

Computing the Parameterized Burrows–Wheeler Transform Online

Daiki Hashimoto¹, Diptarama Hendrian¹, Dominik Köppl², Ryo Yoshinaka¹, and Ayumi Shinohara¹

¹Tohoku University, Japan

²Tokyo Medical and Dental University, Japan

Abstract

Parameterized strings are a generalization of strings in that their characters are drawn from two different alphabets, where one is considered to be the alphabet of static characters and the other to be the alphabet of parameter characters. Two parameterized strings are a parameterized match if there is a bijection over all characters such that the bijection transforms one string to the other while keeping the static characters (i.e., it behaves as the identity on the static alphabet). Ganguly et al. [SODA 2017] proposed the parameterized Burrows–Wheeler transform (pBWT) as a variant of the Burrows–Wheeler transform for space-efficient parameterized pattern matching. In this paper, we propose an algorithm for computing the pBWT online by reading the characters of a given input string one-by-one from right to left. Our algorithm works in $O(|\Pi| \log n / \log \log n)$ amortized time for each input character, where n and Π denote the size of the input string and the alphabet of the parameter characters, respectively.

1 Introduction

The *parameterized matching problem* (*p-matching problem*) [2] is a generalization of the classic pattern matching problem in the sense that we here consider two disjoint alphabets, the set Σ of *static characters* and the set Π of *parameter characters*. We call a string over $\Sigma \cup \Pi$ a *parameterized string* (*p-string*). Two equal-length p-strings X and Y are said to *parameterized match* (*p-match*) if there is a bijection that renames the parameter characters in X so X becomes equal to Y . The *p-matching problem* is, given a text p-string T and pattern p-string P , to output the positions of all substrings of T that p-match P . The p-matching problem is motivated by applications in the software maintenance [1, 2], the plagiarism detection [5], the analysis of gene structures [17], and so on. There exist indexing structures that support p-matching, such as parameterized suffix trees [1, 17], parameterized suffix arrays [7, 10], and so on [4, 6, 13, 14]; see also [12] for a survey. A drawback of these indexing structures is that they have high space requirements.

A more space-efficient indexing structure, the *parameterized Burrows–Wheeler transform* (*pBWT*), was proposed by Ganguly et al. [8]. The pBWT is a variant of the *Burrows–Wheeler transform* (*BWT*) [3] that can be used as an indexing structure for p-matching using only $o(n \log n)$ bits of space. Later on, Kim and Cho [11] improved this indexing structure by changing the encoding of p-strings used for defining the pBWT. Recently, Ganguly et al. [9] augmented this index with capabilities of a suffix tree while keeping the space within $o(n \log n)$ bits. However, as far as we are aware of, none research related to the pBWT [8, 11, 9, 18] has discussed how to construct their pBWT-based data structures in detail. Their construction algorithms mainly rely on the parameterized suffix tree. Given the parameterized suffix tree of a p-string T of length n , the pBWT of T can be constructed in $O(n \log(|\Sigma| + |\Pi|))$ time offline.

In this paper, we propose an algorithm for constructing pBWTs and related data structures used for indexing structures of p-matching. Our algorithm constructs the data structures directly in an online manner by reading the input text from right to left. The algorithm uses the dynamic array

data structures of Navarro and Nekrich [15] to maintain our growing arrays. For each character read, our algorithm takes $O(|\Pi| \log n / \log \log n)$ amortized time, where n is the size of input string. Therefore, we can compute pBWT of a p-string T of length n in $O(n|\Pi| \log n / \log \log n)$ time in total. In comparison, computing the standard BWT on a string T (i.e., the pBWT on a string having no parameter characters) can be done in $O(n \log n / \log \log n)$ time with the dynamic array data structures [15] (see [16] for a description of this online algorithm). Looking at our time complexity, the factor $|\Pi|$ also appears in the time complexity of an offline construction algorithm of parameterized suffix arrays [7] as $O(n|\Pi|)$ and a right-to-left online construction algorithm of parameterized suffix trees [13] as $O(n|\Pi| \log(|\Pi| + |\Sigma|))$. This suggests it would be rather hard to improve the time complexity of the online construction of pBWT to be independent of $|\Pi|$.

2 Preliminaries

We denote the set of nonnegative integers by \mathbb{N} and let $\mathbb{N}_+ = \mathbb{N} \setminus \{0\}$ and $\mathbb{N}_\infty = \mathbb{N}_+ \cup \{\infty\}$. The set of strings over an alphabet A is denoted by A^* . The empty string is denoted by ε . The length of a string $W \in A^*$ is denoted by $|W|$. For a subset $B \subseteq A$, the set of elements of B occurring in $W \in A^*$ is denoted by $B \upharpoonright W$. We count the number of occurrences of characters of B in a string W by $|W|_B$. So, $|W|_A = |W|$. When B is a singleton of b , i.e., $B = \{b\}$, we often write $|W|_b$ instead of $|W|_{\{b\}}$. When W is written as $W = XYZ$, X , Y , and Z are called *prefix*, *factor*, and *suffix* of W , respectively. The i -th character of W is denoted by $W[i]$ for $1 \leq i \leq |W|$. The factor of W that begins at position i and ends at position j is $W[i : j]$ for $1 \leq i \leq j \leq |W|$. For convenience, we abbreviate $W[1 : i]$ to $W[: i]$ and $W[i : |W|]$ to $W[i :]$ for $1 \leq i \leq |W|$. Let $Rot(W, 0) = W$ and $Rot(W, i + 1) = Rot(W, i)[|W|]Rot(W, i)[: |W| - 1]$ be the i -th right rotation of W . Note that $Rot(W, i) = Rot(W, i + |W|)$. For convenience we denote $W_i = Rot(W, i)$. Let $Left_W(a)$ and $Right_W(a)$ be the leftmost and rightmost positions of a character $a \in A$ in W , respectively. If a does not occur in W , define $Left_W(a) = Right_W(a) = 0$.

2.1 Parameterized Burrows–Wheeler transform

Throughout this paper, we fix two disjoint ordered alphabets Σ and Π . We call elements of Σ *static characters* and those of Π *parameter characters*. Elements of Σ^* and $(\Sigma \cup \Pi)^*$ are called *static strings* and *parameterized strings* (or *p-strings* for short), respectively.

Two p-strings S and T of the same length are a *parameterized match* (*p-match*), denoted by $S \approx T$, if there is a bijection f on $\Sigma \cup \Pi$ such that $f(a) = a$ for any $a \in \Sigma$ and $f(S[i]) = T[i]$ for all $1 \leq i \leq |T|$ [2]. We use Kim and Cho’s version of p-string encoding [11], which replaces 0 in Baker’s encoding [1] by ∞ . The *prev-encoding* $\langle T \rangle$ of T is the string over $\Sigma \cup \mathbb{N}_\infty$ of length $|T|$ defined by

$$\langle T \rangle[i] = \begin{cases} T[i] & \text{if } T[i] \in \Sigma, \\ \infty & \text{if } T[i] \in \Pi \text{ and } Right_{T[:i-1]}(T[i]) = 0, \\ i - Right_{T[:i-1]}(T[i]) & \text{if } T[i] \in \Pi \text{ and } Right_{T[:i-1]}(T[i]) \neq 0 \end{cases}$$

for $1 \leq i \leq |T|$. When $T[i] \in \Pi$, $\langle T \rangle[i]$ represents the distance between i and the previous occurrence position of the same parameter character. If $T[i]$ does not occur before the position i , the distance is assumed to be ∞ . We call a string $W \in (\Sigma \cup \mathbb{N}_\infty)^*$ a *pv-string* if $W = \langle T \rangle$ for some p-string T . For any p-strings S and T , $S \approx T$ if and only if $\langle S \rangle = \langle T \rangle$ [2]. For example, given $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ and $\Pi = \{\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}\}$, $S = \mathbf{uvvauvb}$ and $T = \mathbf{xyyaxyb}$ are a p-match by f with $f(\mathbf{u}) = \mathbf{x}$ and $f(\mathbf{v}) = \mathbf{y}$, where $\langle S \rangle = \langle T \rangle = \infty\infty 1\mathbf{a}43\mathbf{b}$.

For defining pBWT, we use another encoding $\llbracket T \rrbracket$ given by

$$\llbracket T \rrbracket[i] = \begin{cases} T[i] & \text{if } T[i] \in \Sigma, \\ |\Pi \upharpoonright T_{n-i}[1 : Left_{T_{n-i}}(T[i])]| & \text{if } T[i] \in \Pi \end{cases}$$

Table 1: The pBWT $pBWT(T) = \mathbf{L}_T = \mathbf{a33131\$22aa}$ of the example string $T = \mathbf{xayzzazyza\$}$ with related arrays, where $\Sigma = \{\mathbf{a}\}$ and $\Pi = \{\mathbf{x, y, z}\}$.

i	T_i	$\langle T_i \rangle$	$RA_T[i]$	$LCP_T^\infty[i]$	$\langle T_{RA_T[i]} \rangle$	$F_T[i]$	$\llbracket T_{RA_T[i]} \rrbracket$	$L_T[i]$
1	$\$xayzzazyza$	$\$xaxxx1a252a$	1	0	$\$xaxxx1a252a$	$\$$	$\$3a211a233a$	\mathbf{a}
2	$\mathbf{a}xayzzazyz$	$\mathbf{a}xaxxx1a252$	2	0	$\mathbf{a}xaxxx1a252$	\mathbf{a}	$\mathbf{a}\$3a211a233$	$\mathbf{3}$
3	$\mathbf{za}xayzzazy$	$\mathbf{xa}xax61a25$	10	2	$\mathbf{axx1a252a}\$x$	\mathbf{a}	$\mathbf{a211a233a}\$3$	$\mathbf{3}$
4	$\mathbf{yza}xayzzaz$	$\mathbf{xxa}\$xa661a2$	6	0	$\mathbf{axx2a}\$xa661$	\mathbf{a}	$\mathbf{a233a}\$3a211$	$\mathbf{1}$
5	$\mathbf{zyza}xayzza$	$\mathbf{xx2a}\$xa661a$	3	1	$\mathbf{xa}\$xax61a25$	$\mathbf{3}$	$\mathbf{3a}\$3a211a23$	$\mathbf{3}$
6	$\mathbf{azyza}xayzz$	$\mathbf{axx2a}\$xa661$	7	1	$\mathbf{xa2x2a}\$xa66$	$\mathbf{1}$	$\mathbf{1a233a}\$3a21$	$\mathbf{1}$
7	$\mathbf{zazyza}xayz$	$\mathbf{xa2x2a}\$xa66$	11	1	$\mathbf{xaxx1a252a}\$$	$\mathbf{3}$	$\mathbf{3a211a233a}\$$	$\mathbf{\$}$
8	$\mathbf{zzazyza}xay$	$\mathbf{x1a2x2a}\$xa6$	8	1	$\mathbf{x1a2x2a}\$xa6$	$\mathbf{1}$	$\mathbf{11a233a}\$3a2$	$\mathbf{2}$
9	$\mathbf{yzzazyza}x\$a$	$\mathbf{xx1a252a}\$xa$	4	2	$\mathbf{xxa}\$xa661a2$	$\mathbf{3}$	$\mathbf{33a}\$3a211a2$	$\mathbf{2}$
10	$\mathbf{ayzzazyza}x\$$	$\mathbf{axx1a252a}\$x$	9	2	$\mathbf{xx1a252a}\$xa$	$\mathbf{2}$	$\mathbf{211a233a}\$3a$	\mathbf{a}
11	$\mathbf{xayzzazyza}\$$	$\mathbf{xaxxx1a252a}\$$	5	0	$\mathbf{xx2a}\$xa661a$	$\mathbf{2}$	$\mathbf{233a}\$3a211a$	\mathbf{a}

for $1 \leq i \leq |T|$. When $T[i] \in \Pi$, $\llbracket T \rrbracket[i]$ counts the number of distinct parameter characters in T between i and the next occurrence of $T[i]$, if $T[i]$ occurs after i . If i is the rightmost occurrence position of $T[i]$, then we continue counting parameter characters from the left end to the right until we find $T[i]$. Since $T[i]$ occurs in T_{n-i} as the last character, $\llbracket T \rrbracket[i]$ cannot be zero. Note that $\llbracket Rot(T, i) \rrbracket = Rot(\llbracket T \rrbracket, i)$ by definition. It is not hard to see that for any p-strings S and T , $S \approx T$ if and only if $\llbracket S \rrbracket = \llbracket T \rrbracket$ (see Proposition 1 in the appendix). For example, the two strings S and T given above are encoded as $\llbracket S \rrbracket = \llbracket T \rrbracket = 212a22b$.

Hereafter in this section, we fix a p-string T of length n which ends with a special static character $\$$ which occurs nowhere else in T . We extend the linear order over Σ to $\Sigma \cup \mathbb{N}_\infty$ by letting $\$ < a < i < \infty$ for any $a \in \Sigma \setminus \{\$\}$ and $i \in \mathbb{N}_+$. The order over \mathbb{N}_+ coincides with the usual numerical order.

The pBWT of T is defined through sorting $\llbracket T_p \rrbracket$ for $p = 1, \dots, n$ using $\langle T_p \rangle$ as keys.

Definition 1 (Parameterized rotation array). *The parameterized rotation array RA_T of T is an array of size n such that $RA_T[i] = p$ with $1 \leq p \leq n$ if and only if $\langle T_p \rangle$ is the i -th lexicographically smallest string in $\{\langle T_p \rangle \mid 1 \leq p \leq n\}$. We denote its inverse by RA_T^{-1} , i.e., $RA_T^{-1}[p] = i$ iff $RA_T[i] = p$.*

Note that RA_T and RA_T^{-1} are well-defined and bijective due to the presence of $\$$ in T . Here, we have $RA_T[i] = n - pSA_T[i] + 1$, where pSA_T refers to the suffix array \mathbf{pSA}_∞ in [11]. The array gives an $n \times n$ square matrix $(\llbracket T_{RA_T[i]} \rrbracket)_{i=1}^n$, which we call the *rotation sort matrix* of T , whose (i, p) entry is $\llbracket T_{RA_T[i]} \rrbracket[p]$. The pBWT of T is formed by the characters in the last column of the matrix.

Definition 2 (pBWT [11]). *The parameterized Burrows–Wheeler transform (pBWT) of a p-string T , denoted by $pBWT(T)$, is a string of length n such that $pBWT(T)[i] = \llbracket T_{RA_T[i]} \rrbracket[n]$.*

An example pBWT can be found in Table 1. We will use \mathbf{L}_T as a synonym of $pBWT(T)$, since it represents the *last* column of the matrix $(\llbracket T_{RA_T[i]} \rrbracket)_{i=1}^n$. When picking up the characters from the *first* column, we obtain another array F_T . That is, $F_T[i] = \llbracket T_{RA_T[i]} \rrbracket[1]$ for all $i \in \{1, \dots, n\}$. Those arrays \mathbf{L}_T and F_T are “linked” by the following mapping.

Definition 3 (LF mapping). *The LF mapping $LF_T : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ for T is defined as $LF_T(i) = j$ if $T_{RA_T[i+1]} = T_{RA_T[j]}$.*

By rotating T_p to the right by one, the last character moves to the first position in T_{p+1} . Roughly speaking, $\mathbf{L}_T[i]$ and $F_T[LF_T(i)]$ “originate” in the same character occurrence of T , which implies $\mathbf{L}_T[i] = F_T[LF_T(i)]$ in particular. One can recover $\llbracket T \rrbracket$ as $\llbracket T \rrbracket[p] = \mathbf{L}_T[LF_T^{-p}(k_T)]$ for $1 \leq p \leq n$ where $k_T = RA_T^{-1}[n]$. \mathbf{L}_T , F_T , and LF_T are used for pattern matching based on pBWT. See [11] for the details.

Our pBWT construction algorithm maintains neither RA_T nor LF_T , but involves some helper data structures in addition to \mathbf{L}_T and F_T . Among those, the array LCP_T^∞ is worth explaining before going

into the algorithmic details. For two pv-strings X and Y , let $lcp^\infty(X, Y) = |W|_\infty$ be the number of ∞ 's in the longest common prefix W of X and Y . The following array counts the number of ∞ 's in the longest common prefixes of two adjacent rows in $(\langle T \rangle_{RA_T[i]})_{i=1}^n$.

Definition 4 (∞ -LCP array). The ∞ -LCP array LCP_T^∞ of T is an array of size n such that $\text{LCP}_T^\infty[n] = 0$ and $\text{LCP}_T^\infty[i] = lcp^\infty(\langle T_{RA_T[i]} \rangle, \langle T_{RA_T[i+1]} \rangle)$ for $1 \leq i < n$.

Table 1 shows an example of a pBWT and related (conceptual) data structures. We can compute $lcp^\infty(\langle T_{RA_T[i]} \rangle, \langle T_{RA_T[j]} \rangle)$ using LCP_T^∞ as follows.

Lemma 1. For $1 \leq i < j \leq n$, $lcp^\infty(\langle T_{RA_T[i]} \rangle, \langle T_{RA_T[j]} \rangle) = \min_{i \leq k < j} \text{LCP}_T^\infty[k]$.

Kim and Cho [11] showed some basic relations among L_T , LF_T , and lcp^∞ . We rephrase Lemma 3 of [11] into a form convenient for our discussions.

Lemma 2. Consider i and j with $1 \leq i < j \leq n$ and $T_{RA_T[i]}[n], T_{RA_T[j]}[n] \in \Pi$. Then, $\text{LF}_T(i) < \text{LF}_T(j)$ iff $\min\{\text{L}_T[i] - 1, lcp^\infty(\langle T_{RA_T[i]} \rangle, \langle T_{RA_T[j]} \rangle)\} < \text{L}_T[j]$.

Corollary 1 ([11]). If $i < j$ and $\text{L}_T[i] = \text{L}_T[j]$, then $\text{LF}_T(i) < \text{LF}_T(j)$.

To maintain L_T , F_T , and LCP_T^∞ dynamically, our algorithm uses the data structure for dynamic arrays by Navarro and Nekrich [15] that supports the following operations on an array Q of size m in $O(\frac{\log m}{\log \log m})$ amortized time.

1. $\text{access}(Q, i)$: returns $Q[i]$ for $1 \leq i \leq m$;
2. $\text{rank}_a(Q, i)$: returns $|Q[:i]|_a$ for $1 \leq i \leq m$;
3. $\text{select}_a(Q, i)$: returns i -th occurrence position of a for $1 \leq i \leq \text{rank}_a(Q, m)$;
4. $\text{insert}_a(Q, i)$: inserts a between $Q[i-1]$ and $Q[i]$ for $1 \leq i \leq m+1$;
5. $\text{delete}(Q, i)$: deletes $Q[i]$ from Q for $1 \leq i \leq m$.

Corollary 1 implies that we can compute $\text{LF}_T(i)$ and its inverse $\text{LF}_T^{-1}(j)$ by

$$\begin{aligned} \text{LF}_T(i) &= \text{select}_x(\text{F}_T, \text{rank}_x(\text{L}_T, i)) \quad \text{where } x = \text{L}_T[i], \\ \text{LF}_T^{-1}(j) &= \text{select}_y(\text{L}_T, \text{rank}_y(\text{F}_T, j)) \quad \text{where } y = \text{F}_T[j]. \end{aligned}$$

3 Computing pBWT online

This section introduces our algorithm computing $p\text{BWT}_T$ in an online manner by reading a p-string T from right to left. Let $T = cS$ for $c \in \Sigma \cup \Pi \setminus \{\$\}$ and $n = |S| \geq 1$. We consider updating L_S to L_T . Hereafter, we assume that Σ is known and $|\Sigma| \leq |T|$ as in [16]. Among the rows of the rotation matrices of S and T , the rows of $\llbracket S \rrbracket = \llbracket S_n \rrbracket$ and $\llbracket T \rrbracket = \llbracket T_{n+1} \rrbracket$ play important roles when updating. Let $k_S = RA_S^{-1}[n]$ and $k_T = RA_T^{-1}[n+1]$. We note that $\text{L}_S[k_S] = \text{L}_T[k_T] = \$$.

First, we observe RA_T is obtained from RA_S just by ‘‘inserting’’ $n+1$ at k_T .

Lemma 3. For $1 \leq i \leq n+1$,

$$RA_T[i] = \begin{cases} RA_S[i] & \text{if } i < k_T, \\ n+1 & \text{if } i = k_T, \\ RA_S[i-1] & \text{if } i > k_T. \end{cases}$$

Algorithm 1: PBWT update algorithm

```

1 Function UpdateAll( $c, n, L, F, \text{Right}, \text{Left}, \text{RM}, C, \text{LCP}^\infty$ )
2    $k = \text{select}_\$(L, 1);$  // =  $k_S$ 
3    $L, F, \text{Right}, \text{Left}, \text{RM} = \text{UpdateLF}(c, n, L, F, \text{Right}, \text{Left}, \text{RM}, k);$ 
4      $// = L_T^\circ, F_T^\circ, \text{Right}_T, \text{Left}_T, \text{RM}_T^\circ$ 
5    $L, F, C, k' = \text{InsertRow}(n, L, F, C, k);$  // =  $L_T, F_T, C_T, k_T$ 
6   foreach  $a \in \Pi$  do
7     if  $\text{RM}[a] \geq k'$  then  $\text{RM}[a] = \text{RM}[a] + 1;$  // =  $\text{RM}_T[a]$ 
8      $x = \text{UpdateLCP}(L, F, \text{LCP}^\infty, k');$  // =  $\text{LCP}_T^\infty[k_T]$ 
9      $\text{LCP}^\infty[k' - 1] = \text{UpdateLCP}(L, F, \text{LCP}^\infty, k' - 1);$  // =  $\text{LCP}_T^\infty[k_T - 1]$ 
10     $\text{insert}_x(\text{LCP}^\infty, k');$ 
11  return  $n + 1, L, F, \text{Right}, \text{Left}, \text{RM}, C, \text{LCP}^\infty;$ 

```

Table 2: An example of our update step for $S = \text{xayzzazyza}\$$ and $T = \text{y}S$. The updated and inserted values are highlighted. In the arrays $\langle T_{RAS}[i] \rangle$, updated/inserted values appear only after $\$$. Lemmas 3 and 8 are immediate consequences of this observation.

$F_S[i]$	$\langle S_{RAS}[i] \rangle$	$L_S[i]$	$F_T^\circ[i]$	$\langle T_{RAS}[i] \rangle$	$L_T^\circ[i]$	$F_T[i]$	$\langle T_{RAT}[i] \rangle$	$L_T[i]$
\$	\$\text{xxxx1a252a}	a	\$	\$\text{xxa3x1a252a}	a	\$	\$\text{xxa3x1a252a}	a
a	a\$\text{xxxx1a252}	3	a	a\$\text{xxa3x1a252}	3	a	a\$\text{xxa3x1a252}	3
a	a\$\text{xx1a252a}\\$	3	a	a\$\text{xx1a252a}\\$4x	3	a	a\$\text{xx1a252a}\\$4x	3
a	a\$\text{xx2a}\\$xa661	1	a	a\$\text{xx2a}\\$4xa361	1	a	a\$\text{xx2a}\\$4xa361	1
3	3\$\text{axxx61a25}	3	3	3\$\text{axxx61a25}	2	3	3\$\text{axxx61a25}	2
1	1\$\text{ax2x2a}\\$xa66	1	1	1\$\text{ax2x2a}\\$4xa36	1	1	1\$\text{ax2x2a}\\$4xa36	1
3	3\$\text{axxx1a252a}\\$	\$	3	3\$\text{axxx1a252a}\\$4	2	3	3\$\text{axxx1a252a}\\$4	2
1	1\$\text{x1a2x2a}\\$xa6	2	1	1\$\text{x1a2x2a}\\$4xa3	2	1	1\$\text{x1a2x2a}\\$4xa3	2
3	3\$\text{xxa}\\$xa661a2	2	2	2\$\text{xxa}\\$4xa361a2	2	2	2\$\text{xxa}\\$4xa361a2	2
2	2\$\text{xxx1a252a}\\$xa	a	2	2\$\text{xxx1a252a}\\$4xa	a	2	2\$\text{xxx3x1a252a}\\$	\$
2	2\$\text{xx2a}\\$xa661a	a	2	2\$\text{xx2a}\\$4xa361a	a	2	2\$\text{xx2a}\\$4xa361a	a

In the BWT, where S and T have no parameter characters, this implies that $L_T[i] = T_{RAT}[i][n+1] = S_{RAS}[i][n] = L_S[i]$ for $i < k_T$ and $L_T[i+1] = T_{RAT}[i+1][n+1] = S_{RAS}[i][n] = L_S[i]$ for $i > k_T$, except when $i = k_S$. Therefore, for computing L_T from L_S , we only need to update $L_S[k_S] = \$$ to c and to find the position $k_T = RA_T^{-1}[n+1]$ where $\$$ should be inserted. However in the pBWT, $RA_T[i] = RA_S[i]$ does not necessarily imply that the values $L_T[i] = \llbracket T_{RAT}[i] \rrbracket[n+1]$ and $L_S[i] = \llbracket S_{RAS}[i] \rrbracket[n]$ coincide, since it is not always true that $\llbracket S \rrbracket = \llbracket T \rrbracket[2 :]$. So we also need to update the values of the encoding.

Algorithm 1 shows our update procedure, which maintains the array F and other auxiliary data structures in addition to L . After getting the key position k_S as the unique occurrence position of $\$$ in L_S at Line 2, to update the values of L and F from L_S and F_S to L_T and F_T , respectively, we compute intermediate arrays L_T° and F_T° of length n , which satisfy

$$L_T^\circ[i] = \llbracket T_{RAS}[i] \rrbracket[n+1] \quad \text{and} \quad F_T^\circ[i] = \llbracket T_{RAS}[i] \rrbracket[1]$$

for $1 \leq i \leq n$ using `UpdateLF` at Line 3. In other words, L_T° and F_T° are extracted from the last and the first columns of the $n \times (n+1)$ matrix $(\llbracket T_{RAS}[i] \rrbracket)_{i=1}^n$, respectively, which can conceptionally be obtained by deleting the k_T -th row of the rotation sort matrix of T . We then find the other key position k_T and inserts appropriate values into L_T° and F_T° at k_T to turn them into L_T and F_T , respectively, by `InsertRow` at Line 4. The rest of the algorithm is devoted to maintaining some of the helper arrays. Particularly, a dedicated function `UpdateLCP` is used to update the ∞ -LCP array. In the remainder of this section, we will explain those functions and involved auxiliary data structures in respective subsections. Table 2 shows an example of our 2-step update.

Algorithm 2: Computing L_T° and F_T°

```

1 Function UpdateLF( $c, n, L, F, \text{Right}, \text{Left}, \text{RM}, k$ )
2   if  $c \in \Sigma$  then  $L[k] = c$ ;
3   else
4     foreach  $a \in \Pi$  with  $\text{Left}[a] \neq 0$  do
5       // Computing  $L_T^\circ[\text{RM}[a]] = F_T^\circ[\text{LF}_S(\text{RM}[a])]$ 
6        $i = \text{RM}[a]$ ;
7        $j = \text{select}_{L[i]}(F, \text{rank}_{L[i]}(L, i));$  //  $j = \text{LF}_S(i)$ 
8       if  $a = c$  then
9          $\text{cnt} = 0$ ;
10        foreach  $b \in \Pi$  with  $\text{Left}[b] \neq 0$  do
11          if  $\text{Left}[a] \geq \text{Right}[b]$  then  $\text{cnt} = \text{cnt} + 1$ ;
12        else
13           $\text{cnt} = L[i]$ ;
14          if  $\text{Left}[c] = 0$  or  $\text{Left}[a] > \text{Left}[c] \geq \text{Right}[c] > \text{Right}[a]$  then
15             $\text{cnt} = \text{cnt} + 1$ ;
16           $L[i] = \text{cnt}; F[j] = \text{cnt};$ 
17          // Computing  $L_T^\circ[k_S] = \llbracket T \rrbracket[1]$ 
18           $\text{cnt} = 1$ ;
19          if  $\text{Left}[c] = 0$  then
20            foreach  $a \in \Pi$  with  $\text{Left}[a] \neq 0$  do  $\text{cnt} = \text{cnt} + 1$ ;
21             $\text{Right}[c] = n + 1; \text{Left}[c] = n + 1; \text{RM}[c] = k$ ;
22          else
23            foreach  $a \in \Pi$  with  $\text{Left}[a] > \text{Left}[c]$  do  $\text{cnt} = \text{cnt} + 1$ ;
24             $\text{Left}[c] = n + 1$ ;
25           $L[k] = \text{cnt};$ 
26   return  $L, F, \text{Right}, \text{Left}, \text{RM}$ ;

```

3.1 Step 1: UpdateLF computes $L_T^\circ[i]$ and $F_T^\circ[i]$

When $c \in \Sigma$, computing L_T° and F_T° from L_S and F_S , respectively, is easy.

Lemma 4. *If $c \in \Sigma$, then for any $i \in \{1, \dots, n\}$, $F_T^\circ[i] = F_S[i]$ and $L_T^\circ[i] = L_S[i]$ except for $L_T^\circ[k_S] = c$.*

Concerning the case $c \in \Pi$, first let us express the values of $\llbracket T \rrbracket$ using $\llbracket S \rrbracket$.

Lemma 5. *Suppose $c \in \Pi$.*

$$\llbracket T \rrbracket[1] = \begin{cases} |\Pi \upharpoonright S| + 1 & \text{if } \text{Left}_S(c) = 0, \\ |\Pi \upharpoonright S[1 : \text{Left}_S(c)]| & \text{otherwise.} \end{cases}$$

For $1 \leq p \leq n$, if $S[p] \in \Sigma$ or $p \neq \text{Right}_S(S[p])$, then $\llbracket T \rrbracket[p + 1] = \llbracket S \rrbracket[p]$. If $S[p] = a \in \Pi$ and $p = \text{Right}_S(a)$, then

$$\llbracket T \rrbracket[p + 1] = \begin{cases} |\Pi \upharpoonright S[p + 1 : n]| + 1 & \text{if } a = c, \\ \llbracket S \rrbracket[p] + 1 & \text{if } \text{Left}_S(c) = 0 \text{ or} \\ & \text{Left}_S(a) < \text{Left}_S(c) \leq \text{Right}_S(c) < \text{Right}_S(a), \\ \llbracket S \rrbracket[p] & \text{otherwise.} \end{cases}$$

Based on Lemmas 4 and 5, Algorithm 2 computes $F_T^\circ[i]$ and $L_T^\circ[i]$ from $F_S[i]$ and $L_S[i]$, as well as other auxiliary data structures. Note that, since the intermediate matrix $(T_{RAS[i]})_{i=1}^n$ misses a

row corresponding to $\llbracket T \rrbracket$, the value $\llbracket T \rrbracket[1]$ does not matter for F_T° , whereas it appears as $L_T^\circ[k_S] = L_T^\circ[RA_S^{-1}[n]]$. When $c \in \Pi$, Lemma 5 implies that, other than $L_T^\circ[k_S] = \llbracket T \rrbracket[1]$, we only need to update the values at the positions in L and F corresponding to the rightmost occurrence position $p = \text{Right}_S(a)$ of each parameter character $a \in \Pi$ in S . By rotating S to the right by $n - p$, that occurrence comes to the right end and appears in the pBWT. That is, the array L needs to be updated only at i such that $RA_S[i] = n - \text{Right}_S(a)$. The algorithm maintains such position i as $RM_S[a]$ for each $a \in \Pi \upharpoonright S$, i.e. $RM_S[a] = RA_S^{-1}[n - \text{Right}_S(a)]$. Similarly, we only need to update F at $LF_S(RM_S[a])$, where $F_T^\circ[LF_S(RM_S[a])] = L_T^\circ[RM_S[a]]$. In our algorithm, as alternatives of $Left_S$ and $Right_S$, we maintain two arrays **Left** and **Right** that store the leftmost and rightmost occurrence positions of parameter characters counting *from the right end*, respectively, i.e., $\text{Left}_S[a] = \text{Right}_{\bar{S}}(a)$ and $\text{Right}_S[a] = \text{Left}_{\bar{S}}(a)$ for each $a \in \Pi$, where \bar{S} is the reverse of S .

Algorithm 2 also updates RM to RM_T° , which indicates the row of L_T° corresponding to the rightmost occurrence of each parameter character in T . That is, $RM_T^\circ[a] = i$ iff $RA_S[i] = \text{Right}_T[a]$, as long as a occurs in T . When $c \in \Pi$ and it appears in the text for the first time, we have $RM_T^\circ[c] = k_S$ (Line 19). Other than that, $RM_T^\circ[a] = RM_S[a]$ for every $a \in \Pi$.

Lemma 6. *Algorithm 2 computes $L_T^\circ[i]$, $F_T^\circ[i]$, Right_T , Left_T , and RM_T° in $O(|\Pi| \frac{\log n}{\log \log n})$ amortized time.*

3.2 Step 2: InsertRow computes L_T and F_T

To transform F_T° and L_T° into F_T and L_T , we insert the values $\llbracket T \rrbracket[1]$ and $\llbracket T \rrbracket[n + 1]$ at the position k_T , respectively. We know those values as $\llbracket T \rrbracket[1] = L_T^\circ[k_S]$ and $\llbracket T \rrbracket[n + 1] = \$$. Therefore, it is enough to discuss how to find the position k_T .

In the case $c \in \Sigma$, the position k_T can be calculated similarly to the case of BWT for static strings thanks to Corollary 1. Define $\Sigma_{<b} = |\{a \in \Sigma \mid a < b\}|$.

Lemma 7. *If $c \in \Sigma$, $k_T = |T|_{\Sigma_{<c}} + |\{i \mid L_T^\circ[i] = c, 1 \leq i \leq k_S\}|$.*

In the case $c \in \Pi$, we will use Lemma 2 for finding k_T in Lemma 9 below. We first observe that one can use LCP_S^∞ to calculate $lcp^\infty(\langle T_p \rangle, \langle T_q \rangle)$ for most cases.

Lemma 8. *For $1 \leq p < q \leq n$, $lcp^\infty(\langle T_p \rangle, \langle T_q \rangle) = lcp^\infty(\langle S_p \rangle, \langle S_q \rangle)$.*

Lemma 9. *Suppose $c \in \Pi$. Let $\ell_i = lcp^\infty(\langle S_{RA_S[i]} \rangle, \langle S_{RA_S[k_S]} \rangle)$ for $1 \leq i \leq n$. Then,*

$$k_T = 1 + |T|_\Sigma + |\{i \mid 1 \leq L_T^\circ[i] \leq L_T^\circ[k_S], 1 \leq i < k_S\}| \quad (1)$$

$$+ |\{i \mid \ell_i < L_T^\circ[k_S] < L_T^\circ[i], 1 \leq i < k_S\}| \quad (2)$$

$$+ |\{i \mid 1 \leq L_T^\circ[i] \leq \min\{L_T^\circ[k_S] - 1, \ell_i\}, k_S < i \leq n\}|. \quad (3)$$

Proof. By definition,

$$k_T = 1 + |T|_\Sigma + |\{j \mid F_T[j] \in \mathbb{N}_+, 1 \leq j < k_T\}|,$$

of which we focus on the last term. Let $h = LF_T^{-1}(k_T)$ and $m_i = lcp^\infty(\langle T_{RA_T[i]} \rangle, \langle T_{RA_T[h]} \rangle)$ for $1 \leq i \leq n + 1$. By Lemma 2, $F_T[j] \in \mathbb{N}_+$ and $1 \leq j < k_T$ iff for $i = LF_T^{-1}(j)$, either

1. $1 \leq i < h$ and $1 \leq L_T[i] \leq L_T[h]$,
2. $1 \leq i < h$ and $m_i < L_T[h] < L_T[i]$, or
3. $h < i \leq n + 1$ and $1 \leq L_T[i] \leq \min\{L_T[h] - 1, m_i\}$.

Those three cases are mutually exclusive. Let $m_i^\circ = lcp^\infty(\langle T_{RA_S[i]} \rangle, \langle T_{RA_S[k_S]} \rangle)$. Counting each of the above cases is equivalent to counting i such that

Algorithm 3: Inserting $\llbracket T \rrbracket[1]$ to F and $\llbracket T \rrbracket[n+1]$ to L

```

1 Function InsertRow( $n, L, F, C, k$ )
2    $x = L[k]$ ; // =  $\llbracket T \rrbracket[1]$ 
3   if  $x \in \Sigma$  then
4      $k' = \text{select}_x(C, 1) - |\Sigma_{<x}| - 1 + \text{rank}_x(L, k)$ ; //  $k_T = |S|_{\Sigma_{<c}} + |\{i \mid L_T^\circ[i] = c, 1 \leq i \leq k_S\}|$ 
5      $\text{insert}_x(C, \text{select}_x(C, 1))$ ;
6   else
7      $k' = 1 + |C| - |\Sigma|$ ; // =  $1 + |S|_\Sigma$ 
8     for  $y = 1$  to  $x$  do  $k' = k' + \text{rank}_y(L, k - 1)$ ; // Term (1)
9      $j = 0$ ;
10    for  $y = 0$  to  $x - 1$  do
11      if  $\text{rank}_y(\text{LCP}^\infty, k - 1) \neq 0$  then
12         $j = \max\{j, \text{select}_y(\text{LCP}^\infty, \text{rank}_y(\text{LCP}^\infty, k - 1))\}$ ;
13        //  $j = \max(\{j\} \cup \{i \mid \text{LCP}^\infty[i] = y \text{ and } 1 \leq i < k_S\})$ 
14      for  $y = x + 1$  to  $|\Pi|$  do  $k' = k' + \text{rank}_y(L, j)$ ; // Term (2)
15       $j = n$ ; //  $j_0 = n$ 
16      for  $y = 1$  to  $x - 1$  do // Term (3)
17        if  $\text{rank}_{y-1}(\text{LCP}^\infty, k - 1) < \text{rank}_{y-1}(\text{LCP}^\infty, n)$  then
18           $j = \min\{j, \text{select}_{y-1}(\text{LCP}^\infty, \text{rank}_{y-1}(\text{LCP}^\infty, k - 1) + 1)\}$ ;
19          //  $j_y = \min(\{j_{y-1}\} \cup \{i \mid \text{LCP}^\infty[i] = y - 1 \text{ and } k_S \leq i \leq n\})$ 
20           $k' = k' + \text{rank}_y(L, j - 1) - \text{rank}_y(L, k)$ ;
21     $\text{insert}_\$(L, k')$ ;  $\text{insert}_x(F, k')$ ;
22    return  $L, F, C, k'$ ;

```

1. $1 \leq i < k_S$ and $1 \leq L_T^\circ[i] \leq L_T^\circ[k_S]$,
2. $1 \leq i < k_S$ and $m_i^\circ < L_T^\circ[k_S] < L_T^\circ[i]$, or
3. $k_S < i \leq n$ and $1 \leq L_T^\circ[i] \leq \min\{L_T^\circ[k_S] - 1, m_i^\circ\}$.

This is because the matrix $(\llbracket T_{RAS[i]} \rrbracket)_{i=1}^n$ can conceptionally be obtained by removing the k_T -th row of the matrix of $(\llbracket T_{RAT[i]} \rrbracket)_{i=1}^{n+1}$, where the row k_S of $(\llbracket T_{RAS[i]} \rrbracket)_{i=1}^n$ corresponds to the row h of $(\llbracket T_{RAT[i]} \rrbracket)_{i=1}^{n+1}$ in particular ($RAS[k_S] = RAT[h] = n$), and $i = k_T = RAT^{-1}[n+1]$ is not counted due to $T[n+1] = \$ \in \Sigma$.

Lemma 8 implies $m_i^\circ = \ell_i$, which completes the proof. \square

Based on Lemmas 7 and 9, Algorithm 3 finds the key position k_T .

For handling the case $c \in \Sigma$, we maintain a dynamic array C by which one can obtain the value $|T|_{\Sigma_{<c}} = |S|_{\Sigma_{<c}}$ quickly. The array C_S can be seen as a string of the form $C_S = a_1^{|S|_{a_1}+1} \dots a_\sigma^{|S|_{a_\sigma}+1}$, where a_1, \dots, a_σ enumerate the static characters of Σ in the lexicographic order ($\sigma = |\Sigma|$) and a^s denotes the sequence of a of length s . Then, $|T|_{\Sigma_{<c}} = \text{select}_c(C_S, 1) - |\Sigma_{<c}| - 1$. The other term $|\{i \mid L_T^\circ[i] = c, 1 \leq i \leq k_S\}|$ in Lemma 7 is calculated as $\text{rank}_c(L_T^\circ, k_S)$. We remark C_S has $a^{|S|_{a+1}}$ rather than $a^{|S|_a}$ so that $\text{select}_c(C, 1)$ is always defined.

Suppose $c \in \Pi$. The term $|T|_\Sigma$ of the equation of Lemma 9 is calculated as $|T|_\Sigma = |S|_\Sigma = |C| - |\Sigma|$. Let $x = L_T^\circ[k_S]$. Term (1) is obtained at Line 8 by

$$(1) = \sum_{y=1}^x \text{rank}_y(L_T^\circ, k_S - 1).$$

Concerning Term (2), we first find the range of $i < k_S$ satisfying $\ell_i < x$. By Lemma 1, $\ell_i = \min_{i \leq j < k_S} \text{LCP}_S^\infty[j]$. Thus, for any $i < k_S$, $\ell_i < x$ iff $i \leq j_* = \max\{j \mid \text{LCP}_S^\infty[j] < x, j < k_S\}$.

The **for** loop of Line 10 computes such j_* . Then, (2) is computed at Line 13 as

$$(2) = |\{i \mid x < L_T^\circ[i], 1 \leq i \leq j_*\}| = \sum_{y=x+1}^{|\Pi|} \text{rank}_y(L_T^\circ, j_*).$$

We compute Term (3) by summing up the numbers of positions $i > k_S$ such that $L_T^\circ[i] = y \leq \ell_i$ for all $y = 1, \dots, x-1$ in the **for** loop of Line 15. To this end, we find the range of $i > k_S$ such that $\ell_i \geq y$. By Lemma 1, $\ell_i = \min_{k_S \leq j < i} \text{LCP}_S^\infty[j]$. Thus, for every $i > k_S$, $\ell_i \geq y$ iff $i < j_y = \min\{j \mid \text{LCP}_S^\infty[j] < y, k_S \leq j \leq n\}$. Note that $j_y = \min(\{j_{y-1}\} \cup \{j \mid \text{LCP}_S^\infty[j] = y-1, k_S \leq j \leq n\})$ for any $y \geq 1$ assuming $j_0 = n$. Line 17 computes j_y as the first occurrence of $y-1$ after those in $\text{LCP}^\infty[1 : k_S - 1]$. Then, (3) is calculated by

$$(3) = |\{i \mid 1 \leq L_T^\circ[i] \leq x-1, k_S < i < j_y\}| \\ = \sum_{y=1}^{x-1} (\text{rank}_y(L_T^\circ, j_y - 1) - \text{rank}_y(L_T^\circ, k_S)).$$

Lemma 10. *Algorithm 3 computes L_T, F_T, C_T , and k_T in $O(|\Pi| \frac{\log n}{\log \log n})$ amortized time.*

3.3 Step 3: Updating LCP^∞ by UpdateLCP

What remains to do is updating the arrays RM and LCP^∞ . On the one hand, updating RM from RM_T° to RM_T is easy. $\text{RM}_T^\circ[a]$ should be incremented by one just if $\text{RM}_T^\circ[a] \geq k_T$. Otherwise, $\text{RM}_T[a] = \text{RM}_T^\circ[a]$. On the other hand, Lemma 8 implies LCP_T^∞ is almost identical to LCP_S^∞ .

Corollary 2. $\text{LCP}_T^\infty[i] = \text{LCP}_S^\infty[i]$ if $i < k_T - 1$, and $\text{LCP}_T^\infty[i] = \text{LCP}_S^\infty[i-1]$ if $i > k_T$.

By Corollary 2, we only need to compute $\text{LCP}_T^\infty[k_T - 1]$ and $\text{LCP}_T^\infty[k_T]$, to which Lemma 8 cannot directly be applied. The following lemma allows us to reduce the calculation of $\text{LCP}_T^\infty[k] = \text{lcp}^\infty(\langle T_{RA_T[k]} \rangle, \langle T_{RA_T[k+1]} \rangle)$ to that of $\text{lcp}^\infty(\langle T_{RA_T[LF_T^{-1}(k)]} \rangle, \langle T_{RA_T[LF_T^{-1}(k+1)]} \rangle)$, to which Lemma 8 may be applied.

Lemma 11. *Let $1 \leq i, j \leq n+1$, $p = RA_T[i]$, $q = RA_T[j]$, $\ell = \text{lcp}^\infty(\langle T_p \rangle, \langle T_q \rangle)$, $i' = LF_T^{-1}(i)$, $j' = LF_T^{-1}(j)$, $p' = RA_T[i']$, $q' = RA_T[j']$, and $\ell' = \text{lcp}^\infty(\langle T_{p'} \rangle, \langle T_{q'} \rangle)$.*

1. If $F_T[i] = F_T[j] \in \Sigma$, then $\ell = \ell'$.
2. If $F_T[i] \neq F_T[j]$ and either $F_T[i] \in \Sigma$ or $F_T[j] \in \Sigma$, then $\ell = 0$.
3. If $F_T[i], F_T[j] \in \mathbb{N}_+$, then

$$\ell = \begin{cases} \ell' + 1 & \text{if } \ell' < \min\{F_T[i], F_T[j]\}, \\ \ell' & \text{if } \ell' \geq F_T[i] = F_T[j], \\ \min\{F_T[i], F_T[j]\} & \text{otherwise.} \end{cases}$$

Proof. Let $T_p = aU$, $T_q = bV$, $T_{p'} = Ua$ and $T_{q'} = Vb$.

1. If $a = b \in \Sigma$, then $\ell = \ell' = \text{lcp}^\infty(\langle U \rangle, \langle V \rangle)$.

2. In the case $a \neq b$ and $\{a, b\} \cap \Sigma \neq \emptyset$, clearly $\langle T_p \rangle[1] \neq \langle T_q \rangle[1]$. Thus $\ell = 0$.

3. In the case $a, b \in \Pi$, let W be the longest common prefix of $\langle T_{p'} \rangle$ and $\langle T_{q'} \rangle$, u and v be the first occurrence positions of a in $T_{p'}$ and b in $T_{q'}$, respectively, and w be the ℓ' -th occurrence position of ∞ in W .

Suppose $\ell' < \min\{F_T[i], F_T[j]\} = \min\{L_T[i'], L_T[j']\}$. That is, $|W|_\infty < \min\{|\langle T_{p'} \rangle[:u]|_\infty, |\langle T_{q'} \rangle[:v]|_\infty\}$. This means $|W| < \min\{u, v\}$ and thus ∞W is the longest common prefix of $\langle T_p \rangle$ and $\langle T_q \rangle$. Thus, we have $\ell = |\infty W|_\infty = \ell' + 1$.

Suppose $\ell' \geq F_T[i] = F_T[j]$, i.e., $|W|_\infty \geq |\langle T_{p'} \rangle[:u]|_\infty = |\langle T_{q'} \rangle[:v]|_\infty$. Then $|W[:u]|_\infty = |W[:v]|_\infty$ and $W[u] = W[v] = \infty$ implies $u = v$. Let Z be the longest common prefix of $\langle T_p \rangle$ and $\langle T_q \rangle$. Then, W and Z can be written as $W = X\infty Y$ and $Z = \infty X u Y$, where $|X| = u - 1$. Therefore, $\ell = \ell'$.

Algorithm 4: Updating $\text{LCP}^\infty[i]$

```
1 Function UpdateLCP(L, F,  $\text{LCP}^\infty, i$ )
2    $j = i + 1; x = 0;$ 
3   if  $F[i] = F[j]$  or  $F[i], F[j] \in \mathbb{N}_+$  then
4      $i' = \text{select}_{F[i]}(L, \text{rank}_{F[i]}(F, i));$  //  $i' = LF_T^{-1}(i)$ 
5      $j' = \text{select}_{F[j]}(L, \text{rank}_{F[j]}(F, j));$  //  $j' = LF_T^{-1}(i + 1)$ 
6     for  $y = |\Pi|$  downto 0 do
7       if  $\text{rank}_y(\text{LCP}^\infty, i' - 1) \neq \text{rank}_y(\text{LCP}^\infty, j' - 1)$  then  $x = y;$ 
8         //  $x = \text{lcp}^\infty(\langle S_{RAS[i']}, \langle S_{RAS[j']}\rangle)$ 
9       if  $F[i], F[j] \in \mathbb{N}_+$  then
10        if  $x < \min\{F[i], F[j]\}$  then  $x = x + 1;$ 
11        else if  $F[i] \neq F[j]$  then  $x = \min\{F[i], F[j]\};$ 
12   return  $x;$ 
```

Otherwise, $F_T[i] \neq F_T[j]$ and $\ell' \geq \min\{F_T[i], F_T[j]\}$. Assume $F_T[i] < F_T[j]$ (the case $F_T[j] < F_T[i]$ is symmetric). Then $u \leq |W|$. Moreover, we have $u < v$, since otherwise, $\langle T_q \rangle[:v]$ had to be a prefix of $\langle T_{p'} \rangle[:u]$, which is impossible by $F_T[i] < F_T[j]$. Let $\langle T_p[:|W| + 1] \rangle = \alpha X u Y$, where $|X| = u - 1$. Then we have $\langle T_q[:|W| + 1] \rangle = \alpha X \alpha Y'$ for some $Y' \in (\Sigma \cup \mathbb{N}_\infty)^*$. Thus $\ell = |\alpha X|_\infty = F_T[i]$. \square

One can compute $\ell' = \text{lcp}^\infty(\langle T_{p'} \rangle, \langle T_{q'} \rangle)$ in Lemma 11 for $1 \leq p' < q' \leq n$ using Lemmas 8 and 1 as

$$\begin{aligned} \text{lcp}^\infty(\langle T_{p'} \rangle, \langle T_{q'} \rangle) &= \text{lcp}^\infty(\langle S_{p'} \rangle, \langle S_{q'} \rangle) = \min\{ \text{LCP}_S^\infty[h] \mid i' \leq h < j' \} \\ &= \min(\{0\} \cup \{y \mid \text{rank}_y(\text{LCP}_S^\infty, i' - 1) \neq \text{rank}_y(\text{LCP}_S^\infty, j' - 1)\}). \end{aligned}$$

Finally, when $q' = n + 1$, we have $F_T[j] = \$ \neq F_T[i]$, and thus $\text{lcp}^\infty(\langle T_p \rangle, \langle T_q \rangle) = 0$. Algorithm 4 computes $\text{LCP}_T^\infty[i]$ using F_T , L_T , and LCP_S^∞ .

Lemma 12. *Algorithm 4 computes $\text{LCP}_T^\infty[i]$ in $O(|\Pi| \frac{\log n}{\log \log n})$ amortized time.*

By Lemmas 6, 10, and 12, we have the following theorem.

Theorem 1. *Given $c \in \Sigma \cup \Pi$, $n = |S|$, $L = L_S$, $F = F_S$, $\text{Right} = \text{Right}_S$, $\text{Left} = \text{Left}_S$, $\text{RM} = \text{RM}_S$, $C = C_S$, and $\text{LCP}^\infty = \text{LCP}_S^\infty$ for some $S \in (\Sigma \cup \Pi)^*$, Algorithm 1 computes $|T|$, L_T , F_T , Right_T , Left_T , RM_T , C_T , and LCP_T^∞ for $T = cS$ in $O(|\Pi| \frac{\log n}{\log \log n})$ amortized time per input character.*

Corollary 3. *For a p -string T of length n , $p\text{BWT}_T$ can be computed in an online manner by reading T from right to left in $O(n|\Pi| \frac{\log n}{\log \log n})$ time.*

References

- [1] Brenda S. Baker. A theory of parameterized pattern matching: algorithms and applications. In *Proceedings of the twenty-fifth annual ACM symposium on Theory of Computing (STOC 1993)*, pages 71–80, 1993.
- [2] Brenda S. Baker. Parameterized pattern matching: Algorithms and applications. *Journal of Computer and System Sciences*, 52(1):28–42, 1996.
- [3] Michael Burrows and David Wheeler. A block-sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994.
- [4] Diptarama, Takashi Katsura, Yuhei Otomo, Kazuyuki Narisawa, and Ayumi Shinohara. Position heaps for parameterized strings. In *Proceedings of the 28th Annual Symposium on Combinatorial Pattern Matching (CPM 2017)*, pages 8:1–8:13, 2017.

- [5] Kimmo Fredriksson and Maxim Mozgovoy. Efficient parameterized string matching. *Information Processing Letters*, 100(3):91 – 96, 2006.
- [6] Noriki Fujisato, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Right-to-left online construction of parameterized position heaps. In *Proceedings of the Prague Stringology Conference 2018 (PSC 2018)*, pages 91–102, 2018.
- [7] Noriki Fujisato, Yuto Nakashima, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Direct linear time construction of parameterized suffix and LCP arrays for constant alphabets. In *Proceedings of the 26th International Symposium on String Processing and Information Retrieval (SPIRE 2019)*, pages 382–391, 2019.
- [8] Arnab Ganguly, Rahul Shah, and Sharma V. Thankachan. pBWT: Achieving succinct data structures for parameterized pattern matching and related problems. In *Proceedings of the 28th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 2017)*, pages 397–407, 2017.
- [9] Arnab Ganguly, Rahul Shah, and Sharma V. Thankachan. Fully functional parameterized suffix trees in compact space. In *Proceedings of the 49th International Colloquium on Automata, Languages, and Programming (ICALP 2022)*, pages 65:1–65:18, 2022.
- [10] Tomohiro I, Satoshi Deguchi, Hideo Bannai, Shunsuke Inenaga, and Masayuki Takeda. Lightweight parameterized suffix array construction. In *Proceedings of the 20th International Workshop on Combinatorial Algorithms (IWOCA 2009)*, pages 312–323, 2009.
- [11] Sung-Hwan Kim and Hwan-Gue Cho. Simpler FM-index for parameterized string matching. *Information Processing Letters*, 165:106026, 2021.
- [12] Juan Mendivelso, Sharma V. Thankachan, and Yoan J. Pinzón. A brief history of parameterized matching problems. *Discrete Applied Mathematics*, 274:103–115, 2020.
- [13] Katsuhito Nakashima, Noriki Fujisato, Diptarama Hendrian, Yuto Nakashima, Ryo Yoshinaka, Shunsuke Inenaga, Hideo Bannai, Ayumi Shinohara, and Masayuki Takeda. DAWGs for parameterized matching: Online construction and related indexing structures. In *Proceedings of the 31st Annual Symposium on Combinatorial Pattern Matching (CPM 2020)*, pages 26:1–26:14, 2020.
- [14] Katsuhito Nakashima, Diptarama Hendrian, Ryo Yoshinaka, and Ayumi Shinohara. An Extension of Linear-size Suffix Tries for Parameterized Strings. In *SOFSEM 2020 Student Research Forum*, pages 97–108, 2020.
- [15] Gonzalo Navarro and Yakov Nekrich. Optimal dynamic sequence representations. *SIAM Journal on Computing*, 43(5):1781–1806, 2014.
- [16] Alberto Policriti and Nicola Prezza. Fast online Lempel–Ziv factorization in compressed space. In *Proceedings of the the 22nd International Symposium on String Processing and Information Retrieval (SPIRE 2015)*, pages 13–20, 2015.
- [17] Tetsuo Shibuya. Generalization of a suffix tree for RNA structural pattern matching. *Algorithmica*, 39(1):1–19, 2004.
- [18] Sharma V. Thankachan. Compact text indexing for advanced pattern matching problems: Parameterized, order-isomorphic, 2D, etc. (invited talk). In *Proceedings of the 33rd Annual Symposium on Combinatorial Pattern Matching (CPM 2022)*, pages 3:1–3:3, 2022.

A Proofs

Proposition 1. For any p -strings S and T , $S \approx T$ if and only if $\llbracket S \rrbracket = \llbracket T \rrbracket$.

Proof. For simplicity, assume that S and T contain no static character. Suppose $S \approx T$. Since $S[i] = S[j]$ iff $T[i] = T[j]$ for any indices i, j , we have

$$\llbracket S \rrbracket[i] = |\Pi \upharpoonright S_{n-i}[1 : \text{Left}_{S_{n-i}}(S[i])]| = |\Pi \upharpoonright T_{n-i}[1 : \text{Left}_{T_{n-i}}(T[i])]| = \llbracket T \rrbracket[i].$$

for all i .

Suppose $S \not\approx T$. Let i be the leftmost position such that $\langle S \rangle[i] \neq \langle T \rangle[i]$. We may assume without loss of generality that $\langle S \rangle[i] < \langle T \rangle[i]$. Let $j = i - \langle S \rangle[i]$. Then,

$$\llbracket S \rrbracket[j] = |\Pi \upharpoonright S[j+1 : i]| = |\Pi \upharpoonright S[j : i-1]| = |\Pi \upharpoonright T[j : i-1]|,$$

since $S[j : i-1] \approx T[j : i-1]$. The fact $S[j] \notin \Pi \upharpoonright S[j+1 : i-1]$ implies $T[j] \notin \Pi \upharpoonright T[j+1 : i-1]$. Moreover, $\langle S \rangle[i] < \langle T \rangle[i]$ implies $T[i] \notin \Pi \upharpoonright T[j : i-1]$. Hence,

$$\llbracket T \rrbracket[j] \geq |(\Pi \upharpoonright T[j : i-1]) \cup \{T[i]\}| > |\Pi \upharpoonright T[j : i-1]| = \llbracket S \rrbracket[j].$$

Lemma 3 is a corollary to the following lemma.

Lemma 13. For any i and j such that $1 \leq i < j \leq n$, $RA_S^{-1}[i] < RA_S^{-1}[j]$ iff $RA_T^{-1}[i] < RA_T^{-1}[j]$.

Proof. Let $S_i = U\$V$ and $S_j = XY\$Z$, where $|U\$| = |X| = i < j \leq n$. We have $T_i = U\$cV$ and $T_j = XY\$cZ$. Since $\$$ does not occur in X , $\langle U\$ \rangle \neq \langle X \rangle$. Thus,

$$RA_S^{-1}[i] < RA_S^{-1}[j] \iff \langle U\$ \rangle < \langle X \rangle \iff RA_T^{-1}[i] < RA_T^{-1}[j].$$

Lemma 4. If $c \in \Sigma$, then for any $i \in \{1, \dots, n\}$, $F_T^\circ[i] = F_S[i]$ and $L_T^\circ[i] = L_S[i]$ except for $L_T^\circ[k_S] = c$.

Proof. If $c \in \Sigma$, then $\llbracket S \rrbracket = \llbracket T \rrbracket[2 :]$ by definition. So, for any $i \in \{1, \dots, n\}$,

$$F_T^\circ[i] = \llbracket T_{RA_S[i]} \rrbracket[1] = \llbracket S_{RA_S[i]} \rrbracket[1] = F_S[i]$$

and

$$L_T^\circ[i] = \llbracket T_{RA_S[i]} \rrbracket[n+1] = \begin{cases} \llbracket S_{RA_S[i]} \rrbracket[n] = L_S[i] & \text{if } RA_S[i] \neq n, \\ c & \text{if } RA_S[i] = n. \end{cases} \quad \square$$

Lemma 5. Suppose $c \in \Pi$.

$$\llbracket T \rrbracket[1] = \begin{cases} |\Pi \upharpoonright S| + 1 & \text{if } \text{Left}_S(c) = 0, \\ |\Pi \upharpoonright S[1 : \text{Left}_S(c)]| & \text{otherwise.} \end{cases}$$

For $1 \leq p \leq n$, if $S[p] \in \Sigma$ or $p \neq \text{Right}_S(S[p])$, then $\llbracket T \rrbracket[p+1] = \llbracket S \rrbracket[p]$. If $S[p] = a \in \Pi$ and $p = \text{Right}_S(a)$, then

$$\llbracket T \rrbracket[p+1] = \begin{cases} |\Pi \upharpoonright S[p+1 : n]| + 1 & \text{if } a = c, \\ \llbracket S \rrbracket[p] + 1 & \text{if } \text{Left}_S(c) = 0 \text{ or} \\ & \text{Left}_S(a) < \text{Left}_S(c) \leq \text{Right}_S(c) < \text{Right}_S(a), \\ \llbracket S \rrbracket[p] & \text{otherwise.} \end{cases}$$

Proof. The claim on value of $\llbracket T \rrbracket[1]$ is clear by definition. For $p \geq 1$, if $T[p+1] = S[p] \in \Sigma$, then $\llbracket T \rrbracket[p+1] = \llbracket S \rrbracket[p]$.

Let us consider the case $S[p] = a \in \Pi$. If $p \neq \text{Right}_S(a)$, then a occurs somewhere after p in S . Let $q > p$ be the first occurrence position of a after p in S . By definition, $\llbracket S \rrbracket[p] = |\Pi \upharpoonright S[p+1 : q]| = |\Pi \upharpoonright T[p+2 : q+1]| = \llbracket T \rrbracket[p+1]$.

Suppose $p = \text{Right}_S(a)$ for some $a \in \Pi$. If $a = c$, since $T[1] = c$ and $c \notin \Pi \upharpoonright S[p+1 : n]$, we have $\llbracket T \rrbracket[p+1] = |\Pi \upharpoonright S[p+1 : n] \cup \{c\}| = |\Pi \upharpoonright S[p+1 : n]| + 1$. If $\text{Right}_S(c) = 0$ or $\text{Left}_S(a) < \text{Left}_S(c) \leq \text{Right}_S(c) < \text{Right}_S(a) = p$, $\llbracket S \rrbracket[p]$ counts the number of distinct p-characters in $S[p+1 :]S[: \text{Left}_S(a)]$, where c does not occur. On the other hand, $\llbracket T \rrbracket[p+1]$ counts the ones in $S[p+1 :]cS[: \text{Left}_S(a)]$. That is, $\llbracket T \rrbracket[p+1] = \llbracket S \rrbracket[p] + 1$. Otherwise, if $\text{Left}_S(c) < \text{Left}_S(a)$ or $\text{Right}_S(a) < \text{Right}_S(c)$, we already have c in $S[p+1 :]S[: \text{Left}_S(a)]$. Thus $\llbracket T \rrbracket[p+1] = \llbracket S \rrbracket[p]$. \square

Lemma 7. If $c \in \Sigma$, $k_T = |T|_{\Sigma_{<c}} + |\{i \mid \mathbf{L}_T^\circ[i] = c, 1 \leq i \leq k_S\}|$.

Proof. By definition, $k_T = |T|_{\Sigma_{<c}} + |\{j \mid \mathbf{F}_T[j] = c, 1 \leq j \leq k_T\}|$. By Corollary 1 and the bijectivity of LF_T , the second term equals

$$|\{i \mid \mathbf{L}_T[i] = c, 1 \leq i \leq LF_T^{-1}(k_T)\}|$$

and further more equals

$$|\{i \mid \mathbf{L}_T^\circ[i] = c, 1 \leq i \leq k_S\}|$$

because \mathbf{L}_T and \mathbf{L}_T° are different only in that \mathbf{L}_T has an extra element $\$ < c$, and the position $LF_T^{-1}(k_T)$ in \mathbf{L}_T corresponds to the position k_S in \mathbf{L}_T° . \square

Lemma 8. For $1 \leq p < q \leq n$, $lcp^\infty(\langle T_p \rangle, \langle T_q \rangle) = lcp^\infty(\langle S_p \rangle, \langle S_q \rangle)$.

Proof. Let $S_p = U\$V$ and $S_q = XY\$Z$, where $|U\$| = |X| = p < q = |XY\$|$. Then, $T_p = U\$cV$ and $T_q = XY\$cZ$. Since $\$$ does not appear in X ,

$$lcp^\infty(\langle S_p \rangle, \langle S_q \rangle) = lcp^\infty(U\$, X) = lcp^\infty(\langle T_p \rangle, \langle T_q \rangle).$$