entries of  $D_{\mathcal{T}}$  are encoded as unary zeros that are separated by '1's with an additional '1' as a sentinel. Then  $D_{\mathcal{T}}[i] = \text{rank}_0(D'_{\mathcal{T}}, \text{select}_1(D'_{\mathcal{T}}, i))$ . Thus,  $D_{\mathcal{T}}$  can be stored in  $O(n)$  bits of space if its binary encoding  $D'_{\mathcal{T}}$  is represented by the data structure of Lemma 2.1 or Lemma 2.3.

In our running example,  $\llbracket T_2 \rrbracket$  is not primitive. Consequently,  $D_{\mathcal{T}} = 111112011201$ , which we encode as  $D'_{\mathcal{T}} = 010101010100110101001101$ . The LF-mapping yields four cycles as before. However, for the cycle (7 11) we get  $D_{\mathcal{T}}[7] = D_{\mathcal{T}}[11] = 0$  indicating that the cycle does not correspond to the first iteration of a root within a  $\llbracket .. \rrbracket$ -encoded indexed text. Thus, we discard the cycle. The cycle (6 10), on the other hand, corresponds to the first iteration of a root within a  $\llbracket .. \rrbracket$ -encoded indexed text, and  $D_{\mathcal{T}}[6] = D_{\mathcal{T}}[10] = 2$  yields the number of iterations to restore the text  $T_2$  up to conjugacy class and p-match equivalence.

### 3.3. Backward Search

The backward search for a pattern  $P \in \Sigma^*$  in  $\mathcal{T}$  consists of updating the conjugate range of  $P[i+1..]$  to that of  $P[i..]$  for each  $i \in [1..|P|]$ , where  $\text{CR}_{\mathcal{T}}(P[|P|+1..]) = \text{CR}_{\mathcal{T}}(\varepsilon) = [1..n]$ . In what follows, we present two different methods to compute these conjugate range updates, where the second one is motivated by our construction algorithm of the epBWT presented in Section 4. The first one has better complexities while the second one allows for an efficient construction and the possibilities to add and remove indexed texts. For the first method, we adapt the findings of Kim and Cho [14].

**Lemma 3.11.** *Let  $\mathcal{T} \subset \Sigma^+$  and  $P \in \Sigma^*$ . Then Algorithm 1 correctly computes  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..])$  for each  $i \in [1..|P|]$ .*

*Proof.* Let  $m = |P|$ ,  $i \in [1..m]$ ,  $h = \pi(P[i..])$ ,  $[\ell..r] = \text{CR}_{\mathcal{T}}(P[i+1..])$  and  $c = |P[i+1..]|_p$ . Algorithm 1 takes  $[\ell..r]$ ,  $h$ ,  $c$ ,  $F_{\mathcal{T}}$ ,  $L_{\mathcal{T}}$  and  $\text{LF}_{\mathcal{T}}$  as input.

By definition,  $h \leq c + 1$ . For the computation of  $r'$  and the length  $\delta = r' - \ell' + 1$  of  $[\ell'..r']$ , we consider two cases depending on  $h$ , with the first being treated from Line 2 through Line 6 and the second from Line 7 through Line 11. If  $\delta > 0$ , then  $\ell' = r' - \delta + 1$ , which is computed in Line 12.

**Case 1.** Assume  $h \leq c$ , i.e.,  $h \in \Sigma_s$  or  $h \notin \Sigma_s$  and there exists some  $j \in [i+1..m]$  such that  $P[i] = P[j]$ . Then  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[j]])^\omega[..m-i+1] \approx P[i..]$  if and only if  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^\omega[..m-i] \approx P[i+1..]$  and  $L_{\mathcal{T}}[j] = h$ , for each  $j \in [1..n]$ . Thus,  $\text{LF}_{\mathcal{T}}[j] \in [\ell'..r']$  if and only if  $j \in [\ell..r]$  and  $L_{\mathcal{T}}[j] = h$ . Then  $\delta = \text{rnkcnt}_{L_{\mathcal{T}}}(\ell, r, h, h)$ , which is computed in Line 3. If  $\delta \leq 0$ , then we return an empty interval in Line 5. If  $\delta > 0$ , then the non-crossing property of Lemma 3.8 already implies  $r' = \text{LF}_{\mathcal{T}}[\text{select}_h(L_{\mathcal{T}}, \text{rank}_h(L_{\mathcal{T}}, r))]$ , which is computed in Line 6.

**Case 2.** Assume  $h = c+1$ , i.e.,  $h \notin \Sigma_s$  and  $P[i] \neq P[j]$  for each  $j \in [i+1..m]$ . Then  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[j]])^\omega[..m-i+1] \approx P[i..]$  if and only if  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^\omega[..m-i+1] \approx P[i..]$ .

---

**Algorithm 1:** Updating  $[\ell..r] = \text{CR}_{\mathcal{T}}(P[i+1..])$  to  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..])$ . Here,  $h = \pi(P[i..])$  and  $c = |P[i+1..]|_p$ .

---

```

1 Function crangeupdrmq( $[\ell..r]$ ,  $h$ ,  $c$ ,  $\text{F}_{\mathcal{T}}$ ,  $\text{L}_{\mathcal{T}}$ ,  $\text{LF}_{\mathcal{T}}$ ):
2   if  $h \leq c$  then                                // Case 1 in Lemma 3.11
3      $\delta \leftarrow \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, r, h, h);$ 
4     if  $\delta \leq 0$  then
5       return  $[r+1..r];$                       // empty interval
6      $r' \leftarrow \text{LF}_{\mathcal{T}}[\text{select}_h(\text{L}_{\mathcal{T}}, \text{rank}_h(\text{L}_{\mathcal{T}}, r))];$ 
7   else                                         // Case 2 in Lemma 3.11
8      $\delta \leftarrow \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, r, h, \sigma_p);$ 
9     if  $\delta \leq 0$  then
10      return  $[r+1..r];$                       // empty interval
11       $r' \leftarrow \text{LF}_{\mathcal{T}}[\text{RMQ}_{\text{LF}_{\mathcal{T}}}(\ell, r)];$ 
12     $\ell' \leftarrow r' - \delta + 1;$ 
13    return  $[\ell'..r'];$ 

```

---

$i] \approx P[i+1..]$  and  $\text{L}_{\mathcal{T}}[j] > c$ , for each  $j \in [1..n]$ . Thus,  $\text{LF}_{\mathcal{T}}[j] \in [\ell'..r']$  if and only if  $j \in [\ell..r]$  and  $\text{L}_{\mathcal{T}}[j] > c$ . Consequently,  $\delta = \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, r, h, \sigma_p)$ , which is computed in Line 8. If  $\delta \leq 0$ , then we return an empty interval in Line 10. By Lemma 3.7,  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[j]]) \prec_{\omega} C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[k]])$  for each  $j, k \in [\ell..r]$  with  $\text{L}_{\mathcal{T}}[j] \leq c$  and  $\text{L}_{\mathcal{T}}[k] > c$ . Hence,  $r' = \max\{\text{LF}_{\mathcal{T}}[j] \mid j \in [\ell..r]\}$ , which we compute in Line 11.  $\square$

Finally, we preprocess the pattern such that we can access the encoding of its suffixes in constant time.

**Lemma 3.12** ([13, Section 5.2]). *Given a  $p$ -string  $P$  of length  $m$ , we can process  $P$  in  $O(m \lg \sigma)$  time such that we can subsequently compute  $\pi(P[i..])$  and  $|P[i+1..]|_p$  in  $O(1)$  time, for every  $i \in [1..m]$ .*

**Theorem 3.13.** *Let  $\mathcal{T} \subset \Sigma^+$  and  $n$  be the accumulated length of all strings of  $\mathcal{T}$ . There exists an index solving EP-COUNT in  $O(m \lg \sigma)$  time, where  $m$  is the length of the queried pattern  $p$ -string. The index takes  $2n \lg \sigma + 2n + O(n \frac{\lg \lg n}{\lg n})$  bits of space.*

*Proof.* The index consists of representations of  $\text{F}_{\mathcal{T}}$ ,  $\text{L}_{\mathcal{T}}$  and  $\text{LF}_{\mathcal{T}}$ . We represent  $\text{F}_{\mathcal{T}}$  and  $\text{L}_{\mathcal{T}}$  by the data structure of Lemma 2.1 and  $\text{LF}_{\mathcal{T}}$  by the

data structure of Lemma 2.2, which leads to the claimed space complexity. Let  $P \in \Sigma^*$  with  $m = |P|$ . We preprocess  $P$  with Lemma 3.12. Then the claim regarding count queries follows from Lemma 2.1, Lemma 2.2, and Lemma 3.11.  $\square$

For the second method, we took inspiration from Hashimoto et al. [39] and Iseri et al. [15]. Let  $\text{LCP}_{\mathcal{T}}^\infty$  denote a string of length  $n$  over  $[0..\sigma_p]$  such that  $\text{LCP}_{\mathcal{T}}^\infty[1] = 0$  and  $\text{LCP}_{\mathcal{T}}^\infty[i] = \text{lcp}^\infty(C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i]), C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i-1]))$  for  $i \in [2..n]$ .

**Lemma 3.14.** *Let  $\mathcal{T} \subset \Sigma^+$  and  $n$  be the accumulated length of all strings of  $\mathcal{T}$ . Then for each  $i, j \in [1..n]$  with  $i < j$ ,*

$$\text{lcp}^\infty(C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i]), C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])) = \text{RNV}_{\text{LCP}_{\mathcal{T}}^\infty}(i+1, j, -1).$$

*Proof.* Let  $U, V, W \in \Sigma^+$  and  $z = \max\{|U|, |V|, |W|\}$ . If  $\langle U^\omega[..3z] \rangle \leq \langle V^\omega[..3z] \rangle \leq \langle W^\omega[..3z] \rangle$ , then we conclude that  $\text{lcp}(\langle U^\omega[..3z] \rangle, \langle W^\omega[..3z] \rangle) = \min\{\text{lcp}(\langle U^\omega[..3z] \rangle, \langle V^\omega[..3z] \rangle), \text{lcp}(\langle V^\omega[..3z] \rangle, \langle W^\omega[..3z] \rangle)\}$ . By Lemma 3.4,  $U \preceq_\omega V \preceq_\omega W$  implies  $\text{lcp}^\infty(U, W) = \min\{\text{lcp}^\infty(U, V), \text{lcp}^\infty(V, W)\}$ .

Let  $i, j \in [1..n]$  and  $i < j$ . Then  $\text{lcp}^\infty(C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i]), C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])) = \min\{\text{LCP}_{\mathcal{T}}^\infty[k] \mid k \in [i+1..j]\} = \text{RNV}_{\text{LCP}_{\mathcal{T}}^\infty}(i+1, j, -1)$ .  $\square$

**Lemma 3.15.** *Let  $\mathcal{T} \subset \Sigma^+$  and  $P \in \Sigma^*$ . Then Algorithm 2 correctly computes  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..])$  for each  $i \in [1..|P|]$ .*

*Proof.* Let  $m = |P|$ ,  $i \in [1..m]$ ,  $h = \pi(P[i..])$ ,  $[\ell..r] = \text{CR}_{\mathcal{T}}(P[i+1..])$  and  $c = |P[i+1..]|_p$ . Algorithm 2 takes  $[\ell..r]$ ,  $h$ ,  $c$ ,  $\mathcal{F}_{\mathcal{T}}$ ,  $\mathcal{L}_{\mathcal{T}}$  and  $\text{LCP}_{\mathcal{T}}^\infty$  as input.

By definition,  $h \leq c + 1$ . The case where  $h \leq c$  is treated from Line 2 through Line 6, and the proof for correctness is the same as that of Case 1 in the proof of Lemma 3.11. Thus, assume  $h = c + 1$ , which we treat from Line 7 through Line 16. Then  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[\text{LF}_{\mathcal{T}}[j]])^\omega[..m-i+1] \approx P[i..]$  if and only if  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^\omega[..m-i] \approx P[i+1..]$  and  $\mathcal{L}_{\mathcal{T}}[j] > c$ , for each  $j \in [1..n]$ . Hence,  $\text{LF}_{\mathcal{T}}[j] \in [\ell'..r']$  if and only if  $j \in [\ell..r]$  and  $\mathcal{L}_{\mathcal{T}}[j] > c$ . Consequently,  $\delta = |\text{CR}_{\mathcal{T}}(P[i..])| = \text{rnkcnt}_{\mathcal{L}_{\mathcal{T}}}(\ell, r, h, \sigma_p)$ , which is computed in Line 8. If  $\delta \leq 0$ , then we return an empty interval in Line 10.

If  $\delta > 0$ , then we compute the highest index  $x$  of  $\mathcal{L}_{\mathcal{T}}$  such that  $x \in [\ell..r]$  and  $\mathcal{L}_{\mathcal{T}}[x]$  takes the lowest value  $v$  in  $\mathcal{L}_{\mathcal{T}}[\ell..r]$  higher than  $c$  in Line 11 and Line 12. We then proceed to compute  $y = \text{LF}_{\mathcal{T}}[x] - \ell'$ , i.e., we count the entries  $j \in [\ell..r]$  with  $\mathcal{L}_{\mathcal{T}}[j] \geq h$  and  $j \neq x$  such that  $C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[j]) \preceq_\omega$

---

**Algorithm 2:** Updating  $[\ell..r] = \text{CR}_{\mathcal{T}}(P[i+1..])$  to  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..])$ . Here,  $h = \pi(P[i..])$  and  $c = |P[i+1..]|_p$ .

---

```

1 Function crangeupdLCP( $[\ell..r]$ ,  $h$ ,  $c$ ,  $\text{F}_{\mathcal{T}}$ ,  $\text{L}_{\mathcal{T}}$ ,  $\text{LCP}_{\mathcal{T}}^{\infty}$ ):
2   if  $h \leq c$  then
3      $\delta \leftarrow \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, r, h, h);$ 
4     if  $\delta \leq 0$  then
5       return  $[r+1..r];$  // empty interval
6      $r' \leftarrow \text{LF}_{\mathcal{T}}[\text{select}_h(\text{L}_{\mathcal{T}}, \text{rank}_h(\text{L}_{\mathcal{T}}, r))];$ 
7   else // Case in Lemma 3.15
8      $\delta \leftarrow \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, r, h, \sigma_p);$ 
9     if  $\delta \leq 0$  then
10      return  $[r+1..r];$  // empty interval
11       $v \leftarrow \text{RNV}_{\text{L}_{\mathcal{T}}}(\ell, r, c);$ 
12       $x \leftarrow \text{select}_v(\text{L}_{\mathcal{T}}, \text{rank}_v(\text{L}_{\mathcal{T}}, r));$ 
13       $[\ell''..r''] \leftarrow \text{MI}_{\text{LCP}_{\mathcal{T}}^{\infty}}(x, v);$ 
14       $y \leftarrow \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\max\{\ell, \ell''\}, x-1, v, v);$ 
15       $y \leftarrow y + \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, \ell''-1, v, \sigma_p);$ 
16       $r' \leftarrow \text{LF}_{\mathcal{T}}[x] + \delta - y - 1;$ 
17       $\ell' \leftarrow r' - \delta + 1;$ 
18      return  $[\ell'..r'];$ 

```

---

$C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[x])$ . To that end, let  $[\ell''..r''] = \text{MI}_{\text{LCP}_{\mathcal{T}}^{\infty}}(x, v)$  as computed in Line 13. Then  $\text{lcp}^{\infty}(C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[x]), C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])) \geq v$  if and only if  $j \in [\ell''..r'']$ .

By Lemma 3.7 and choice of  $x$ ,  $C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[x]) \prec_{\omega} C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[j])$  for each  $j \in [x+1..r]$  with  $\text{L}_{\mathcal{T}}[j] \geq h$ , i.e., we do not need to consider the entries of  $[x+1..r]$  in our computation of  $y$ . For each  $j \in [\max\{\ell, \ell''\}..x-1]$  with  $\text{L}_{\mathcal{T}}[j] \geq h$ ,  $C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[j]) \preceq_{\omega} C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[x])$  if and only if  $\text{L}_{\mathcal{T}}[j] = v$  by Lemma 3.7 and choice of  $[\ell''..r'']$ . Again by Lemma 3.7 and choice of  $[\ell''..r'']$ ,  $C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[j]) \preceq_{\omega} C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[x])$  for each  $j \in [\ell..l''-1]$  with  $\text{L}_{\mathcal{T}}[j] \geq h$ .

Then the non-crossing property implies  $y = \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\max\{\ell, \ell''\}, x-1, v, v) + \text{rnkcnt}_{\text{L}_{\mathcal{T}}}(\ell, \ell''-1, v, \sigma_p)$ , which is computed in Line 14 and Line 15. The computation of  $r'$  and  $\ell'$  in Line 16 and Line 17, respectively, is like before.  $\square$

**Theorem 3.16.** Let  $\mathcal{T} \subset \Sigma^+$  and  $n$  be the accumulated length of all strings

$i$	$\text{CA}_{\mathcal{T}}[i]$	$C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])$	$F_{\mathcal{T}}[i]$	$L_{\mathcal{T}}[i]$	$\text{LCP}_{\mathcal{T}}^{\infty}[i]$	$\text{LF}_{\mathcal{T}}[i]$	$\text{FL}_{\mathcal{T}}[i]$
1	2	abB	a	1	0	4	2
2	3	bBa	b	a	0	1	4
3	10	bBBCA	b	3	1	5	8
4	1	Bab	1	b	0	2	1
5	9	AbBBC	3	3	1	9	3
6	5	BBAA	1	2	1	10	10
7	7	AABB	1	2	2	11	11
8	11	BBCAb	1	b	2	3	12
9	8	CAbBB	3	3	1	12	5
10	4	ABBA	2	1	2	6	6
11	6	BAAB	2	1	2	7	7
12	12	BCAbB	3	1	2	8	9

Figure 7: The epBWT index of our running example  $\mathcal{T} = \{\text{Bab}, \text{ABBA}, \text{CAbBB}\}$ .

of  $\mathcal{T}$ . There exists an index solving EP-COUNT in  $O(m \frac{\lg \sigma \lg n}{\lg \lg n})$  time, where  $m$  is the length of the queried pattern  $p$ -string. The index takes  $O(n \lg \sigma)$  bits of space.

*Proof.* The index consists of representations of  $F_{\mathcal{T}}$ ,  $L_{\mathcal{T}}$  and  $\text{LCP}_{\mathcal{T}}^{\infty}$ . We represent  $L_{\mathcal{T}}$  and  $\text{LCP}_{\mathcal{T}}^{\infty}$  by the data structure of Lemma 2.4, and  $F_{\mathcal{T}}$  by the data structure of Lemma 2.3, which leads to the claimed space complexity. Let  $P \in \Sigma^*$  with  $m = |P|$ . We preprocess  $P$  with Lemma 3.12. Then the claim regarding count queries follows from Lemma 2.3, Lemma 2.4, and Lemma 3.15.  $\square$

We call the index described in Theorem 3.16 the *epBWT index of  $\mathcal{T}$* . See Figure 7 for the epBWT index of our running example. Note that the epBWT index of  $\{V\}$  is simply the pBWT index of  $V$  if  $V \in \Sigma^+$ ,  $\text{rank}_{\$}(V, |V|) = 1$  and  $V[|V|] = \$$ .

### 3.4. Showcase

We give an exemplary backward search for  $P = \text{bCCA}$  with the epBWT index of our running example  $\mathcal{T} = \{\text{Bab}, \text{ABBA}, \text{CAbBB}\}$ , which is presented in Figure 7.

Since  $\varepsilon$  is a prefix of every conjugate considered in  $\mathcal{T}$ , we initially have  $[\ell..r] = \text{CR}_{\mathcal{T}}(P[5..]) = \text{CR}_{\mathcal{T}}(\varepsilon) = [1..12]$ . For the computation of  $\text{CR}_{\mathcal{T}}(\text{A})$ , we

observe that the p-symbol  $\mathbf{A}$  did not yet appear in  $P[5..]$ . Computationally, we can infer the appearance of  $\mathbf{A}$  in  $P[5..]$  from  $h = \pi(P[4..]) = \pi(\mathbf{A}) = 1$ ,  $c = |P[5..]|_p = |\varepsilon|_p = 0$  and  $h = 1 = 0 + 1 = c + 1$ . Algorithm 2 consequently enters Line 7. Next, the proof of Lemma 3.15 tells us that the conjugates corresponding to entries  $i \in [\ell..r] = [1..12]$  with  $L_{\mathcal{T}}[i] \geq h = 1$  are the left rotations of conjugates whose prefix of length 1 p-matches  $\mathbf{A}$ . Thus, we compute  $\delta = \text{rnkcnt}_{L_{\mathcal{T}}}(1, 12, 1, 3) = 9$  in Line 8 of Algorithm 2 to find the length of  $\text{CR}_{\mathcal{T}}(\mathbf{A})$ . In particular, we find that  $\text{CR}_{\mathcal{T}}(\mathbf{A})$  is not empty. It remains to find either the left or right boundary of  $\text{CR}_{\mathcal{T}}(\mathbf{A})$  and compute the other from the length. To that end, we compute  $v = \text{RNV}_{L_{\mathcal{T}}}(1, 12, 0) = 1$ ,  $x = \text{select}_1(L_{\mathcal{T}}, \text{rank}_1(L_{\mathcal{T}}, 12)) = 12$  and  $[\ell''..r''] = \text{MI}_{\text{LCP}_{\mathcal{T}}^{\infty}}(12, 1) = [4..12]$ , which are the lowest value in  $L_{\mathcal{T}}[1..12]$  larger than  $c$ , the highest index at which the value  $v$  occurs and the interval whose entries correspond to prev-encoded conjugates that have at least  $v$  occurrences of  $\infty$  in the longest common prefix with the prev-encoded conjugate corresponding to  $x$ , respectively. With  $v$ ,  $x$  and  $[\ell''..r'']$  we can leverage Lemma 3.7 as shown in the proof of Lemma 3.15 to compute the number  $y$  of entries  $j \in [1..12] - \{12\}$  with  $L_{\mathcal{T}}[j] \geq 1$  such that  $C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[j]) \preceq_{\omega} C_{\mathcal{T}}(\text{LF}_{\mathcal{T}}[12])$ , i.e.,  $\text{LF}_{\mathcal{T}}[x]$  is the  $(y + 1)$ -th smallest element of  $\text{CR}_{\mathcal{T}}(\mathbf{A})$  by choice of  $x$ . From Line 14 through Line 15 of Algorithm 2 we obtain  $y = \text{rnkcnt}_{L_{\mathcal{T}}}(\max\{1, 4\}, 12 - 1, 1, 1) + \text{rnkcnt}_{L_{\mathcal{T}}}(1, 4 - 1, 1, 3) = 2 + 2 = 4$ . Then  $r' = \text{LF}_{\mathcal{T}}[12] + 9 - 4 - 1 = 8 + 4 = 12$  and  $\text{CR}_{\mathcal{T}}(\mathbf{A}) = [\ell'..r'] = [12 - 9 + 1..12] = [4..12]$ , which one can confirm by looking at the range of integer values of  $F_{\mathcal{T}}$  in Figure 7.

Next, we compute  $\text{CR}_{\mathcal{T}}(P[3..]) = \text{CR}_{\mathcal{T}}(\mathbf{CA})$ . We set  $[\ell..r] = \text{CR}_{\mathcal{T}}(\mathbf{A}) = [4..12]$ , and compute  $h = \pi(\mathbf{CA}) = 2$  and  $c = |\mathbf{A}|_p = 1$ . Consequently,  $\mathbf{C}$  did not appear in  $P[4..] = \mathbf{A}$  and we have to proceed as in the previous update of the conjugate range. We get  $\delta = \text{rnkcnt}_{L_{\mathcal{T}}}(4, 12, 2, 3) = 4$ ,  $v = \text{RNV}_{L_{\mathcal{T}}}(4, 12, 1) = 2$ ,  $x = \text{select}_2(L_{\mathcal{T}}, \text{rank}_2(L_{\mathcal{T}}, 12)) = 7$ ,  $[\ell''..r''] = \text{MI}_{\text{LCP}_{\mathcal{T}}^{\infty}}(7, 2) = [6..8]$  and  $y = \text{rnkcnt}_{L_{\mathcal{T}}}(\max\{4, 6\}, 7 - 1, 2, 2) + \text{rnkcnt}_{L_{\mathcal{T}}}(4, 6 - 1, 2, 3) = 1 + 1 = 2$ . Then  $r' = \text{LF}_{\mathcal{T}}[7] + 4 - 2 - 1 = 11 + 1 = 12$  and  $\text{CR}_{\mathcal{T}}(\mathbf{CA}) = [\ell'..r'] = [12 - 4 + 1..12] = [9..12]$ , which is again readily confirmed by, e.g., Figure 5 and  $\langle \mathbf{CA} \rangle = \infty\infty$  with Lemma 2.5.

Next, we compute  $\text{CR}_{\mathcal{T}}(P[2..]) = \text{CR}_{\mathcal{T}}(\mathbf{CCA})$ . We set  $[\ell..r] = \text{CR}_{\mathcal{T}}(\mathbf{CA}) = [9..12]$ , and compute  $h = \pi(\mathbf{CCA}) = 1$  and  $c = |\mathbf{CA}|_p = 2$ . This time, we have  $h \leq c$ , i.e.,  $\mathbf{C}$  already appeared in  $\mathbf{CA}$  and we enter Line 2 of Algorithm 2. Then the proof of Lemma 3.15 tells us that the conjugates corresponding to entries  $i \in [\ell..r] = [9..12]$  with  $L_{\mathcal{T}}[i] = h = 1$  are the left rotations of conjugates whose prefix of length 3 p-matches  $\mathbf{CCA}$ . Then

$\delta = \text{rnkcnt}_{L_{\mathcal{T}}}(9, 12, 1, 1) = 3$  is the length of  $\text{CR}_{\mathcal{T}}(\text{CCA})$ . Moreover, by the non-crossing property coined in Lemma 3.8, the LF-mapping of the largest element  $x$  in  $\text{CR}_{\mathcal{T}}(\text{CA})$  such that  $L_{\mathcal{T}}[x] = h$  is already the right boundary of  $\text{CR}_{\mathcal{T}}(\text{CCA})$ , i.e.,  $r' = \text{LF}_{\mathcal{T}}[\text{select}_1(L_{\mathcal{T}}, \text{rank}_1(L_{\mathcal{T}}, 12))] = \text{LF}_{\mathcal{T}}[12] = 8$ . Then  $\text{CR}_{\mathcal{T}}(\text{CCA}) = [8 - 3 + 1..8] = [6..8]$ .

Last, we compute  $\text{CR}_{\mathcal{T}}(P[1..]) = \text{CR}_{\mathcal{T}}(\text{bCCA})$ . We set  $[\ell..r] = \text{CR}_{\mathcal{T}}(\text{CCA}) = [6..8]$ , and compute  $h = \pi(\text{bCCA}) = b$  and  $c = |\text{CCA}|_p = 2$ . Here,  $h < c$  since we stipulated that s-symbols are smaller than any integer. We proceed similarly to the previous case leveraging the non-crossing property of entries of  $L_{\mathcal{T}}[\ell..r]$  evaluating to  $\text{b}$ . As  $\delta = \text{rnkcnt}_{L_{\mathcal{T}}}(6, 8, b, b) = 1$  indicates that  $\text{CR}_{\mathcal{T}}(\text{bCCA})$  is non-empty, we may compute the right boundary  $r' = \text{LF}_{\mathcal{T}}[\text{select}_b(L_{\mathcal{T}}, \text{rank}_b(L_{\mathcal{T}}, 8))] = \text{LF}_{\mathcal{T}}[8] = 3$  of  $\text{CR}_{\mathcal{T}}(\text{bCCA})$ . Then  $\text{CR}_{\mathcal{T}}(\text{bCCA}) = [3 - 1 + 1..3] = [3..3]$ . Thus, the count query for  $P$  returns 1, namely  $P = \text{bCCA} \approx \text{BBC} = \text{Rot}(T_3, 2)^{\omega}[..4]$ .

## 4. Constructing the epBWT in Compact Space

Our final goal is an iterative construction of the epBWT index. To this end, let  $\mathcal{R} \subset \Sigma^+$  be a set of p-strings and  $S \in \Sigma^+$  a p-string we want to add to  $\mathcal{R}$ , where we assume that  $S$  is at least as long as the longest p-string in  $\mathcal{R}$ . First, we give techniques to extend the epBWT index of  $\mathcal{R}$  to that of  $\mathcal{S} = \mathcal{R} \cup \{S\}$ . We subsequently generalize this extension such that we can construct the epBWT index of  $\mathcal{T}$  by iteratively extending the epBWT index of  $\{T_1, \dots, T_{k-1}\}$  to that of  $\{T_1, \dots, T_k\}$  for each  $k \in [1..d]$  in ascending order, where  $F_{\emptyset} = L_{\emptyset} = \text{LCP}_{\emptyset}^{\infty} = \varepsilon$ . In this section, we assume  $|T_k| = n_k \leq n_{k+1} = |T_{k+1}|$  for each  $k \in [1..d-1]$  to obtain the claimed time bounds.

### 4.1. Extending an epBWT Index

*Idea.* Consider two decks of cards, each with a different back color: red (for  $\mathcal{R}$ ) and slate (for  $\{S\}$ ). The first deck has  $|\text{CA}_{\mathcal{R}}|$  cards with *red-colored* backs and whose card ranks are  $C_{\mathcal{R}}(1)$  through  $C_{\mathcal{R}}(|\text{CA}_{\mathcal{R}}|)$ . The second deck has  $|\text{CA}_{\{S\}}|$  cards with *slate-colored* backs and whose card ranks are  $C_{\{S\}}(1)$  through  $C_{\{S\}}(|\text{CA}_{\{S\}}|)$ . So the red cards and the slate cards correspond one-to-one to the conjugates considered in  $\mathcal{R}$  and  $\{S\}$ , respectively. We stack the cards of each deck such that the stack of red cards and the stack of slate cards represent the conjugate arrays of  $\mathcal{R}$  and  $\{S\}$ , respectively. We then label the cards with the corresponding values of  $F$ ,  $L$  and  $\text{LCP}^{\infty}$ , e.g.,

the topmost card of the stack of red cards is labeled with  $F_{\mathcal{R}}[1]$ ,  $L_{\mathcal{R}}[1]$  and  $LCP_{\mathcal{R}}^{\infty}[1]$ . Then the stack of red cards and the stack of slate cards represent the epBWT index of  $\mathcal{R}$  and  $\{S\}$ , respectively.

Now, in order to obtain a representation of the conjugate array of  $\mathcal{S}$ , we have to correctly place the previously labeled slate cards into the stack of red cards representing the epBWT index of  $\mathcal{R}$ . To obtain a representation of the epBWT index of  $\mathcal{S}$ , we also need to update the  $LCP^{\infty}$  label of any slate card that has a red card directly on *top* of it, and we need to update the  $LCP^{\infty}$  labels of red cards if the card is directly *underneath* a slate card.

*Approach.* Per assumption, we already have the epBWT index of  $\mathcal{R}$ . Thus, our tasks are fourfold.

- (i) For each conjugate  $S'$  of  $S$ , compute the number of conjugates considered in  $\mathcal{R}$  that are not larger according to  $\omega$ -preorder.
- (ii) For each conjugate  $S'$  of  $S$ , compute the  $lcp^{\infty}$ -value with the largest conjugate considered in  $\mathcal{R}$  that is not larger than  $S'$  according to  $\omega$ -preorder.
- (iii) For each conjugate  $S'$  of  $S$ , compute the  $lcp^{\infty}$ -value with the smallest conjugate considered in  $\mathcal{R}$  that is larger than  $S'$  according to  $\omega$ -preorder.
- (iv) Compute the epBWT index of  $\{S\}$ .

Then (i) and (iv) allow us to infer the conjugate array of  $\mathcal{S}$ . In addition, (iv) yields the entries of  $F_{\mathcal{S}}$  and  $L_{\mathcal{S}}$  corresponding to conjugates of  $S$  as well as  $LCP_{\mathcal{S}}^{\infty}[i]$  if both  $CA_{\mathcal{S}}[i]$  and  $CA_{\mathcal{S}}[i - 1]$  correspond to conjugates of  $S$ , where  $i \in [2..|CA_{\mathcal{S}}|]$ . Moreover, (ii) yields  $LCP_{\mathcal{S}}^{\infty}[i]$  if  $CA_{\mathcal{S}}[i]$  and  $CA_{\mathcal{S}}[i - 1]$  correspond to a conjugate considered in  $\{S\}$  and  $\mathcal{R}$ , respectively, for  $i \in [2..|CA_{\mathcal{S}}|]$ . Finally, (iii) yields  $LCP_{\mathcal{S}}^{\infty}[i]$  if  $CA_{\mathcal{S}}[i]$  and  $CA_{\mathcal{S}}[i - 1]$  correspond to a conjugate considered in  $\mathcal{R}$  and  $\{S\}$ , respectively, for  $i \in [2..|CA_{\mathcal{S}}|]$ . The remaining entries of  $F_{\mathcal{S}}$ ,  $L_{\mathcal{S}}$  and  $LCP_{\mathcal{S}}^{\infty}$  are inferred from the epBWT index of  $\mathcal{R}$ . Below, we will define the data structures used to store the helper values presented in (i), (ii) and (iii) and rigorously prove that they suffice to compute the epBWT index of  $\mathcal{S}$  if those of  $\mathcal{R}$  and  $\{S\}$  are given. Subsequently, in Section 4.2, we will delve into the computation of both the epBWT index of  $\{S\}$  and the helper values.

---

**Algorithm 3:** Extending  $F_{\mathcal{R}}$  and  $L_{\mathcal{R}}$  to  $F_{\mathcal{S}}$  and  $L_{\mathcal{S}}$ , respectively.  
Here,  $\mathcal{R} \subset \Sigma^+$ ,  $S \in \Sigma^+$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$  and  $\lambda = |S|$ .

---

```

1 Function extend1( $P_S$ ,  $F_{\mathcal{R}}$ ,  $L_{\mathcal{R}}$ ,  $F_{\{S\}}$ ,  $L_{\{S\}}$ ):
2   for  $i \leftarrow 1$  to  $\lambda$  do
3      $pos \leftarrow \text{select}_1(P_S, i)$ ;
4      $\text{insert}_{F_{\mathcal{R}}}(pos, F_{\{S\}}[i])$ ;
5      $\text{insert}_{L_{\mathcal{R}}}(pos, L_{\{S\}}[i])$ ;
6   return  $F_{\mathcal{R}}$ ,  $L_{\mathcal{R}}$ ;

```

---

*Helper Values.* Let  $\rho = |L_{\mathcal{R}}|$  and  $\lambda = |S| = \max\{|V| \mid V \in \mathcal{S}\}$  be, respectively, the length of the current epBWT and the longest string we want to index, which is  $S$ . Let  $CNT_{\mathcal{S}}$  denote a string of length  $\lambda$  over  $[1..\rho]$  such that

$$CNT_{\mathcal{S}}[i] = |\{j \in [1..\rho] \mid C_{\mathcal{R}}(j) \preceq_{\omega} C_{\{S\}}(\text{CA}_{\{S\}}[i])\}| \text{ for each } i \in [1..\lambda],$$

i.e.,  $CNT_{\mathcal{S}}$  stores the helper values described in (i). Note that  $CNT_{\mathcal{S}}$  needs  $O((\rho + \lambda) \lg(\rho + \lambda))$  bits of space. In view of a construction in compact space, we want to represent  $CNT_{\mathcal{S}}$  in  $O((\rho + \lambda) \lg \sigma)$  bits of space. To that end, let  $P_S$  denote the bit string of length  $\rho + \lambda$  satisfying  $\text{rank}_1(P_S, \rho + \lambda) = \lambda$  and  $CNT_{\mathcal{S}}[i] = \text{rank}_0(P_S, \text{select}_1(P_S, i))$  for each  $i \in [1..\lambda]$ . Then  $P_S$  takes  $O(\rho + \lambda)$  bits of space and allows us access to  $CNT_{\mathcal{S}}[i]$  in  $O(\frac{\lg(\rho + \lambda)}{\lg \lg(\rho + \lambda)})$  time for each  $i \in [1..\lambda + \rho]$  if represented by the data structure of Lemma 2.3. For our running example, we have  $P_{\{T_1\}} = 111$ ,  $P_{\{T_1, T_2\}} = 0001111$  and  $P_{\mathcal{T}} = 001010011001$ , where  $S$  is defined as  $T_i$  for  $\mathcal{R} = \{T_1, \dots, T_{i-1}\}$ . We can construct  $P_S$  by zeroing  $P_{\mathcal{R}}$  and calling  $\text{insert}_{P_S}(\text{select}_0(P_{\mathcal{R}}, CNT_{\mathcal{S}}[i]) + 1, 1)$  for each  $i \in [1..\lambda]$ . Thus, we only need to compute each entry of  $CNT_{\mathcal{S}}$  once and do not need to store the entry explicitly. The following result is immediate.

**Lemma 4.1.** *Let  $\mathcal{R} \subset \Sigma^+$ ,  $S \in \Sigma^+$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$  and  $\lambda = |S|$ . Given  $P_S$ ,  $F_{\mathcal{R}}$ ,  $L_{\mathcal{R}}$ ,  $F_{\{S\}}$  and  $L_{\{S\}}$ , Algorithm 3 correctly computes  $F_{\mathcal{S}}$  and  $L_{\mathcal{S}}$ .*

Consequently, the extension of both  $F_{\mathcal{R}}$  and  $L_{\mathcal{R}}$  boils down to the computation of  $P_S$  and the epBWT index of  $\{S\}$ . We tackle the extension of  $LCP_{\mathcal{R}}^{\infty}$  to  $LCP_{\mathcal{S}}^{\infty}$  similarly. Let  $\text{TOP}_{\mathcal{S}}$  and  $\text{UDN}_{\mathcal{S}}$  denote strings of length  $\lambda$

$i$	$F_S$	$L_S$	$LCP_S^\infty$	$P_S$
1	a	1	0	0
2	b	a	0	0
3	b	2	1	1
4	b	3	2	0
5	1	b	0	0
6	2	1	1	1
7	3	3	2	0
8	1	2	1	1
9	1	b	1	1
10	1	2	2	0
11	1	2	2	0
12	1	b	2	0
13	3	3	1	0
14	2	1	2	1
15	2	1	2	0
16	2	1	2	0
17	3	1	2	0

(a) The epBWT index of  $\mathcal{S} = \mathcal{T} \cup \{S\}$ .

$i$	$F_{\{S\}}[i]$	$L_{\{S\}}[i]$	$LCP_{\{S\}}^\infty[i]$
1	b	2	0
2	2	1	0
3	1	2	1
4	1	b	1
5	2	1	1

(b) Auxiliary strings for the extension.

Figure 8: The extension of the epBWT index of  $\mathcal{R} = \mathcal{T}$  by  $S = \text{BBCCb}$ .

over  $[-1..\sigma_p]$  such that

$$\text{TOP}_{\mathcal{S}}[i] = \begin{cases} \text{lcp}^\infty(C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[j]), C_{\{S\}}(\text{CA}_{\{S\}}[i])) & \text{if } j \geq 1, \\ -1 & \text{otherwise, and} \end{cases}$$

$$\text{UDN}_{\mathcal{S}}[i] = \begin{cases} \text{lcp}^\infty(C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[j+1]), C_{\{S\}}(\text{CA}_{\{S\}}[i])) & \text{if } j+1 \leq \rho, \\ -1 & \text{otherwise,} \end{cases}$$

for each  $i \in [1..\lambda]$  with  $j = \text{CNT}_{\mathcal{S}}[i]$ , i.e.,  $\text{TOP}_{\mathcal{S}}$  and  $\text{UDN}_{\mathcal{S}}$  store the helper values presented in (ii) and (iii), respectively. For an example, see Figure 8.

**Lemma 4.2.** *Let  $\mathcal{R} \subset \Sigma^+$ ,  $S \in \Sigma^+$ , and  $\mathcal{S} = \mathcal{R} \cup \{S\}$ . Then Algorithm 4 correctly computes the epBWT index of  $\mathcal{S}$ .*

*Proof.* Let  $\lambda = |S|$  and  $\rho = |\text{L}_{\mathcal{R}}|$ . The input of Algorithm 4 consists of  $P_{\mathcal{S}}$ ,  $\text{TOP}_{\mathcal{S}}$ ,  $\text{UDN}_{\mathcal{S}}$  and the epBWT indexes of  $\mathcal{R}$  and  $\{S\}$ . The statement

---

**Algorithm 4:** Extending the epBWT index of  $\mathcal{R}$  to that of  $\mathcal{S}$ .

Here,  $\mathcal{R} \subset \Sigma^+$ ,  $S \in \Sigma^+$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$  and  $\lambda = |S|$ .

---

```

1 Function extend2( $P_{\mathcal{S}}$ ,  $\text{TOP}_{\mathcal{S}}$ ,  $\text{UDN}_{\mathcal{S}}$ ,  $F_{\mathcal{R}}$ ,  $L_{\mathcal{R}}$ ,  $\text{LCP}_{\mathcal{R}}^{\infty}$ ,  $F_{\{S\}}$ ,  $L_{\{S\}}$ ,  

    $\text{LCP}_{\{S\}}^{\infty}$ ):  

2    $F_{\mathcal{R}}, L_{\mathcal{R}} \leftarrow \text{extend1}(P_{\mathcal{S}}, F_{\mathcal{R}}, L_{\mathcal{R}}, F_{\{S\}}, L_{\{S\}})$ ; // Alg. 3  

3   for  $i \leftarrow 1$  to  $\lambda$  do  

4      $pos \leftarrow \text{select}_1(P_{\mathcal{S}}, i)$ ;  

5     if  $\text{TOP}_{\mathcal{S}}[i] \geq 0$  and  $P_{\mathcal{S}}[pos - 1] = 0$  then  

6        $\text{insert}_{\text{LCP}_{\mathcal{R}}^{\infty}}(pos, \text{TOP}_{\mathcal{S}}[i])$ ;  

7     else  $\text{insert}_{\text{LCP}_{\mathcal{R}}^{\infty}}(pos, \text{LCP}_{\{S\}}^{\infty}[i])$ ;  

8     if  $\text{UDN}_{\mathcal{S}}[i] \geq 0$  and  $P_{\mathcal{S}}[pos + 1] = 0$  then  

9        $\text{LCP}_{\mathcal{R}}^{\infty}[pos + 1] \leftarrow \text{UDN}_{\mathcal{S}}[i]$ ;  

10  return  $F_{\mathcal{R}}, L_{\mathcal{R}}, \text{LCP}_{\mathcal{R}}^{\infty}$ ;
```

---

regarding  $F_{\mathcal{S}}$  and  $L_{\mathcal{S}}$  follows from Lemma 4.1. By definition, initially

$$\begin{aligned} \text{LCP}_{\mathcal{R}}^{\infty}[..\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, 1)) - 1] &= \text{LCP}_{\mathcal{S}}^{\infty}[..\text{select}_1(P_{\mathcal{S}}, 1) - 1] \text{ and} \\ \text{LCP}_{\mathcal{R}}^{\infty}[\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, \lambda)) + 2..] &= \text{LCP}_{\mathcal{S}}^{\infty}[\text{select}_1(P_{\mathcal{S}}, \lambda) + 2..]. \end{aligned}$$

Moreover, for each  $i \in [1..\lambda - 1]$ , initially

$$\begin{aligned} \text{LCP}_{\mathcal{R}}^{\infty}[\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, i)) + 2..\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, i + 1)) - 1] \\ = \text{LCP}_{\mathcal{S}}^{\infty}[\text{select}_1(P_{\mathcal{S}}, i) + 2..\text{select}_1(P_{\mathcal{S}}, i + 1) - 1]. \end{aligned}$$

Now, assume  $\text{LCP}_{\mathcal{R}}^{\infty}$  was updated such that

$$\begin{aligned} \text{LCP}_{\mathcal{R}}^{\infty}[..\text{select}_1(P_{\mathcal{S}}, i) - 1] &= \text{LCP}_{\mathcal{S}}^{\infty}[..\text{select}_1(P_{\mathcal{S}}, i) - 1], \\ \text{LCP}_{\mathcal{R}}^{\infty}[\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, \lambda)) + 1 + i..] &= \text{LCP}_{\mathcal{S}}^{\infty}[\text{select}_1(P_{\mathcal{S}}, \lambda) + 2..], \end{aligned}$$

and  $\text{LCP}_{\mathcal{R}}^{\infty}[\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, j)) + 1 + i..\text{rank}_0(P_{\mathcal{S}}, \text{select}_1(P_{\mathcal{S}}, j + 1)) - 2 + i] = \text{LCP}_{\mathcal{S}}^{\infty}[\text{select}_1(P_{\mathcal{S}}, i) + 2..\text{select}_1(P_{\mathcal{S}}, i + 1) - 1]$ , for some fixed  $i \in [1..\lambda]$  and each  $j \in [i..\lambda - 1]$ .

- If  $\text{TOP}_{\mathcal{S}}[i] \geq 0$  and  $P_{\mathcal{S}}[\text{select}_1(P_{\mathcal{S}}, i) - 1] = 0$ , then  $\text{CA}_{\mathcal{S}}^{-1}[\text{select}_1(P_{\mathcal{S}}, i) - 1] \in [1..\rho]$  and thus  $\text{LCP}_{\mathcal{S}}^{\infty}[\text{select}_1(P_{\mathcal{S}}, i)] = \text{TOP}_{\mathcal{S}}[i]$ .
- If  $\text{TOP}_{\mathcal{S}}[i] \geq 0$  and  $P_{\mathcal{S}}[\text{select}_1(P_{\mathcal{S}}, i) - 1] = 1$ , then  $\text{CA}_{\mathcal{S}}^{-1}[\text{select}_1(P_{\mathcal{S}}, i) - 1] \in [\rho + 1..\rho + \lambda]$ , i.e.,  $\text{LCP}_{\mathcal{S}}^{\infty}[\text{select}_1(P_{\mathcal{S}}, i)] = \text{LCP}_{\{S\}}^{\infty}[i]$ .

- If  $\text{TOP}_S[i] < 0$  and  $i = 1$ , then  $\text{LCP}_S^\infty[\text{select}_1(P_S, i)] = 0 = \text{LCP}_{\{S\}}^\infty[i]$ .
- If  $\text{TOP}_S[i] < 0$  and  $i > 1$ , then  $\text{LCP}_S^\infty[\text{select}_1(P_S, i)] = \text{LCP}_{\{S\}}^\infty[i]$ .

These cases are treated from Line 5 through Line 7 and we updated  $\text{LCP}_R^\infty$  such that  $\text{LCP}_R^\infty[..\text{select}_1(P_S, i)] = \text{LCP}_S^\infty[..\text{select}_1(P_S, i)]$ .

If  $\text{select}_1(P_S, i) < \rho + \lambda$  and  $P_S[\text{select}_1(P_S, i) + 1] = 1$ , or  $\text{UDN}_S[i] < 0$ , then by assumption we already have

$$\begin{aligned}\text{LCP}_R^\infty[..\text{select}_1(P_S, i + 1) - 1] &= \text{LCP}_R^\infty[..\text{select}_1(P_S, i)] \\ &= \text{LCP}_S^\infty[..\text{select}_1(P_S, i)] \\ &= \text{LCP}_S^\infty[..\text{select}_1(P_S, i + 1) - 1]\end{aligned}$$

for  $i + 1 \leq \lambda$  and  $\text{LCP}_R^\infty = \text{LCP}_S^\infty$  for  $i = \lambda$ . Otherwise,  $\text{CA}_S^{-1}[\text{select}_1(P_S, i) + 1] \in [1..\rho]$  and consequently  $\text{LCP}_S^\infty[\text{select}_1(P_S, i)] = \text{UDN}_S[i]$ . By assumption,  $\text{LCP}_R^\infty[..\text{select}_1(P_S, i+1)-1] = \text{LCP}_S^\infty[..\text{select}_1(P_S, i+1)-1]$  for  $i+1 \leq \lambda$  and  $\text{LCP}_R^\infty = \text{LCP}_S^\infty$  for  $i = \lambda$ . We treat this part from Line 8 through Line 9. We readily confirm that after the loop our assumption holds for  $i + 1$  if  $i + 1 \leq \lambda$ . Since  $i$  was chosen arbitrarily, the statement follows.  $\square$

#### 4.2. Building the Remaining Parts

We conclude that the extension of the epBWT of  $R$  to that of  $S$  can be reduced to computing  $\text{TOP}_S$ ,  $\text{UDN}_S$ ,  $P_S$  and the epBWT index of  $\{S\}$ . For their computation, we assume  $R, S \subset (\Sigma - \{\$\})^+$  in this subsection. This is without loss of generality because we can always add a new s-symbol to  $\Sigma$  taking over the role of  $\$$ , which is stipulated to be the smallest s-symbol.

Let  $z = \max\{|V| \mid V \in \mathcal{S}\}$ ,  $S' = S^\omega[..4z]\$$  and denote by  $Y_S$  a bit string of length  $4z + 1$  such that  $Y_S[i] = 1$  if and only if  $\text{CA}_{\{S'\}}[i] \in [2..\lambda + 1]$  for each  $i \in [1..4z + 1]$ . Here,  $Y_S$  induces a one-to-one mapping of the conjugates of  $S$  to the conjugates of  $S'$  such that  $\text{Rot}(S, i)^\omega[..3z] = S'[i + 1..3z + i + 1]$  for each  $i \in [1..\lambda]$ , which allows us to draw conclusions about the  $\omega$ -preorder and the  $[\dots]$ -encoding of the conjugates of  $S$  from the  $\omega$ -preorder and the  $[\dots]$ -encoding of the conjugates of  $S'$  by application of Lemma 3.4 and Corollary 3.5.

The convenience of working with  $S'$  instead of  $S$  will become clear shortly. For what follows, recall that we assumed  $S$  to be at least as long as the longest string in  $R$ , i.e.,  $\lambda = z$  and  $S' = S^\omega[..4\lambda]\$ = S^4\$$ .

**Lemma 4.3.** *Algorithm 5 correctly computes the epBWT index of  $\{S\}$  for  $S \in (\Sigma - \{\$\})^+$ .*

---

**Algorithm 5:** Computing the epBWT of the singleton set  $\{S\}$ .

Here,  $S \in (\Sigma - \{\$\})^+$ ,  $S' = S^4\$$  and  $\lambda = |S|$ .

---

```

1 Function singlonepBWT( $Y_S$ ,  $F_{\{S'\}}$ ,  $L_{\{S'\}}$ ,  $LCP_{\{S'\}}^\infty$ ):
2    $F_{\{S\}}, L_{\{S\}} \leftarrow$  zeroed string of length  $\lambda$  over  $\Sigma_s \cup [1..|\sigma_p|]$ ;
3    $LCP_{\{S\}}^\infty \leftarrow$  zeroed string of length  $\lambda$  over  $[0..|\sigma_p|]$ ;
4    $prev \leftarrow 0$ ;
5   for  $i \leftarrow 1$  to  $\lambda$  do
6      $curr \leftarrow \text{select}_1(Y_S, i)$ ;
7      $F_{\{S\}}[i] \leftarrow F_{\{S'\}}[curr]$ ;
8      $L_{\{S\}}[i] \leftarrow L_{\{S'\}}[curr]$ ;
9      $LCP_{\{S\}}^\infty[i] \leftarrow \text{RNV}_{LCP_{\{S'\}}^\infty}(prev + 1, curr, -1)$ ;
10     $prev \leftarrow curr$ ;
11  return  $F_{\{S\}}, L_{\{S\}}, LCP_{\{S\}}^\infty$ ;

```

---

*Proof.* Let  $S' = S^4\$$ , and  $\lambda = |S|$ . Algorithm 5 takes as its input  $Y_S$  and the epBWT index of  $\{S'\}$ . We have

$$\langle C_{\{S\}}(\text{CA}_{\{S\}}[i])^3 \rangle = \langle C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{select}_1(Y_S, i)])[..3\lambda] \rangle$$

for each  $i \in [1..\lambda]$  by construction. Thus,  $\text{select}_1(Y_S, i) \in \text{CR}_{\{S'\}}(C_{\{S\}}(i)^3)$  for each  $i \in [1..\lambda]$ , which in turn implies

$$\begin{aligned} \llbracket C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{select}_1(Y_S, i)])[..3\lambda] \rrbracket &= \llbracket C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{select}_1(Y_S, i)]) \rrbracket[..3\lambda] \\ &= \llbracket C_{\{S\}}(\text{CA}_{\{S\}}[i]) \rrbracket \end{aligned}$$

by choice of  $S'$  and the backward search. Also, by choice of  $S'$ ,

$$\llbracket C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{select}_1(Y_S, i)]) \rrbracket[\lambda] = \llbracket C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{select}_1(Y_S, i)]) \rrbracket[4\lambda + 1]$$

for  $i \in [1..\lambda]$ . Then the statement follows from Lemma 3.4 and Lemma 3.14.  $\square$

By choice of  $S'$ , the epBWT index of  $\{S'\}$  is the pBWT index of  $S'$ , and we can leverage Lemma 2.8 for its construction. Moreover, during construction of the pBWT index of  $S'$  the string  $Y_S$  can be built similarly to the bit string indicating sample positions of the suffix array in the proof of Lemma 2.8. The epBWT index of  $\{S'\}$  and  $Y_S$  now come in handy if we consider the set of strings  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$ .

**Lemma 4.4.** Let  $\mathcal{R} \subset (\Sigma - \{\$\})^+$  and  $S \in (\Sigma - \{\$\})^+$ . Define  $\mathcal{S} = \mathcal{R} \cup \{S\}$ ,  $S' = S^4\$$ ,  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$ ,  $\lambda = |S| = \max\{|V| \mid V \in \mathcal{S}\}$  and  $\rho = |\mathcal{L}_{\mathcal{R}}|$ . Then  $\text{CNT}_{\mathcal{S}}[i] = \text{CNT}_{\mathcal{S}'}[j] + |\text{CR}_{\mathcal{R}}(S^3)|$ ,

$$\begin{aligned} \text{TOP}_{\mathcal{S}}[i] &= \begin{cases} \text{TOP}_{\mathcal{S}'}[j] & \text{if } |\text{CR}_{\mathcal{R}}(S^3)| = 0, \\ \text{UDN}_{\mathcal{S}'}[j] & \text{otherwise, and} \end{cases} \\ \text{UDN}_{\mathcal{S}}[i] &= \begin{cases} \text{UDN}_{\mathcal{S}'}[j] & \text{if } |\text{CR}_{\mathcal{R}}(S^3)| = 0, \\ \text{LCP}_{\mathcal{R}}^\infty[\text{CNT}_{\mathcal{S}}[i] + 1] & \text{if } \text{CNT}_{\mathcal{S}}[i] < \rho, \\ -1 & \text{otherwise,} \end{cases} \end{aligned}$$

for each  $i \in [1..|\lambda|]$  with  $j = \text{select}_1(\text{Y}_S, i)$ .

*Proof.* By choice of  $S'$ ,  $C_{\mathcal{R}}(x) \neq_{\omega} C_{\{S'\}}(y)$  for each  $x \in [1..\rho]$  and  $y \in [1..4\lambda + 1]$ . In particular,  $C_{\{S'\}}(y) \prec_{\omega} C_{\mathcal{R}}(x)$  for each  $x \in [1..\rho]$  and  $y \in [1..4\lambda + 1]$  with  $\langle C_{\{S'\}}(y)[..3\lambda] \rangle = \langle C_{\mathcal{R}}(x)^\omega[..3\lambda] \rangle$ . Since

$$\langle C_{\{S\}}(\text{CA}_{\{S\}}[i])^3 \rangle = \langle C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{select}_1(\text{Y}_S, i)])[..3\lambda] \rangle \text{ for each } i \in [1..|\lambda|]$$

by construction, the statement regarding  $\text{CNT}_{\mathcal{S}}[i]$  then follows from the definition of a conjugate range, Lemma 3.4 and Corollary 3.5.

Fix some  $i \in [1..|\lambda|]$  with  $j = \text{select}_1(\text{Y}_S, i)$ . Then the statement regarding  $\text{TOP}_{\mathcal{S}}[i]$  and  $\text{UDN}_{\mathcal{S}}[i]$  is immediate from their definition if  $|\text{CR}_{\mathcal{R}}(S^3)| = 0$ . Assume  $|\text{CR}_{\mathcal{R}}(S^3)| \neq 0$ . Then

$$\text{UDN}_{\mathcal{S}'}[j] = \text{lcp}^\infty(C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{\mathcal{S}'}[i] + k]), C_{\{S'\}}(\text{CA}_{\{S'\}}[i]))$$

for each  $k \in [1..|\text{CR}_{\mathcal{R}}(S^3)|]$  by the definition of a conjugate range and the proof of Lemma 3.14. In particular,  $\text{UDN}_{\mathcal{S}'}[j] = \text{TOP}_{\mathcal{S}}[i]$ . For the statement regarding  $\text{UDN}_{\mathcal{S}}[i]$ , we simply apply its definition to the fact that  $\langle C_{\{S\}}(\text{CA}_{\{S\}}[i])^3 \rangle = \langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{\mathcal{S}}[i]])^\omega[..3\lambda] \rangle$ .  $\square$

We give an example in Figure 9. It remains to show how to compute  $\text{TOP}_{\mathcal{S}'}$ ,  $\text{UDN}_{\mathcal{S}'}$  and the individual values of  $\text{CNT}_{\mathcal{S}'}$  for our extension. By choice of  $S'$ ,  $\text{CNT}_{\mathcal{S}'}[1] = \text{CNT}_{\mathcal{S}'}[\text{CA}_{\{S'\}}^{-1}[4\lambda + 1]] = 0$ ,  $\text{TOP}_{\mathcal{S}'}[1] = -1$ ,  $\text{UDN}_{\mathcal{S}'}[1] = 0$  and  $\llbracket S' \rrbracket$  is primitive, with the latter implying that  $\text{LF}_{\{S'\}}$  decomposes into a single cycle. In what follows, we will show how to compute  $\text{CNT}_{\mathcal{S}'}[j]$ ,  $\text{TOP}_{\mathcal{S}'}[j]$  and  $\text{UDN}_{\mathcal{S}'}[j]$  from  $\text{CNT}_{\mathcal{S}'}[i]$ ,  $\text{TOP}_{\mathcal{S}'}[i]$  and  $\text{UDN}_{\mathcal{S}'}[i]$  for each  $i \in [1..4\lambda + 1]$  with  $j = \text{LF}_{\{S'\}}[i]$ . Then application of this result  $4\lambda$

$i$	$\text{CA}_{\{S'\}}[i]$	$F_{\{S'\}}[i]$	$L_{\{S'\}}[i]$	$LCP_{\{S'\}}^{\infty}[i]$	$Y_S[i]$	$CNT_{S'}[i]$	$TOP_{S'}[i]$	$UDN_{S'}[i]$
1	21	\$	b	0	0	0	-1	0
2	20	b	2	0	0	1	0	0
3	15	b	2	0	0	2	1	2
4	10	b	2	2	0	2	1	2
5	5	b	2	2	1	2	1	2
6	19	2	1	0	0	4	1	1
7	14	2	1	1	0	4	1	2
8	9	2	1	2	0	4	1	2
9	4	2	1	2	1	4	1	2
10	18	1	2	1	0	5	1	1
11	13	1	2	1	0	5	1	1
12	8	1	2	2	0	5	1	1
13	3	1	2	2	1	5	1	1
14	16	1	b	1	0	5	1	2
15	11	1	b	2	0	5	1	2
16	6	1	b	2	1	5	1	2
17	1	1	\$	2	0	5	1	2
18	17	2	1	1	0	9	2	2
19	12	2	1	2	0	9	2	2
20	7	2	1	2	0	9	2	2
21	2	2	1	2	1	9	2	2

Figure 9: Auxiliary data structures for extending the epBWT index of our running example  $\mathcal{T} = \mathcal{R}$  by  $S = \text{BBCCb}$ . Here,  $\mathcal{S} = \mathcal{R} \cup \{S\}$ ,  $S' = S^4\$$ ,  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$ , and  $|\text{CR}_{\mathcal{R}}(S^3)| = 0$ . By Lemma 4.3 and Lemma 4.4,  $F_{\{S\}}$ ,  $L_{\{S\}}$ ,  $CNT_{\mathcal{S}}$ ,  $TOP_{\mathcal{S}}$  and  $UDN_{\mathcal{S}}$  is the subsequence of entries  $i$  in  $F_{\{S'\}}$ ,  $L_{\{S'\}}$ ,  $CNT_{\mathcal{S}'}$ ,  $TOP_{\mathcal{S}'}$  and  $UDN_{\mathcal{S}'}$ , respectively, such that  $Y_S[i] = 1$ . Moreover, we use RNV queries on  $LCP_{\{S'\}}^{\infty}$  as presented in Algorithm 5 to obtain  $LCP_{\{S\}}^{\infty} = 00111$ . The extracted data structures are presented in Figure 8.

times starting from  $i = 1$  will yield all values of  $\text{CNT}_{\mathcal{S}'}$ ,  $\text{TOP}_{\mathcal{S}'}$  and  $\text{UDN}_{\mathcal{S}'}$  since  $\text{LF}_{\{\mathcal{S}'\}}$  decomposes into a single cycle.

The following result is an adaptation of the techniques presented by Iseri et al. [15, Section 3.1].

**Lemma 4.5.** *Let  $\mathcal{R} \subset (\Sigma - \{\$\})^+$ ,  $S \in (\Sigma - \{\$\})^+$ ,  $S' = S^4\$$ ,  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$  and  $\lambda = |S| \geq \max\{|V| \mid V \in \mathcal{R}\}$ . Algorithm 6 correctly computes  $\text{CNT}_{\mathcal{S}'}[\text{LF}_{\{\mathcal{S}'\}}[i]]$  for each given  $i \in [1..4\lambda + 1]$ .*

*Proof.* Let  $\rho = |\mathcal{L}_{\mathcal{R}}|$ . The input of Algorithm 6 for some  $i \in [1..4\lambda + 1]$  consists of  $\pi(C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{LF}_{\{S'\}}[i]]))$ ,  $\text{CNT}_{\mathcal{S}'}[i]$ ,  $\text{TOP}_{\mathcal{S}'}[i]$ ,  $\text{UDN}_{\mathcal{S}'}[i]$  and the epBWT index of  $\mathcal{R}$ .

Fix some  $i \in [1..4\lambda + 1]$  and let  $j = \text{LF}_{\{S'\}}[i]$  and  $h = \text{F}_{\{S'\}}[j] = \pi(C_{\{S'\}}(\text{CA}_{\{S'\}}[j]))$ . We assume  $\mathcal{R} \neq \emptyset$  throughout the proof since the statement is immediate if  $\mathcal{R} = \emptyset$ .

**Case 1.** Assume  $h \leq 0$ , i.e.,  $h \in \Sigma_s$ . Then Lemma 3.7 implies  $\text{CNT}_{\mathcal{S}'}[j] = \text{rnkcnt}_{\mathcal{L}_{\mathcal{R}}}(1, \rho, \$, h) - \text{rnkcnt}_{\mathcal{L}_{\mathcal{R}}}(i+1, \rho, h, h)$ . The case is treated from Line 3 through Line 4.

**Case 2.** Assume  $h \geq 1$  and  $\text{rnkcnt}_{\mathcal{L}_{\mathcal{R}}}(1, \rho, 1, \sigma_p) = 0$ , i.e.,  $C_{\{S'\}}(\text{CA}_{\{S'\}}[j])[1]$  is a p-symbol and  $C_{\mathcal{R}}(\alpha) \in \Sigma_s^+$  for each  $\alpha \in [1..\rho]$ . Then  $\text{CNT}_{\mathcal{S}'}[j] = \rho$  since  $\alpha$  is larger than any s-symbol. The case is treated in Line 5.

**Case 3.** Assume  $h \geq 1$  and  $\text{rnkcnt}_{\mathcal{L}_{\mathcal{R}}}(1, \rho, 1, \sigma_p) \geq 1$ , i.e.,  $C_{\{S'\}}(\text{CA}_{\{S'\}}[j])[1]$  is a p-symbol and  $C_{\mathcal{R}}(\alpha) \notin \Sigma_s^+$  for some  $\alpha \in [1..\rho]$ . The case is treated from Line 6. We will follow the proof of Lemma 12 by Iseri et al. [15] in less detail and adapt it where necessary.

For notational convenience, let  $V_\alpha = C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{LF}_{\mathcal{R}}[\alpha]])$  for each  $\alpha \in [1..\rho]$  and  $U = C_{\{S'\}}(\text{CA}_{\{S'\}}[j])$ . We define

$$\gamma = \max\{\text{lcp}(\langle V_\alpha^\omega[..4\lambda + 1] \rangle, \langle U \rangle) \mid \alpha \in [1..\rho]\}, \text{ and}$$

$$t_e = \begin{cases} \text{select}_x(\langle \text{Rot}(U, 1) \rangle, e) & \text{if } e \in [1.. \min\{h, |U|_p\}], \text{ or} \\ 4\lambda + 1 & \text{if } e \in [\min\{h, |U|_p\} + 1.. \sigma_p + 1]. \end{cases}$$

In other words, the left rotation of  $V_\alpha$ ,  $\text{Rot}(V_\alpha, 1)$ , is the  $\alpha$ -smallest conjugate in  $\mathcal{R}$ ,  $U$  and  $\text{Rot}(U, 1)$  are the  $j$ -smallest and the  $i$ -smallest conjugate in  $\{S'\}$ , respectively, and  $\gamma$  is the maximum length of the prefixes of the pre-encodings of  $U$  with all  $V_\alpha$ 's. We will not compute  $\gamma$  directly, but  $e \in [1.. \min\{h, |U|_p\}]$  such that  $\gamma \in [t_e..t_{e+1}]$ , and then infer  $\text{CNT}_{\mathcal{S}'}[j]$ . The following is immediate by definitions.

---

**Algorithm 6:** Computing  $\text{CNT}_{S'}[j]$  from  $\text{CNT}_{S'}[i]$ ,  $\text{TOP}_{S'}[i]$  and  $\text{UDN}_{S'}[i]$  for  $i \in [1..4\lambda + 1]$  with  $j = \text{LF}_{\{S'\}}[i]$ . Here,  $\mathcal{R}, \mathcal{S} \subset (\Sigma - \{\$\})^+$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$ ,  $S' = S^4\$$ ,  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$ ,  $\rho = |\mathcal{L}_\mathcal{R}|$ ,  $\lambda = |S| = \max\{|V| \mid V \in \mathcal{S}\}$  and  $h = \pi(C_{\{S'\}}(\text{CA}_{\{S'\}}[j])) = F_{\{S'\}}[j]$ .

---

```

1 Function nextcnt( $h$ ,  $\text{CNT}_{S'}[i]$ ,  $\text{TOP}_{S'}[i]$ ,  $\text{UDN}_{S'}[i]$ ,  $F_\mathcal{R}$ ,  $L_\mathcal{R}$ ,  

2    $\text{LCP}_\mathcal{R}^\infty$ ):  

3     if  $\rho = 0$  then return 0 ; //  $\mathcal{R} = \emptyset$   

4     if  $h < 1$  then // Case 1 of Lemma 4.5  

5       return rnkcnt $_{\mathcal{L}_\mathcal{R}}$ (1,  $\rho$ ,  $\$$ ,  $h$ ) - rnkcnt $_{\mathcal{L}_\mathcal{R}}$ ( $i + 1$ ,  $\rho$ ,  $h$ ,  $h$ );  

6     if rnkcnt $_{\mathcal{L}_\mathcal{R}}$ (1,  $\rho$ , 1,  $\sigma_p$ ) = 0 then return  $\rho$  ; // Case 2  

7      $k \leftarrow \text{CNT}_{S'}[i]$ ;  

8     for  $e \leftarrow \min\{h, \max\{\text{TOP}_{S'}[i], \text{UDN}_{S'}[i]\}\}$  to 0 do // Case 3  

9        $[\ell..r] \leftarrow \text{MI}_{\text{LCP}_\mathcal{R}^\infty}(\max\{1, k\}, e)$ ;  

10      if  $e = h$  then // Case 3.1  

11        if  $x \leftarrow \text{FPQ}_e(\mathcal{L}_\mathcal{R}, k) \in [\ell..r]$  then return  $\text{LF}_\mathcal{R}[x]$  ;  

12        if  $x \leftarrow \text{FNQ}_e(\mathcal{L}_\mathcal{R}, k) \in [\ell..r]$  then return  $\text{LF}_\mathcal{R}[x] - 1$  ;  

13        if  $x \leftarrow \text{FNQ}_{\geq e+1}(\mathcal{L}_\mathcal{R}, \ell) \in [\ell..r]$  then // check for (H1)  

14          return min  $\text{MI}_{\text{LCP}_\mathcal{R}^\infty}(\text{LF}_\mathcal{R}[x], e + 1) - 1$ ;  

15      else // Case 3.2  

16        if  $x \leftarrow \text{FPQ}_{\geq e+1}(\mathcal{L}_\mathcal{R}, k) \in [\ell..r]$  then // Case 3.2.1  

17           $[\ell'..r'] \leftarrow \text{MI}_{\text{LCP}_\mathcal{R}^\infty}(x, e + 1)$ ;  

18          if  $x' \leftarrow \text{FPQ}_{\geq e+2}(\mathcal{L}_\mathcal{R}, r') \in [\ell'..r']$  then  

19            return max  $\text{MI}_{\text{LCP}_\mathcal{R}^\infty}(\text{LF}_\mathcal{R}[x'], e + 2)$ ;  

20          else return  $\text{LF}_\mathcal{R}[x]$  ;  

21        if  $x \leftarrow \text{FNQ}_{\geq e+1}(\mathcal{L}_\mathcal{R}, k) \in [\ell..r]$  then // Case 3.2.2  

22           $[\ell'..r'] \leftarrow \text{MI}_{\text{LCP}_\mathcal{R}^\infty}(x, e + 1)$ ;  

23          if  $x' \leftarrow \text{FNQ}_{e+1}(\mathcal{L}_\mathcal{R}, \ell') \in [\ell'..r']$  then  

24            return  $\text{LF}_\mathcal{R}[x'] - 1$ ;  

25          else return min  $\text{MI}_{\text{LCP}_\mathcal{R}^\infty}(\text{LF}_\mathcal{R}[x], e + 2) - 1$  ;  

26        if  $x \leftarrow \text{FPQ}_e(\mathcal{L}_\mathcal{R}, r) \in [\ell..r]$  then // check for (H1)  

27          return  $\text{LF}_\mathcal{R}[x]$ ;
```

---

**Fact.** Let  $e \in [1.. \min\{h, |U|_p\}]$  and  $\alpha, \beta \in [1..\rho]$ . If  $\text{lcp}(\langle U \rangle, \langle V_\alpha^\omega[..4\lambda+1] \rangle) = t_e$ , then either

- (H1)  $\text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V_\alpha, 1)) \geq e \wedge (h > \pi(V_\alpha) = e \vee \pi(V_\alpha) > h = e)$ , or
- (H2)  $\text{lcp}(\langle \text{Rot}(U, 1) \rangle, \langle \text{Rot}(V_\alpha, 1)^\omega[..4\lambda+1] \rangle) = t_e - 1 \wedge \min\{h, \pi(V_\alpha)\} \geq e$ .

Moreover, if (H1) holds for  $V_\alpha$ ,  $\text{lcp}(\langle U \rangle, \langle V_\beta^\omega[..4\lambda+1] \rangle) = t_e$  and (H2) holds for  $V_\beta$ , then  $\neg(V_\alpha \prec_\omega V_\beta \prec_\omega U)$  and  $\neg(U \prec_\omega V_\beta \prec_\omega V_\alpha)$ .

Fix some  $e \in [1.. \min\{h, \max\{\text{TOP}_{S'}[i], \text{UDN}_{S'}[i]\}\}]$  and let  $[\ell..r] \subseteq [1..\rho]$  such that  $\text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V_\alpha, 1)) \geq e$  if and only if  $\alpha \in [\ell..r]$ , i.e.,  $[\ell..r]$  is the non-empty interval computed in Line 8. We will check if either  $\gamma \in [t_e + 1..t_{e+1}]$  or if there exists some  $\alpha \in [1..\rho]$  such that  $\text{lcp}(\langle U \rangle, \langle V_\alpha^\omega[..4\lambda+1] \rangle) = t_e$  and  $V_\alpha$  satisfies (H1), and show to compute  $\text{CNT}_{S'}[j]$  if either is the case. Note that if neither is the case, then for each  $\alpha \in [1..\rho]$  either  $\text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V_\alpha, 1)) \leq e - 1$  or  $\pi(V_\alpha) \leq e - 1$  by the fact.

**Case 3.1.** Assume  $e = h$ . This case is treated from Line 9 through Line 13. For each  $\alpha \in [1..\rho]$ ,  $\text{lcp}(\langle U \rangle, \langle V_\alpha^\omega[..4\lambda+1] \rangle) > t_e$  if and only if  $\text{L}_R[\alpha] = h$  and  $\alpha \in [\ell..r]$ . By Lemma 3.7,  $\text{Rot}(U, 1) \prec_\omega \text{Rot}(V_\alpha, 1)$  if and only if  $U \prec_\omega V_\alpha$  for each  $\alpha \in [\ell..r]$  with  $\text{L}_R[\alpha] = h$ . Thus, it suffices to compute either the largest  $\alpha \in [\ell.. \text{CNT}_{S'}[i]]$  or smallest  $\alpha \in [\text{CNT}_{S'}[i] + 1..r]$  satisfying  $\text{L}_R[\alpha] = h$ , resulting in  $\text{CNT}_{S'}[j] = \text{LF}_R[\alpha]$  or  $\text{CNT}_{S'}[j] = \text{LF}_R[\alpha] - 1$ , respectively, if they exist. The computation is done in Line 10 and Line 11, respectively, and if neither returns a value, we conclude that  $\gamma \in [1..t_e]$ .

From Line 12 through Line 13 we check if there exists some  $\alpha \in [1..\rho]$  such that  $\text{lcp}(\langle U \rangle, \langle V_\alpha^\omega[..4\lambda+1] \rangle) = t_e$  and  $V_\alpha$  satisfies (H1). But this condition is equivalent to  $\alpha \in [\ell..r]$  and  $\text{L}_R[\alpha] > h$ . If some  $\alpha \in [\ell..r]$  with  $\text{L}_R[\alpha] > h$  exists, then  $\text{CNT}_{S'}[j] = \min \text{MI}_{\text{LCP}_R^\infty}(\text{LF}_R[\alpha], e + 1) - 1$  by the fact.

**Case 3.2.** Assume that either  $\text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V_\alpha, 1)) \leq e$  or that  $\pi(V_\alpha) \leq e$  for each  $\alpha \in [1..\rho]$ . Note that this assumption is fulfilled by definition for  $e = \max\{\text{TOP}_{S'}[i], \text{UDN}_{S'}[i]\} < h$ . Also,  $\gamma \in [1..t_{e+1}]$ .

We first check if  $\gamma \in [t_e + 1..t_{e+1}]$ . For each  $\alpha \in [1..\rho]$ ,  $\text{lcp}^\infty(U, V_\alpha) \geq e + 1$  if and only if  $\alpha \in [\ell..r]$  and  $\text{L}_R[\alpha] \geq e + 1$  by Lemma 3.7. Consequently,  $\gamma \in [t_e + 1..t_{e+1}]$  if some  $\alpha \in [\ell..r]$  with  $\text{L}_R[\alpha] \geq e + 1$  exists.

**Case 3.2.1.** Assume we have an  $x \in [\ell.. \text{CNT}_{S'}[i]]$  such that  $\text{L}_R[x] \geq e + 1$  and  $\alpha \leq x$  for any  $\alpha \in [\ell.. \text{CNT}_{S'}[i]]$  with  $\text{L}_R[\alpha] \geq e + 1$ . We treat this case from Line 15 through Line 19. Let  $t' = \text{select}_\infty(\langle \text{Rot}(V_x, 1) \rangle, e + 1)$  and let  $[\ell'..r'] \subseteq [1..\rho]$  such that  $\text{lcp}(\langle \text{Rot}(V_\alpha, 1) \rangle, \langle \text{Rot}(V_\beta, 1) \rangle) \geq t'$  if and only if

$\alpha, \beta \in [\ell'..r']$ , i.e.,  $[\ell'..r']$  is the interval computed in Line 16. Then  $V_x \prec_{\omega} U$  and  $r' \geq \text{CNT}_{S'}[i]$ . By assumption and Lemma 3.7,  $\text{CNT}_{S'}[j] \in [\ell'..x]$ . Then  $\text{lcp}(\langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[j]])^{\omega}[..4\lambda+1] \rangle, \langle U \rangle) \leq t'$ , which in turn implies that  $\langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[j]])^{\omega}[..4\lambda+1] \rangle$  is prefixed by  $W = \langle V_x^{\omega}[..t'] \rangle$ . If there exists some  $\alpha \in [\ell'..r']$  such that  $L_{\mathcal{R}}[\alpha] \geq e+2$ , then  $\langle V_{\alpha} \rangle$  is prefixed by  $W \cdot \infty$  and consequently  $\text{CNT}_{S'}[j] = \max \text{MI}_{\text{LCP}_{\mathcal{R}}^{\infty}}(\text{LF}_{\mathcal{R}}[\alpha], e+2)$ , with the condition being checked in Line 17 and computation of  $\text{CNT}_{S'}[j]$  taking place in Line 18. Otherwise  $\langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[j]])^{\omega}[..4\lambda+1] \rangle$  is prefixed by  $W \cdot t'$  and consequently  $\text{CNT}_{S'}[j] = \text{LF}_{\mathcal{R}}[x]$ , which is computed in Line 19.

**Case 3.2.2.** Assume we have an  $x \in [\text{CNT}_{S'}[i]+1..r]$  such that  $L_{\mathcal{R}}[x] \geq e+1$  and  $\alpha \geq x$  for any  $\alpha \in [\ell..r]$  with  $L_{\mathcal{R}}[\alpha] \geq e+1$ . We treat this case from Line 20 through Line 24. Let  $t' = \text{select}_{\omega}(\langle \text{Rot}(V_x, 1) \rangle, e+1)$  and let  $[\ell'..r'] \subseteq [1..\rho]$  such that  $\text{lcp}(\langle \text{Rot}(V_{\alpha}, 1) \rangle, \langle \text{Rot}(V_{\beta}, 1) \rangle) \geq t'$  if and only if  $\alpha, \beta \in [\ell'..r']$ , i.e.,  $[\ell'..r']$  is the interval computed in Line 21. Then  $U \prec_{\omega} V_x$  and  $\ell' \leq \text{CNT}_{S'}[j]+1$ . By assumption and Lemma 3.7,  $\text{CNT}_{S'}[j]+1 \in [x..r']$ . Then  $\text{lcp}(\langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[j]+1])^{\omega}[..4\lambda+1] \rangle, \langle U \rangle) \leq t'$ , which in turn implies that  $\langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[j]+1])^{\omega}[..4\lambda+1] \rangle$  is prefixed by  $W = \langle V_x^{\omega}[..t'] \rangle$ . If  $\alpha \in [\ell'..r']$  such that  $L_{\mathcal{R}}[\alpha] = e+1$ , then  $\langle V_{\alpha} \rangle$  is prefixed by  $W \cdot t'$ . Then  $\text{CNT}_{S'}[j]+1 = \text{LF}_{\mathcal{R}}[x']$  if  $x' \in [\ell'..r']$  minimal such that  $L_{\mathcal{R}}[x'] = e+1$  by Lemma 3.7, with the condition checked in Line 22 and the value of  $\text{CNT}_{S'}[j]$  computed in Line 23. If on the other hand  $L_{\mathcal{R}}[\alpha] \neq e+1$  for each  $\alpha \in [\ell'..r']$ , then  $\langle C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[j]+1])^{\omega}[..4\lambda+1] \rangle$  is prefixed by  $W \cdot \infty$  and consequently  $\text{CNT}_{S'}[j]+1 = \min \text{MI}_{\text{LCP}_{\mathcal{R}}^{\infty}}(\text{LF}_{\mathcal{R}}[x], e+2)$ , computed in Line 24. This concludes Case 3.2.2.

If  $L_{\mathcal{R}}[\alpha] \leq e$  for each  $\alpha \in [\ell..r]$ , then  $\gamma \in [1..t_e]$ , which is the case if we reach Line 25. It remains to check if there exists some  $\alpha \in [1..\rho]$  such that  $\text{lcp}(\langle U \rangle, \langle V_{\alpha}^{\omega}[..4\lambda+1] \rangle) = t_e$  and  $V_{\alpha}$  satisfies (H1). But this condition is equivalent to  $\alpha \in [\ell..r]$  and  $L_{\mathcal{R}}[\alpha] = e$ . Then  $\text{CNT}_{S'}[j] = \text{LF}_{\mathcal{R}}[\alpha]$  for the largest  $\alpha \in [\ell..r]$  with  $L_{\mathcal{R}}[\alpha] = e$  by Lemma 3.7 if such an  $\alpha$  exists. This concludes Case 3.

Note that a value is guaranteed to be returned before the for-loop from Line 7 terminates since  $\mathcal{R} \neq \emptyset$ , with correctness of the returned value following from Case 3.  $\square$

For the computation of the entries of  $\text{TOP}_{S'}$  and  $\text{UDN}_{S'}$ , we leverage Lemma 2.7 to arrive at the following.

**Corollary 4.6.** *Let  $U, V \in \Sigma^+$  and  $e = \text{lcp}^{\infty}(\text{Rot}(U, 1), \text{Rot}(V, 1))$ . Then, given  $e$ ,  $\pi(U)$  and  $\pi(V)$ , Algorithm 7 correctly computes  $\text{lcp}^{\infty}(U, V)$ .*

---

**Algorithm 7:** Computing  $\text{lcp}^\infty(U, V)$  for  $U, V \in \Sigma^+$  with  $e = \text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V, 1))$  and  $\text{Rot}(U, 1) \prec_\omega \text{Rot}(V, 1)$ .

---

```

1 Function rotlcp( $e, \pi(U), \pi(V)$ ):
2   if  $\pi(U) = \pi(V) \leq e$  then return  $e$  ;
3   if  $\min\{\pi(U), \pi(V)\} \leq 0$  then return 0 ;
4   if  $\pi(U) \leq \min\{e, \pi(V)\}$  then return  $\pi(U)$  ;
5   if  $\pi(V) \leq \min\{e, \pi(U)\}$  then return  $\pi(V)$  ;
6   return  $e + 1$ ;

```

---

**Lemma 4.7.** Let  $\mathcal{R} \subset (\Sigma - \{\$\})^+$ ,  $S \in (\Sigma - \{\$\})^+$ ,  $S' = S^4\$$ ,  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$  and  $\lambda = |S| \geq \max\{|V| \mid V \in \mathcal{R}\}$ . Algorithm 8 correctly computes  $\text{TOP}_{S'}[\text{LF}_{\{S'\}}[i]]$  and  $\text{UDN}_{S'}[\text{LF}_{\{S'\}}[i]]$  for each  $i \in [1..4\lambda + 1]$ .

*Proof.* The input of Algorithm 8 consists of  $\pi(C_{\{S'\}}(\text{LF}_{\{S'\}}[i]))$ ,  $\text{CNT}_{S'}[i]$ ,  $\text{CNT}_{S'}[\text{LF}_{\{S'\}}[i]]$ ,  $\text{TOP}_{S'}[i]$ ,  $\text{UDN}_{S'}[i]$  for some  $i \in [1..4\lambda + 1]$  and the epBWT index of  $\mathcal{R}$ . We show the statement for  $\text{TOP}_{S'}[\text{LF}_{\{S'\}}[i]]$  since the proof for  $\text{UDN}_{S'}[\text{LF}_{\{S'\}}[i]]$  is similar.

Fix some arbitrary  $i \in [1..4\lambda + 1]$ . We catch the border case where  $\text{CNT}_{S'}[\text{LF}_{\{S'\}}[i]] = 0$  in Line 2. Thus, assume  $\text{CNT}_{S'}[\text{LF}_{\{S'\}}[i]] > 0$ . Let  $U = C_{\mathcal{R}}(\text{CA}_{\mathcal{R}}[\text{CNT}_{S'}[\text{LF}_{\{S'\}}[i]]])$  and  $V = C_{\{S'\}}(\text{CA}_{\{S'\}}[\text{LF}_{\{S'\}}[i]])$ . Then we utilize Lemma 3.14 and the FL-mapping of the epBWT index of  $\mathcal{R}$  to obtain  $\text{lcp}^\infty(\text{Rot}(U, 1), \text{Rot}(V, 1))$  in Line 4 through Line 9, where the different cases are due to the  $\omega$ -preorder of  $\text{Rot}(U, 1)$  and  $\text{Rot}(V, 1)$ . Then the first statement follows from Corollary 4.6 with the computation of  $\text{TOP}_{S'}[\text{LF}_{\{S'\}}[i]] = \text{lcp}^\infty(U, V)$  in Line 10.  $\square$

We give a summarizing result.

**Lemma 4.8.** Let  $\mathcal{R} \subset (\Sigma - \{\$\})^+$ ,  $S \in (\Sigma - \{\$\})^+$ ,  $S' = S^4\$$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$  and  $\lambda = |S| = \max\{|V| \mid V \in \mathcal{S}\}$ . Algorithm 9 correctly computes  $P_{\mathcal{S}}$  and the epBWT index of  $\mathcal{S}$ .

*Proof.* Let  $\rho = |\text{L}_{\mathcal{R}}|$ . The input of Algorithm 9 consists of  $P_{\mathcal{R}}$ ,  $|\text{CR}_{\mathcal{R}}(S^3)|$ ,  $Y_S$ , and the epBWT indexes of  $\mathcal{R}$  and  $\{S'\}$ .

We correctly compute the epBWT index of  $\{S\}$  in Line 2 by Lemma 4.3. If  $\mathcal{R} = \emptyset$ , we simply construct  $P_{\mathcal{S}} = 1^\lambda$  and return the epBWT for the singleton set  $\{S\}$  in Line 3 through Line 5.

---

**Algorithm 8:** Computing  $\text{TOP}_{S'}[j]$  and  $\text{UDN}_{S'}[j]$  from  $\text{CNT}_{S'}[i]$ ,  $\text{CNT}_{S'}[j]$ ,  $\text{TOP}_{S'}[i]$  and  $\text{UDN}_{S'}[i]$  for  $i \in [1..4\lambda + 1]$  with  $j = \text{LF}_{\{S'\}}[i]$ . Here,  $\mathcal{R}, \mathcal{S} \subset (\Sigma - \{\$\})^+$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$ ,  $S' = S^4\$$ ,  $\mathcal{S}' = \mathcal{R} \cup \{S'\}$ ,  $\lambda = |S|$ ,  $\rho = |\mathcal{L}_{\mathcal{R}}|$  and  $h = \pi(C_{\{S'\}}(j)) = F_{\{S'\}}[j]$ .

---

```

1 Function nexttopbot( $h$ ,  $\text{CNT}_{S'}[i]$ ,  $\text{CNT}_{S'}[j]$ ,  $\text{TOP}_{S'}[i]$ ,  $\text{UDN}_{S'}[i]$ ,
2    $F_{\mathcal{R}}$ ,  $L_{\mathcal{R}}$ ,  $LCP_{\mathcal{R}}^{\infty}$ ):
3   if  $\text{CNT}_{S'}[j] = 0$  then  $top \leftarrow -1$ ;                                // compute  $\text{TOP}_{S'}[j]$ 
4   else
5      $k \leftarrow \text{FL}_{\mathcal{R}}[\text{CNT}_{S'}[j]]$ ;
6     if  $k = \text{CNT}_{S'}[i] + 1$  then  $e \leftarrow \text{UDN}_{S'}[i]$  ;
7     else if  $k > \text{CNT}_{S'}[i] + 1$  then
8        $e \leftarrow \min\{\text{UDN}_{S'}[i], \text{RNV}_{LCP_{\mathcal{R}}^{\infty}}(\text{CNT}_{S'}[i] + 2, k, -1)\}$  ;
9     else if  $k = \text{CNT}_{S'}[i]$  then  $e \leftarrow \text{TOP}_{S'}[i]$  ;
10    else  $e \leftarrow \min\{\text{TOP}_{S'}[i], \text{RNV}_{LCP_{\mathcal{R}}^{\infty}}(k + 1, \text{CNT}_{S'}[i], -1)\}$  ;
11     $top \leftarrow \text{rotlcp}(e, h, F_{\mathcal{R}}[\text{CNT}_{S'}[j]])$ ;                      // Alg. 7
12    if  $\text{CNT}_{S'}[j] = \rho$  then  $bot \leftarrow -1$ ;                                // compute  $\text{UDN}_{S'}[j]$ 
13    else
14       $k \leftarrow \text{FL}_{\mathcal{R}}[\text{CNT}_{S'}[j] + 1]$ ;
15      if  $k = \text{CNT}_{S'}[i] + 1$  then  $e \leftarrow \text{UDN}_{S'}[i]$  ;
16      else if  $k > \text{CNT}_{S'}[i] + 1$  then
17         $e \leftarrow \min\{\text{UDN}_{S'}[i], \text{RNV}_{LCP_{\mathcal{R}}^{\infty}}(\text{CNT}_{S'}[i] + 2, k, -1)\}$  ;
18        else if  $k = \text{CNT}_{S'}[i]$  then  $e \leftarrow \text{TOP}_{S'}[i]$  ;
19        else  $e \leftarrow \min\{\text{TOP}_{S'}[i], \text{RNV}_{LCP_{\mathcal{R}}^{\infty}}(k + 1, \text{CNT}_{S'}[i], -1)\}$  ;
20         $bot \leftarrow \text{rotlcp}(e, h, F_{\mathcal{R}}[\text{CNT}_{S'}[j] + 1])$ ;                  // Alg. 7
21
22   return  $top, bot$ ;

```

---

Now, assume  $\mathcal{R} \neq \emptyset$ . Then the condition of Line 3 does not evaluate to true. We initialize  $\text{TOP}_{\mathcal{S}}$  and  $\text{UDN}_{\mathcal{S}}$  and set the initial values of  $\text{cnt} = \text{CNT}_{\mathcal{S}'}[1] = 0$ ,  $\text{bot} = \text{UDN}_{\mathcal{S}'}[1] = 0$ ,  $\text{top} = \text{TOP}_{\mathcal{S}'}[1] = -1$  and the current index in Line 6 through Line 10. Correctness of the initial values follows by choice of  $S'$ . We proceed to compute  $P_{\mathcal{S}}$ ,  $\text{TOP}_{\mathcal{S}}$  and  $\text{UDN}_{\mathcal{S}}$  in Line 11 through Line 24. Correctness follows from Lemma 4.4, Lemma 4.5 and Lemma 4.7. Last, we correctly compute the epBWT index of  $\mathcal{S}$  by Lemma 4.2 in Line 25.

□

The input of Algorithm 9 is computed as follows. We construct the pBWT index of  $S'$  with Lemma 2.8, which is the epBWT index of  $\{S'\}$  by choice of  $S'$ . During construction of the pBWT index,  $Y_{\mathcal{S}}$  is built similarly to the bit string indicating sample positions of the suffix array in the proof of Lemma 2.8 without impacting space or time complexity of the construction algorithm. For the computation of  $|\text{CR}_{\mathcal{R}}(S^3)|$ , we use a count query for  $S^3$  in  $\mathcal{R}$ , which is supported by the epBWT index of  $\mathcal{R}$  as shown in Theorem 3.16. Since  $P_{\mathcal{R}}$  and the epBWT index of  $\mathcal{R}$  are known a priori, we covered all arguments of the function `extend3` of Algorithm 9.

*Remark 4.9.* For the extension of the epBWT from  $\mathcal{R}$  to  $\mathcal{S}$  we assumed that  $\lambda = |S| = \max\{|V| \mid V \in \mathcal{S}\} = z$  for simplicity and because the results are sufficient for our construction algorithm. However, the extension in case of  $\lambda \leq z - 1$  is done similarly with  $S' = S^\omega[..4z] \cdot \$$ , i.e., the epBWT index is extendable in general.

#### 4.3. Constructing the epBWT Index of $\mathcal{T}$

We finally analyze the time and space complexities of the epBWT construction.

**Theorem 4.10.** *Let  $T_1, \dots, T_d \in (\Sigma - \{\$\})^+$ ,  $n = |T_1 \cdots T_d|$  and  $n_k = |T_k|$  for each  $k \in [1..d]$  with  $d \in \mathbb{N}_{\geq 1}$ . If  $n_i \leq n_j$  for each  $i, j \in [1..d]$  with  $i \leq j$ , then the epBWT index of  $\mathcal{T} = \{T_1, \dots, T_d\}$  can be computed in  $O(n \lg \sigma)$  bits of space and  $O(n \frac{\lg \sigma \lg n}{\lg \lg n})$  time.*

*Proof.* Let  $n_i \leq n_j$  for each  $i, j \in [1..d]$  with  $i \leq j$ . For notational convenience, let  $\mathcal{T}_k = \{T_1, \dots, T_k\}$ ,  $T'_k = T_k^4 \$$ ,  $\mathcal{T}'_k = \{T_1, \dots, T_{k-1}, T'_k\}$  for each  $k \in [1..d]$  and  $t_{\text{query}} = \frac{\lg \sigma \lg n}{\lg \lg n}$ . Then for each  $k \in [1..d]$ , the pBWT index of  $T'_k$  is the epBWT index of  $\{T'_k\}$ .

Fix some  $k \in [1..d]$ . We leverage Lemma 2.8 to construct the epBWT index of  $\{T'_k\}$  and  $Y_{T_k}$  in  $O(n_k \lg \sigma)$  bits of space and  $O(n_k t_{\text{query}})$  time. Here,

---

**Algorithm 9:** Extending the epBWT index of  $\mathcal{R}$  to that of  $\mathcal{S}$ .  
 Here,  $\mathcal{R} \subset (\Sigma - \{\$\})^+$ ,  $S \in (\Sigma - \{\$\})^+$ ,  $\mathcal{S} = \mathcal{R} \cup \{S\}$ ,  $S' = S^4\$$ ,  
 $\mathcal{S}' = \mathcal{R} \cup \{S'\}$ ,  $\rho = |\mathcal{L}_{\mathcal{R}}|$  and  $\lambda = |S| = \max\{|V| \mid V \in \mathcal{S}\}$ .

---

```

1 Function extend3( $F_{\mathcal{R}}$ ,  $L_{\mathcal{R}}$ ,  $LCP_{\mathcal{R}}^{\infty}$ ,  $P_{\mathcal{R}}$ ,  $|CR_{\mathcal{R}}(S^3)|$ ,  $Y_S$ ,  $F_{\{S'\}}$ ,  

    $L_{\{S'\}}$ ,  $LCP_{\{S'\}}^{\infty}$ ):  

2    $F_{\{S\}}, L_{\{S\}}, LCP_{\{S\}}^{\infty} \leftarrow \text{singletonepBWT}(Y_S, F_{\{S'\}}, L_{\{S'\}}, LCP_{\{S'\}}^{\infty})$ ;  

   // Alg. 5  

3   if  $P_{\mathcal{R}} = \varepsilon$  then //  $\mathcal{R} = \emptyset$   

4      $P_S \leftarrow$  non-zeroed bit string of length  $\lambda$ ;  

5     return  $F_{\{S\}}, L_{\{S\}}, LCP_{\{S\}}^{\infty}, P_S$ ;  

6   for  $i \leftarrow 1$  to  $\text{rank}_1(P_{\mathcal{R}}, |P_{\mathcal{R}}|)$  do  $P_{\mathcal{R}}[\text{select}_1(P_{\mathcal{R}}, 1)] \leftarrow 0$  ;  

7    $\text{TOP}_{\mathcal{S}}, \text{UDN}_{\mathcal{S}} \leftarrow$  zeroed string of length  $\lambda$  over  $[-1..0]$ ;  

8    $index \leftarrow 1$ ;  

9    $cnt, bot \leftarrow 0$ ; // initialize  $CNT_{\mathcal{S}'}[1] = UDN_{\mathcal{S}'}[1] = 0$   

10   $top \leftarrow -1$ ; // initialize  $TOP_{\mathcal{S}'}[1] = -1$   

11  for  $j \leftarrow 1$  to  $4\lambda$  do //  $CNT_{\mathcal{S}'}, TOP_{\mathcal{S}'}, UDN_{\mathcal{S}'}$  with Alg. 6 and Alg. 8  

12     $index \leftarrow LF_{\{S'\}}[index]$ ;  

13     $nxtcnt \leftarrow \text{nextcnt}(F_{\{S'\}}[index], cnt, top, bot, L_{\mathcal{R}}, LCP_{\mathcal{R}}^{\infty})$ ;  

14     $top, bot \leftarrow \text{nexttopbot}(F_{\{S'\}}[index]), cnt, nxtcnt, top, bot, F_{\mathcal{R}}$ ,  

       $L_{\mathcal{R}}, LCP_{\mathcal{R}}^{\infty}$ ;  

15    if  $Y_S[index] = 1$  then //  $CNT_{\mathcal{S}}, TOP_{\mathcal{S}}, UDN_{\mathcal{S}}$  from Lemma 4.4  

16       $\text{insert}_{P_{\mathcal{R}}}(\text{select}_0(P_{\mathcal{R}}, nxtcnt + |CR_{\mathcal{R}}(S^3)|) + 1, 1)$ ; //  $P_S$   

      from extending zeroed  $P_{\mathcal{R}}$  with  $CNT_{\mathcal{S}}$  values  

17      if  $|CR_{\mathcal{R}}(S^3)| = 0$  then  

18         $TOP_{\mathcal{S}}[\text{rank}_1(Y_S, index)] \leftarrow top$ ;  

19         $UDN_{\mathcal{S}}[\text{rank}_1(Y_S, index)] \leftarrow bot$ ;  

20      else  

21         $TOP_{\mathcal{S}}[\text{rank}_1(Y_S, index)] \leftarrow bot$ ;  

22        if  $nxtcnt = \rho$  then  $UDN_{\mathcal{S}}[\text{rank}_1(Y_S, index)] \leftarrow -1$  ;  

23        else  $UDN_{\mathcal{S}}[\text{rank}_1(Y_S, index)] \leftarrow LCP_{\mathcal{R}}^{\infty}[nxtcnt + 1]$  ;  

24       $cnt \leftarrow nxtcnt$ ;  

25     $F_{\mathcal{S}}, L_{\mathcal{S}}, LCP_{\mathcal{S}}^{\infty} \leftarrow \text{extend2}(P_S, TOP_{\mathcal{S}}, UDN_{\mathcal{S}}, F_{\mathcal{R}}, L_{\mathcal{R}}, LCP_{\mathcal{R}}^{\infty}, F_{\{S\}},$   

       $L_{\{S\}}, LCP_{\{S\}}^{\infty})$ ; // Alg. 4  

26    return  $F_{\mathcal{S}}, L_{\mathcal{S}}, LCP_{\mathcal{S}}^{\infty}, P_{\mathcal{R}}$ ;
```

---

$Y_{T_k}$  is built similarly to the bit string indicating sample positions of the suffix array in the proof of Lemma 2.8,  $L_{\{T'_k\}}$  and  $LCP_{\{T'_k\}}^\infty$  are represented by the data structure of Lemma 2.4 and  $F_{\{T'_k\}}$  and  $Y_{T_k}$  are represented by that of Lemma 2.3. By Theorem 3.16, we compute  $CR_{\mathcal{T}_{k-1}}(T_k^3)$  and its length in  $O(n_k t_{\text{query}})$  time.

By Lemma 4.3, Algorithm 5 correctly computes the epBWT index of  $\{T_k\}$ . We represent  $L_{\{T_k\}}$  and  $LCP_{\{T_k\}}^\infty$  by the data structure of Lemma 2.4 and  $F_{\{T_k\}}$  by that of Lemma 2.3, which implies these occupy  $O(n_k \lg \sigma)$  bits of space in total and their construction takes  $O(n_k t_{\text{query}})$  time.

Assume  $k = 1$ . We initialize  $P_{\{T_1\}} = 1^{n_1}$  and represent it by the data structure of Lemma 2.3, which implies it occupies  $O(n_k)$  bits of space and is constructed in  $O(n_1 t_{\text{query}})$  time. In summary, we computed  $P_{\{T_1\}}$  and the epBWT index of  $\mathcal{T}_1$  in  $O(n_1 t_{\text{query}})$  time and  $O(n_1 \lg \sigma)$  bits of space, with correctness following by Lemma 4.8.

Assume  $k \geq 2$  and that  $P_{\mathcal{T}_{k-1}}$ ,  $F_{\mathcal{T}_{k-1}}$ ,  $L_{\mathcal{T}_{k-1}}$  and  $LCP_{\mathcal{T}_{k-1}}^\infty$  have already been computed with  $P_{\mathcal{T}_{k-1}}$  and  $F_{\mathcal{T}_{k-1}}$  represented by the data structure of Lemma 2.3, and  $L_{\mathcal{T}_{k-1}}$  and  $LCP_{\mathcal{T}_{k-1}}^\infty$  by that of Lemma 2.4. We use Algorithm 9 to extend  $P_{\mathcal{T}_{k-1}}$ ,  $F_{\mathcal{T}_{k-1}}$ ,  $L_{\mathcal{T}_{k-1}}$  and  $LCP_{\mathcal{T}_{k-1}}^\infty$  to  $P_{\mathcal{T}_k}$ ,  $F_{\mathcal{T}_k}$ ,  $L_{\mathcal{T}_k}$  and  $LCP_{\mathcal{T}_k}^\infty$ , respectively, by insert operations, i.e., the extension is done in  $O(n_k \lg \sigma)$  bits of space.

We claim that Algorithm 9 takes  $O(n_k t_{\text{query}})$  time, which is readily confirmed with the exception of the function call in Line 13 due to the for-loop from Line 14 of Algorithm 6. Algorithm 6 computes  $CNT_{\mathcal{T}'_k}[\text{LF}_{\{T'_k\}}[i]]$  for each  $i \in [1..4n_k + 1]$  by invoking  $O(2 + e - e')$  queries and operations, where we assigned  $e' = \max\{\text{TOP}_{\mathcal{T}'_k}[\text{LF}_{\{T'_k\}}[i]], \text{UDN}_{\mathcal{T}'_k}[\text{LF}_{\{T'_k\}}[i]]\}$  and  $e = \max\{\text{TOP}_{\mathcal{T}'_k}[i], \text{UDN}_{\mathcal{T}'_k}[i]\}$ . The value  $e$  can be seen as a potential held by the current conjugate  $C_{\{T'_k\}}(\text{CA}_{\{T'_k\}}[i])$  of  $T'_k$ , which upper bounds the number of queries and operations for a call of Algorithm 6. Since a single call increases the potential for the next conjugate  $C_{\{T'_k\}}(\text{CA}_{\{T'_k\}}[\text{LF}_{\{T'_k\}}[i]])$  of  $T'_k$  by at most 1 according to Lemma 2.7, the total number of invoked queries and operations by the function call in Line 13 of Algorithm 9 is bounded by  $O(n_k)$ . Our claim follows.

Since  $k \in [1..d]$  was chosen arbitrarily and  $n_i \leq n_j$  for each  $i, j \in [1..d]$  with  $i \leq j$ , the construction of the epBWT index of  $\mathcal{T}$  takes  $O((n+n_d) \lg \sigma) = O(n \lg \sigma)$  bits of space and  $O((\sum_{k=1}^d n_k) t_{\text{query}}) = O(nt_{\text{query}})$  time, which concludes the proof.  $\square$

The following result is then a consequence of Theorem 4.10 and Lemma 2.1

through Lemma 2.4.

**Corollary 4.11.** *Let  $T_1, \dots, T_d \in (\Sigma - \{\$\})^+$ ,  $n = |T_1 \cdots T_d|$  and  $n_k = |T_k|$  for each  $k \in [1..d]$  with  $d \in \mathbb{N}_{\geq 1}$ . If  $n_i \leq n_j$  for each  $i, j \in [1..d]$  with  $i \leq j$ , then the static epBWT index of  $\mathcal{T} = \{T_1, \dots, T_d\}$  presented in Theorem 3.13 can be computed in  $O(n \lg \sigma)$  bits of space and  $O(n \frac{\lg \sigma \lg n}{\lg \lg n})$  time.*

Finally, we join together the main results of all sections of this article to affirm the theorem we stated in the very beginning.

*Proof of Theorem 1.1.* The space complexity of the epBWT index follows from Theorem 3.16 and the claims regarding its construction from Theorem 4.10. We can switch to static data structures by Corollary 4.11 and the space complexity of the static epBWT index follows from Theorem 3.13.

Both variants support count queries in the claimed time complexities by Theorem 3.13 and Theorem 3.16. The time complexity for inserting or deleting a string of the set of strings indexed by the dynamic epBWT is due to Remark 4.9.

If  $\llbracket T_1 \rrbracket, \dots, \llbracket T_d \rrbracket$  are primitive, then the claim regarding restoration of the input text p-strings follows from Corollary 3.10. Otherwise, we augment the index with the auxiliary bit string  $D'_{\mathcal{T}}$  as described at the end of Section 3.2. The bit string  $D'_{\mathcal{T}}$  is represented by the data structure of Lemma 2.3 and constructed alongside the epBWT index of  $\mathcal{T}$  without impacting the claimed complexities since  $\exp(\llbracket T_k \rrbracket)$  can be determined by a backward search for  $T_k^3$  on the epBWT index of  $\{T_k^4 \$\}$  for each  $k \in [1..d]$ .  $\square$

*Remark 4.12.* In case we want to report the locations of occurrences, we can construct  $\text{CA}_{\mathcal{T}}$  alongside the epBWT index and store it in  $O(n + d \lg n)$  bits of space by means of sampling each  $(\lg n)$ -th text position [3, Sect. 3.2], where the latter term is due to the  $O(d)$  cycles of the LF-mapping corresponding to the  $\llbracket .. \rrbracket$ -encoded roots of the input text p-strings, which require at least a sample each.

## 5. Matching Statistics

As an application, we give a generalization of the matching statistics defined by Chang and Lawler [16]. The *circular matching statistics* between  $\mathcal{T} \subset \Sigma^+$  and a pattern p-string  $P \in \Sigma^+$  of length  $m$  is an array  $M[1..m]$ , whose entries are triples of a prefix length  $M[i].q \in [0..m]$ , a left interval

$i$	1	2	3	4	5	6	7	8	9	10	11	12
$E_{\mathcal{T}}$	0	0	2	0	1	1	15	3	1	2	15	2

Figure 10: The string  $E_{\mathcal{T}}$  of our running example. Its values are readily confirmed by, e.g., the last column of Figure 5. Here, since  $n' = n_3 = 5$ , an entry  $i$  of  $E_{\mathcal{T}}$  evaluating to 15 indicates  $\omega$ -equivalence of  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])$  and  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i-1])$  by Lemma 3.4.

boundary  $M[i].\ell \in [1..n]$  and a right interval boundary  $M[i].r \in [1..n]$  such that for each  $i \in [1..m]$

- (i)  $M[i].q = \max\{\text{lcp}(\langle C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^{\omega}[..m] \rangle, \langle P[i..] \rangle) \mid j \in [1..n]\}$ ,
- (ii)  $M[i].\ell \leq M[i].r$  and
- (iii)  $j \in [M[i].\ell..M[i].r] \Leftrightarrow \text{lcp}(\langle C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^{\omega}[..m] \rangle, \langle P[i..] \rangle) = M[i].q$ .

See Figure 11 for an example. We adapt the results of Ohlebusch et al. [46] and compute the circular matching statistics with the backward search of the epBWT index. Their main idea is to compute sufficiently many conjugate ranges of  $P[i..j]$  with  $i, j \in [1..m]$ . We start with  $i = j = m$  and decrease  $i$  until  $\text{CR}_{\mathcal{T}}(P[i..j])$  is empty, and then decrease  $j$  until  $\text{CR}_{\mathcal{T}}(P[i..j])$  is non-empty, called a *contraction* operation. We repeat until  $i = 1$  and either  $j = 0$  or  $\text{CR}_{\mathcal{T}}(P[1..j])$  is non-empty. For the contraction, we use an auxiliary string that helps us to find the correct range if we end up with a mismatch. Let  $n' = \max\{n_k \mid k \in [1..d]\}$  and denote by  $E_{\mathcal{T}}$  a string of length  $n$  over  $[0..3n']$  such that  $E_{\mathcal{T}}[1] = 0$  and

$$E_{\mathcal{T}}[i] = \text{lcp}(\langle C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i-1])^{\omega}[..3n'] \rangle, \langle C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])^{\omega}[..3n'] \rangle)$$

for each  $i \in [2..n]$ . Note that by Lemma 3.4 and for each  $i \in [2..n]$ ,  $C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i-1]) =_{\omega} C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[i])$  if and only if  $E_{\mathcal{T}}[i] = 3n'$ . See Figure 10 for an example.

**Lemma 5.1.** *Algorithm 10 correctly computes the circular matching statistics between  $\mathcal{T} \subset \Sigma^+$  and a pattern  $p$ -string  $P \in \Sigma^+$ .*

The statement is a consequence of Ohlebusch et al. [46, Section 3], but for completeness we give a detailed proof.

*Proof.* Algorithm 10 takes a pattern  $P \in \Sigma^+$ , the static epBWT index of  $\mathcal{T}$  presented in Theorem 3.13 and  $E_{\mathcal{T}}$  as input. We address the computation in

---

**Algorithm 10:** Computing the circular matching statistics between  $\mathcal{T} \subset \Sigma^+$  and a pattern p-string  $P \in \Sigma^+$ .

---

```

1 Function cms( $P, F_{\mathcal{T}}, L_{\mathcal{T}}, LF_{\mathcal{T}}, E_{\mathcal{T}}$ ):
2   initialize  $M$  and the binary search tree of Lemma 3.12;
3    $(q, \ell, r, i) \leftarrow (0, 1, n, |P|)$ ;
4   while  $i \geq 1$  do
5     compute  $\pi(P[i..i+q])$  and  $|P[i+1..i+q]|_p$ ;
6      $[\ell'..r'] \leftarrow \text{crangeupdrmq}([\ell..r], \pi(P[i..i+q]), |P[i+1..i+q]|_p,$ 
7        $F_{\mathcal{T}}, L_{\mathcal{T}}, LF_{\mathcal{T}})$ ;                                // Alg. 1
8     if  $\ell' \leq r'$  then                                // expand previous match
9        $q \leftarrow q + 1$ ;
10       $M[i] \leftarrow (q, \ell', r')$ ;
11       $[\ell..r] \leftarrow [\ell'..r']$ ;
12       $i \leftarrow i - 1$ ;
13    else if  $q = 0$  then                                // reset
14       $M[i] \leftarrow (0, 1, n)$ ;
15       $i \leftarrow i - 1$ ;
16    else                                              // contract previous match
17       $q \leftarrow q - 1$ ;
18       $[\ell..r] \leftarrow MI_{E_{\mathcal{T}}}(r, q)$ ;
19   return  $M$ ;

```

---

Line 5 in the proof of Corollary 5.2. All branches except Line 15 decrement  $i$ . However, this branch decrements  $q$ . Since  $q \leq m$ , the algorithm terminates.

For correctness, we consider the loop invariant that  $M[i+1..]$  has been correctly computed and  $[\ell..r] = CR_{\mathcal{T}}(P[i+1..i+q])$ . Clearly, this invariant holds initially. Assume the invariant holds for a tuple  $(i, q)$  with  $i \in [1..m]$ ,  $q \in [0..m-i]$  and  $[\ell..r] = CR_{\mathcal{T}}(P[i+1..i+q])$ , i.e., if  $i < m$ , then  $P[i+1..i+M[i].q]$  is the longest substring of  $P$  starting at position  $i$  that p-matches some prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$ .

**Case 1.** Assume  $i = m$ . Then  $q = 0$ . If  $[\ell'..r'] = CR_{\mathcal{T}}(P[m])$  is non-empty, then  $P[m]$  is the longest substring of  $P[m..] = P[m]$  starting at position  $i = m$  that p-matches some prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$ , i.e., we enter Line 7 where we set  $M[m] = (q+1, \ell', r')$ , adjust  $i, q$  and update  $[\ell..r]$ . If on the other hand  $[\ell'..r']$  is empty, then  $\varepsilon$

is the longest substring of  $P[m..] = P[m]$  starting at position  $i = m$  that p-matches some prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$ , i.e., we enter Line 12, where we set  $M[m] = (0, 1, n)$ , decrement  $i$  by 1 and leave  $q = 0$  and  $[\ell..r] = [1..n]$  as is. In either case, the invariant holds after the iteration of the loop.

**Case 2.** Assume  $i < m$  and  $q = M[i + 1].q$ . If  $q = 0$ , then the proof is similar to that of Case 1. So assume  $q \geq 1$ . If  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..i + q])$  is non-empty, then  $P[i..i + q]$  is the longest substring of  $P$  starting at position  $i$  that p-matches some prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$  by assumption on  $M[i + 1..]$ , i.e., we enter Line 7, where we set  $M[i]$  and maintain the invariant. If  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..i + q])$  is empty, then we will consider  $P[i..i + q - 1]$  as a candidate for the longest substring of  $P$  starting at position  $i$  that p-matches some prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$  in the next iteration, i.e., we enter Line 15 and maintain the invariant by adjusting  $q$  and  $[\ell..r]$ .

**Case 3.** Assume  $i < m$  and  $0 \leq q < M[i + 1].q$ . Note that we only enter this case if  $\text{CR}_{\mathcal{T}}(P[i..i + q + k])$  does not p-match any prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$  for each  $k \in \mathbb{N}$ . If  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..i + q])$  is empty, then we argue similarly to Case 1 if  $q = 0$  and similarly to Case 2 if  $q \neq 0$ . Thus, assume  $[\ell'..r'] = \text{CR}_{\mathcal{T}}(P[i..i + q])$  is non-empty. Then  $P[i..i + q]$  is the longest substring of  $P$  starting at position  $i$  that p-matches some prefix of an infinitely iterated conjugate considered in  $\mathcal{T}$  by our remark, i.e., we enter Line 7, where we set  $M[i]$  and maintain the invariant.

□

**Corollary 5.2.** *There exists a  $2n \lg \sigma + n \lg 3n' + 4n + O(n \frac{\lg \lg n}{\lg n})$  bits of space data structure that supports the computation of circular matching statistics between  $\mathcal{T} \subset \Sigma^+$  and any pattern p-string in  $O(m(\lg \sigma + \lg n))$  time, where  $m$  is the pattern length and  $n'$  is the maximum length of all texts in  $\mathcal{T}$ .*

*Proof.* The data structure consists of the static epBWT index from Theorem 3.13,  $E_{\mathcal{T}}$  represented by the data structure of Lemma 2.2 and an additional copy of  $E_{\mathcal{T}}$  as a simple string for access occupying  $n \lg 3n'$  bits of space, which results in the claimed space complexity.

We examine Algorithm 10 for the claim regarding time complexity. To obtain the values in Line 5 in  $O(\lg \sigma)$  time, we proceed similarly to the proof of Lemma 3.12, but possibly remove nodes in the binary search tree if we enter Line 12 or Line 15. Then Line 6 is computed in  $O(\lg \sigma)$  time by the proof of Theorem 3.13. We can perform the MI query of Line 17

$i$	1	2	3	4	5	6	7	8
$P[i]$	C	B	A	A	b	C	\$	b
$M[i].q$	3	3	2	3	2	1	0	1
$M[i].\ell$	12	10	6	5	2	4	1	2
$M[i].r$	12	11	8	5	3	12	12	3

Figure 11: The circular matching statistics  $M$  between the pattern p-string  $P = \text{CBAAbC\$b}$  and our running example  $\mathcal{T} = \{\text{Bab}, \text{ABBA}, \text{CAbBB}\}$ .

by exponential search on the RmQ values of ranges in  $E_{\mathcal{T}}$ . The exponential search requires  $O(\lg n)$  RmQ queries, and thus the computation takes  $O(\lg n)$  time by Lemma 2.2. As each of the cases in the while-loop is entered at most  $m$  times until termination, the claim follows.  $\square$

We compute the circular matching statistics  $M$  between our running example  $\mathcal{T} = \{\text{Bab}, \text{ABBA}, \text{CAbBB}\}$  and the pattern p-string  $P = \text{CBAAbC\$a}$ . The solution is also presented in Figure 11.

We begin by computing  $\text{CR}_{\mathcal{T}}(P[8..]) = \text{CR}_{\mathcal{T}}(\text{b}) = [2..3]$ , which is non-empty. Thus,  $M[8] = (1, 2, 3)$ . Next,  $\text{CR}_{\mathcal{T}}(P[7..]) = \text{CR}_{\mathcal{T}}(\$b) = \emptyset$ , and consequently we contract the previous match and compute  $\text{CR}_{\mathcal{T}}(P[7..7]) = \text{CR}_{\mathcal{T}}(\$) = \emptyset$ . We conclude that  $\varepsilon$  is the longest prefix of  $P[7..]$  p-matching an infinitely iterated conjugate considered in  $\mathcal{T}$ , set  $M[7] = (0, 1, 12)$ , and reset. Afterwards, we compute  $\text{CR}_{\mathcal{T}}(P[6..6]) = \text{CR}_{\mathcal{T}}(\text{C}) = [4..12]$ ,  $\text{CR}_{\mathcal{T}}(P[5..6]) = \text{CR}_{\mathcal{T}}(\text{bC}) = [2..3]$  and  $\text{CR}_{\mathcal{T}}(P[4..6]) = \text{CR}_{\mathcal{T}}(\text{AbC}) = [5..5]$ , which result in  $M[6] = (1, 4, 12)$ ,  $M[5] = (2, 5, 6)$  and  $M[4] = (3, 5, 5)$ , respectively. We get empty intervals for  $\text{CR}_{\mathcal{T}}(P[3..6]) = \text{CR}_{\mathcal{T}}(\text{AAbC})$  and  $\text{CR}_{\mathcal{T}}(P[3..5]) = \text{CR}_{\mathcal{T}}(\text{AAb})$ , but then  $\text{CR}_{\mathcal{T}}(P[3..4]) = \text{CR}_{\mathcal{T}}(\text{AA}) = [6..8]$ , which yields  $M[3] = (2, 6, 8)$ . We have  $\text{CR}_{\mathcal{T}}(P[2..4]) = \text{CR}_{\mathcal{T}}(\text{BAA}) = [10..11]$  and consequently  $M[2] = (3, 10, 11)$ . The conjugate range  $\text{CR}_{\mathcal{T}}(P[1..4]) = \text{CR}_{\mathcal{T}}(\text{CBAA})$  is empty again, but  $\text{CR}_{\mathcal{T}}(P[1..3]) = \text{CR}_{\mathcal{T}}(\text{CBA}) = [12..12]$  and thus  $M[1] = (3, 12, 12)$ .

*Remark 5.3.* We can construct  $E_{\mathcal{T}}$  alongside the epBWT index of  $\mathcal{T}$  as follows. We represent  $E_{\mathcal{T}}$  by the data structure of Lemma 2.4, occupying  $O(n \lg n)$  bits of space at the expense of compact space. Using the setting of Section 4, we assume to have already built the epBWT index of  $\mathcal{R} \subset (\Sigma - \{\$\})^+$  and  $E_{\mathcal{R}}$  and want to extend it by some  $S \in (\Sigma - \{\$\})^+$  with  $S$  not shorter than the longest text in  $\mathcal{R}$ . Initially  $E_{\emptyset} = \varepsilon$ . Given  $S' = S^4\$$ , we compute both  $E_{\{S'\}}$  and the circular matching statistics between  $\mathcal{R}$  and  $S'$ , where we note that the former is simply the parameterized LCP array [30]

$i$	1	2	3	4	5	$i$	1	2	3	4	5
$P[i]$	A	A	C	C	b	$P[i]$	A	A	C	C	b
$M[i].q$	4	3	2	2	1	$M'[i].q$	4	3	2	4	4
$M[i].\ell$	7	10	6	5	2	$M'[i].\ell$	7	10	6	5	3
$M[i].r$	8	11	8	5	3	$M'[i].r$	8	11	8	5	3

Figure 12: Comparison of the matching statistics  $M$  and  $M'$  between the pattern p-string  $P = \text{AACCb}$  and our running example  $\mathcal{T} = \{\text{Bab}, \text{ABBA}, \text{CAbBB}\}$ .

of  $S'$  by choice of  $S'$ . We use the circular matching statistics to compute the equivalent of  $\text{TOP}_{\mathcal{S}}$  and  $\text{UDN}_{\mathcal{S}}$ , and then proceed similarly to Line 3 through Line 9 of Algorithm 4 to obtain  $E_{\mathcal{S}}$ .

### 5.1. Circular Pattern p-Strings

We can also consider a pattern p-string to be circular, which motivates the following definition. Given  $\mathcal{T} \subset \Sigma^+$  and a pattern p-string  $P \in \Sigma^+$  of length  $m$ , the *modified circular matching statistics*  $M'$  denotes an array whose entries are triples of a prefix length  $M'[i].q \in [0..m]$ , a left interval boundary  $M'[i].\ell \in [1..n]$  and a right interval boundary  $M'[i].r \in [1..n]$  such that for each  $i \in [1..m]$

- (i)  $M'[i].q = \max\{\text{lcp}(\langle C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^{\omega}[..m] \rangle, \langle P[i..] \cdot P[..i-1] \rangle) \mid j \in [1..n]\}$ ,
- (ii)  $M'[i].\ell \leq M'[i].r$  and
- (iii)  $j \in [M'[i].\ell..M'[i].r] \Leftrightarrow \text{lcp}(\langle C_{\mathcal{T}}(\text{CA}_{\mathcal{T}}[j])^{\omega}[..m] \rangle, \langle P[i..] \cdot P[..i-1] \rangle) = M'[i].q$ .

We give an example in Figure 12. To compute  $M'$ , we slightly modify Algorithm 10. In detail, we first compute  $M$  and then initialize the values in Line 3 with  $(M[1].q, M[1].\ell, M[1].r, |P|)$ . Also, in the while-loop we first check for  $q = |P|$  and contract the match as we do in Line 15 if the statement holds.

As a side product, we also solve the circular pattern matching problem in the parameterized setting, where we consider texts and the pattern as circular strings.

## 6. Conclusion

We introduced the epBWT-index as a new type of index capable to search a pattern in a set of circular texts in the parameterized pattern matching setting. The main goal of the epBWT-index is to solve the problem of counting the number of pattern occurrences, which we formalized as EP-COUNT. The main ideas stem from combining core concepts of (a) the pBWT, which is a parameterized index for a single non-circular text, and (b) the eBWT, which is an index for multiple circular texts. We could show that it is possible to solve EP-COUNT with our epBWT within the same space and time complexities as Kim and Cho’s pBWT variant [14]. Further, we considered the input set of circular texts to be dynamic, supporting the insertion and the removal of a text. To that end, we followed the approach of Iseri et al. [15] who proposed a construction algorithm of Kim and Cho’s pBWT, and extended this construction algorithm to construct and extend the epBWT. We inherited the time complexity of the technique of Iseri et al., which makes the epBWT construction faster for general parameterized alphabets to construction algorithms for other parameterized indexes, whose time complexities usually have a multiplicative dependency on the input length and the size of the parameterized alphabet. Finally, as an application of the epBWT, we show how to compute generalized versions of matching statistics to parameterized and circular matching. We conclude this article with future work spawning two different research directions: lowering the space and adaptation to other generalized pattern matchings.

### 6.1. Lowering Space

First, we did not address the pBWT of Ganguly et al. [13] in detail, who can escape the requirement of storing  $F_{\mathcal{T}}$  by having a compressed suffix tree topology at hand. It is unclear whether we can make use of their approach, in particular if we want to emphasize on construction.

Second, to actually locate all the reported occurrences within compact space, we want to shrink the conjugate array sampling addressed in Remark 4.12 to fit into  $O(n)$  bits of space. The main question here is how to get rid of the  $O(d \lg n)$  term, which is due to the LF-mapping having at least one cycle for each of the indexed texts.

Finally, for matching statistics to work in compact space, the string  $E_{\mathcal{T}}$  remains the bottleneck. An idea is to investigate new compression techniques for  $E_{\mathcal{T}}$ . However, a major problem is that the values of  $E_{\mathcal{T}}$  are based on the

*prev-encoded* conjugates, for which a permuted LCP-like compression [47] seems not to work. Also, we are unaware of how to compute  $E_{\mathcal{T}}$  efficiently since we rely on the construction of the parameterized LCP array, whose time is only linear for constant alphabets [34].

### 6.2. Practical Considerations

Due to the fact that the standard FM-index is a well-perceived data structure for practical application, especially in bioinformatics, it is natural to investigate the practical performance of our epBWT-index. For that it looks promising to first gather more ideas on algorithmic improvements on the practical side, which may not have an impact on the asymptotic time complexity — there are certainly several screws to turn, such as improving Lemma 3.4 to compare  $2z$  instead of  $3z$  characters. (This was pointed out by one reviewer of this manuscript.)

*Acknowledgements.* This research was supported by JSPS KAKENHI with grant numbers JP23H04378 and JP25K21150.

## Appendix A. Variables

The table below gives an overview of variables used throughout this article.

variable name	description
$i, j, k, v, x, y, \alpha, \beta$	index or natural number
$\Sigma, \Sigma_s, \Sigma_p$	alphabet
$\psi$	permutation of $\Sigma_p$
$\sigma, \sigma_s, \sigma_p$	size of $\Sigma, \Sigma_s, \Sigma_p$
$\$, \infty, f, g$	symbol
$\varepsilon$	empty string
$U, W, V, V_1, ..$	string
$S, T, T_1, ..$	text string
$\lambda, n, n_1, ..$	length of $S, T, T_1, ..$
$S', T', T'_1, ..$	$S^4\$, T^4\$, T_1^4\$, ..$
$\mathcal{S}, \mathcal{R}, \mathcal{T}, \mathcal{T}_1, ..$	set of text strings
$d$	number of texts in $\mathcal{T}$
$\rho$	total length of all texts in $\mathcal{R}$
$z, n'$	maximum length of text in a set of texts

$\mathcal{S}', \mathcal{T}'_1, \mathcal{T}'_2, \dots$	$\mathcal{R} \cup \{S'\}, \{T'_1\}, \mathcal{T}_1 \cup \{T'_2\}, \dots$
$P, P_1, \dots$	pattern string
$m$	length of $P$
$\text{SA}_T, \text{SA}_T^{-1}$	parameterized suffix array and its inverse
$\text{SR}_T(P)$	suffix range of $P$ in $T$
$\ell, \ell', \ell'', r, r', r''$	interval boundary
$\text{F}_T, \text{L}_T, \text{LCP}_T^\infty$	the pBWT index of $T$
$\text{LF}_T$	LF-mapping on the suffix array of $T$
$\text{A}_T, \text{B}_T$	augmenting pBWT index to solve EP-COUNT
$C_{\mathcal{T}}(i)$	a conjugate of a text in $\mathcal{T}$
$\text{CA}_{\mathcal{T}}, \text{CA}_{\mathcal{T}}^{-1}$	conjugate array of $\mathcal{T}$ and its inverse
$\text{CR}_{\mathcal{T}}(P)$	conjugate range of $P$ in $\mathcal{T}$
$\delta$	length of a conjugate range
$\text{LEN}_{\mathcal{T}}, \text{PRV}_{\mathcal{T}}, \text{ENC}_{\mathcal{T}}$	conceptual string for defining the epBWT of $\mathcal{T}$
$\text{LF}_{\mathcal{T}}, \text{FL}_{\mathcal{T}}$	LF-/FL-mapping of the conjugate array of $\mathcal{T}$
$\text{F}_{\mathcal{T}}, \text{L}_{\mathcal{T}}, \text{LCP}_{\mathcal{T}}^\infty$	the epBWT index of $\mathcal{T}$
$h$	a return value of $\pi(\dots)$
$c$	a return value of $ .. _p$
$\text{CNT}_{\mathcal{S}}$	conceptual string for extending epBWT index
$\text{P}_{\mathcal{S}}, \text{Y}_{\mathcal{S}}, \text{TOP}_{\mathcal{S}}, \text{UDN}_{\mathcal{S}}$	auxiliary string for epBWT index extension
$q, \gamma, t', t_1, \dots$	length of a prefix
$e, e'$	number of $\infty$ in a prefix
$\text{D}_{\mathcal{T}}, \text{D}'_{\mathcal{T}}$	auxiliary string for restoring indexed texts
$M, M'$	matching statistics
$\text{E}_{\mathcal{T}}$	conceptual string for contraction operation

## References

- [1] E. M. Osterkamp, D. Köppl, Extending the parameterized Burrows–Wheeler transform, in: Proc. DCC, 2024, pp. 143–152.
- [2] G. Bell, T. Hey, A. Szalay, Beyond the data deluge, *Science* 323 (5919) (2009) 1297–1298. doi:[10.1126/science.1170411](https://doi.org/10.1126/science.1170411).
- [3] P. Ferragina, G. Manzini, Opportunistic data structures with applications, in: Proc. FOCS, 2000, pp. 390–398. doi:[10.1109/SFCS.2000.892127](https://doi.org/10.1109/SFCS.2000.892127).
- [4] M. Burrows, D. J. Wheeler, A block sorting lossless data compression algorithm, Tech. Rep. 124, Digital Equipment Corporation, Palo Alto, California (1994).
- [5] R. Grossi, J. S. Vitter, B. Xu, Wavelet trees: From theory to practice, in: Proc. CCP, 2011, pp. 210–221. doi:[10.1109/CCP.2011.16](https://doi.org/10.1109/CCP.2011.16).
- [6] C. S. Iliopoulos, M. S. Rahman, Indexing circular patterns, in: Proc. WALCOM, Vol. 4921 of LNCS, 2008, pp. 46–57. doi:[10.1007/978-3-540-77891-2\\_5](https://doi.org/10.1007/978-3-540-77891-2_5).
- [7] B. L. Strang, N. D. Stow, Circularization of the herpes simplex virus type 1 genome upon lytic infection, *Journal of virology* 79 (19) (2005) 12487–12494.
- [8] T. B. H. Karoline K. Ebbesen, J. Kjems, Insights into circular RNA biology, *RNA Biology* 14 (8) (2017) 1035–1045. doi:[10.1080/15476286.2016.1271524](https://doi.org/10.1080/15476286.2016.1271524).
- [9] A. F. Nielsen, A. Bindereif, I. Bozzoni, M. Hanan, T. B. Hansen, M. Irimia, S. Kadener, L. S. Kristensen, I. Legnini, M. Morlando, M. T. Jarlstad Olesen, R. J. Pasterkamp, S. Preibisch, N. Rajewsky, C. Suenkel, J. Kjems, Best practice standards for circular rna research, *Nature Methods* 19 (10) (2022) 1208–1220. doi:[10.1038/s41592-022-01487-2](https://doi.org/10.1038/s41592-022-01487-2).
- [10] S. Mantaci, A. Restivo, G. Rosone, M. Sciortino, An extension of the Burrows–Wheeler transform, *Theor. Comput. Sci.* 387 (3) (2007) 298–312. doi:[10.1016/j.tcs.2007.07.014](https://doi.org/10.1016/j.tcs.2007.07.014).

- [11] B. S. Baker, A theory of parameterized pattern matching: algorithms and applications, in: Proc. STOC, 1993, pp. 71–80. doi:10.1145/167088.167115.
- [12] T. Shibuya, Generalization of a suffix tree for RNA structural pattern matching, *Algorithmica* 39 (1) (2004) 1–19. doi:10.1007/S00453-003-1067-9.
- [13] A. Ganguly, R. Shah, S. V. Thankachan, pBWT: Achieving succinct data structures for parameterized pattern matching and related problems, in: Proc. SODA, 2017, pp. 397–407. doi:10.1137/1.9781611974782.25.
- [14] S. Kim, H. Cho, Simpler FM-index for parameterized string matching, *Inf. Process. Lett.* 165 (2021) 106026. doi:10.1016/j.ipl.2020.106026.
- [15] K. Iseri, T. I, D. Hendrian, D. Köppl, R. Yoshinaka, A. Shinohara, Breaking a barrier in constructing compact indexes for parameterized pattern matching, in: Proc. ICALP, Vol. 297 of LIPIcs, 2024, pp. 89:1–89:19. doi:10.4230/LIPICS.ICALP.2024.89.
- [16] W. I. Chang, E. L. Lawler, Sublinear approximate string matching and biological applications, *Algorithmica* 12 (4/5) (1994) 327–344. doi:10.1007/BF01185431.
- [17] P. Weiner, Linear pattern matching algorithms, in: Proc. SWAT, 1973, pp. 1–11. doi:10.1109/SWAT.1973.13.
- [18] I. M. Gessel, C. Reutenauer, Counting permutations with given cycle structure and descent set, *J. Comb. Theory, Ser. A* 64 (2) (1993) 189–215.
- [19] W. Hon, C. Lu, R. Shah, S. V. Thankachan, Succinct indexes for circular patterns, in: Proc. ISAAC, Vol. 7074 of LNCS, 2011, pp. 673–682. doi:10.1007/978-3-642-25591-5\\_69.
- [20] W. Hon, T. Ku, C. Lu, R. Shah, S. V. Thankachan, Efficient algorithm for circular Burrows–Wheeler transform, in: Proc. CPM, Vol. 7354 of LNCS, 2012, pp. 257–268. doi:10.1007/978-3-642-31265-6\\_21.

- [21] N. Cotumaccio, Improved circular dictionary matching (2025). arXiv: 2504.03394.  
URL <https://arxiv.org/abs/2504.03394>
- [22] S. Bonomo, S. Mantaci, A. Restivo, G. Rosone, M. Sciortino, Sorting conjugates and suffixes of words in a multiset, *Int. J. Found. Comput. Sci.* 25 (8) (2014) 1161. doi:10.1142/S0129054114400309.
- [23] A. Policriti, N. Prezza, Computing LZ77 in run-compressed space, in: Proc. DCC, 2016, pp. 23–32. doi:10.1109/DCC.2016.30.
- [24] T. Ohno, Y. Takabatake, T. I, H. Sakamoto, A faster implementation of online run-length Burrows–Wheeler transform, in: Proc. IWOCA, Vol. 10765 of LNCS, 2017, pp. 409–419. doi:10.1007/978-3-319-78825-8\\_33.
- [25] H. Bannai, J. Kärkkäinen, D. Köppl, M. Piątkowski, Constructing and indexing the bijective and extended Burrows–Wheeler transform, *Inf. Comput.* 297 (105153) (2024) 1–30. doi:10.1016/j.ic.2024.105153.
- [26] C. Boucher, D. Cenzato, Z. Lipták, M. Rossi, M. Sciortino, r-indexing the eBWT, *Inf. Comput.* 298 (2024) 105155. doi:10.1016/j.ic.2024.105155.
- [27] J. Olbrich, E. Ohlebusch, T. Büchler, Generic non-recursive suffix array construction, *ACM Trans. Algorithms* 20 (2) (2024) 18. doi:10.1145/3641854.
- [28] B. S. Baker, Parameterized duplication in strings: Algorithms and an application to software maintenance, *SIAM J. Comput.* 26 (5) (1997) 1343–1362. doi:10.1137/S0097539793246707.
- [29] J. Mendivilso, S. V. Thankachan, Y. J. Pinzón, A brief history of parameterized matching problems, *Discret. Appl. Math.* 274 (2020) 103–115. doi:10.1016/j.dam.2018.07.017.
- [30] S. Deguchi, F. Higashijima, H. Bannai, S. Inenaga, M. Takeda, Parameterized suffix arrays for binary strings, in: Proc. PSC, 2008, pp. 84–94.
- [31] T. I, S. Deguchi, H. Bannai, S. Inenaga, M. Takeda, Lightweight parameterized suffix array construction, in: Proc. IWOCA, Vol. 5874 of LNCS, 2009, pp. 312–323. doi:10.1007/978-3-642-10217-2\\_31.

- [32] R. Beal, D. A. Adjeroh, Variations of the parameterized longest previous factor, *J. Discrete Algorithms* 16 (2012) 129–150. doi:[10.1016/j.jda.2012.05.004](https://doi.org/10.1016/j.jda.2012.05.004).
- [33] D. Hendrian, T. Katsura, Y. Otomo, K. Narisawa, A. Shinohara, Position heaps for parameterized strings, in: Proc. CPM, Vol. 78 of LIPIcs, 2017, pp. 8:1–8:13. doi:[10.4230/LIPICS.CPM.2017.8](https://doi.org/10.4230/LIPICS.CPM.2017.8).
- [34] N. Fujisato, Y. Nakashima, S. Inenaga, H. Bannai, M. Takeda, Direct linear time construction of parameterized suffix and LCP arrays for constant alphabets, in: Proc. SPIRE, Vol. 11811 of LNCS, 2019, pp. 382–391. doi:[10.1007/978-3-030-32686-9\\_27](https://doi.org/10.1007/978-3-030-32686-9_27).
- [35] K. Nakashima, D. Hendrian, R. Yoshinaka, A. Shinohara, An extension of linear-size suffix tries for parameterized strings, in: Proc. SOFSEM, Vol. 2568 of CEUR Workshop Proceedings, 2020, pp. 97–108.
- [36] N. Fujisato, Y. Nakashima, S. Inenaga, H. Bannai, M. Takeda, The parameterized suffix tray, in: Proc. CIAC, Vol. 12701 of LNCS, 2021, pp. 258–270. doi:[10.1007/978-3-030-75242-2\\_18](https://doi.org/10.1007/978-3-030-75242-2_18).
- [37] A. Ganguly, R. Shah, S. V. Thankachan, Fully functional parameterized suffix trees in compact space, in: Proc. ICALP, Vol. 229 of LIPIcs, 2022, pp. 65:1–65:18. doi:[10.4230/LIPIcs.ICALP.2022.65](https://doi.org/10.4230/LIPIcs.ICALP.2022.65).
- [38] K. Nakashima, N. Fujisato, D. Hendrian, Y. Nakashima, R. Yoshinaka, S. Inenaga, H. Bannai, A. Shinohara, M. Takeda, Parameterized DAWGs: Efficient constructions and bidirectional pattern searches, *Theor. Comput. Sci.* 933 (2022) 21–42. doi:[10.1016/J.TCS.2022.09.008](https://doi.org/10.1016/J.TCS.2022.09.008).
- [39] D. Hashimoto, D. Hendrian, D. Köppl, R. Yoshinaka, A. Shinohara, Computing the parameterized Burrows–Wheeler transform online, in: Proc. SPIRE, Vol. 13617 of LNCS, 2022, pp. 70–85. doi:[10.1007/978-3-031-20643-6\\_6](https://doi.org/10.1007/978-3-031-20643-6_6).
- [40] E. Ohlebusch, J. Fischer, S. Gog, CST++, in: Proc. SPIRE, Vol. 6393 of LNCS, 2010, pp. 322–333. doi:[10.1007/978-3-642-16321-0\\_34](https://doi.org/10.1007/978-3-642-16321-0_34).

- [41] V. Mäkinen, G. Navarro, Position-restricted substring searching, in: Proc. LATIN, Vol. 3887 of LNCS, 2006, pp. 703–714. doi:10.1007/11682462\\_64.
- [42] P. G. S. da Fonseca, I. B. F. da Silva, Online construction of wavelet trees, in: 16th International Symposium on Experimental Algorithms, SEA 2017, June 21-23, 2017, London, UK, Vol. 75 of LIPICS, 2017, pp. 16:1–16:14. doi:10.4230/LIPIcs.SEA.2017.16.
- [43] J. Fischer, V. Heun, Space-efficient preprocessing schemes for range minimum queries on static arrays, SIAM J. Comput. 40 (2) (2011) 465–492. doi:10.1137/090779759.
- [44] J. I. Munro, Y. Nekrich, J. S. Vitter, Dynamic data structures for document collections and graphs, in: Proc. PODS, 2015, pp. 277–289. doi:10.1145/2745754.2745778.
- [45] M. Crochemore, W. Rytter, Jewels of Stringology: Text Algorithms, World Scientific, 2002.
- [46] E. Ohlebusch, S. Gog, A. Kügel, Computing matching statistics and maximal exact matches on compressed full-text indexes, in: Proc. SPIRE, Vol. 6393 of LNCS, 2010, pp. 347–358. doi:10.1007/978-3-642-16321-0\\_36.
- [47] K. Sadakane, Succinct representations of lcp information and improvements in the compressed suffix arrays, in: Proc. SODA, 2002, pp. 225–232.