



**HAL**  
open science

# Linking BWT and XBW via Aho-Corasick Automaton: Applications to Run-Length Encoding

Bastien Cazaux, Eric Rivals

► **To cite this version:**

Bastien Cazaux, Eric Rivals. Linking BWT and XBW via Aho-Corasick Automaton: Applications to Run-Length Encoding. CPM 2019 - 30th Annual Symposium on Combinatorial Pattern Matching, University of Pisa, Jun 2019, Pise, Italy. pp.24:1–24:20, 10.4230/LIPIcs.CPM.2019.24. lirmm-02382066

**HAL Id: lirmm-02382066**

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-02382066v1>

Submitted on 27 Nov 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Linking BWT and XBW via Aho-Corasick Automaton: Applications to Run-Length Encoding

**Bastien Cazaux** 

Department of Computer Science, University of Helsinki, Finland  
L.I.R.M.M., CNRS, Université Montpellier, France  
<https://www.mv.helsinki.fi/home/cazaux/>  
bastien.cazaux@helsinki.fi

**Eric Rivals** 

L.I.R.M.M., CNRS, Université Montpellier, France  
<https://www.lirmm.fr/~rivals/>  
rivals@lirmm.fr

---

## Abstract

The boom of genomic sequencing makes compression of sets of sequences inescapable. This underlies the need for multi-string indexing data structures that helps compressing the data. The most prominent example of such data structures is the Burrows-Wheeler Transform (BWT), a reversible permutation of a text that improves its compressibility. A similar data structure, the eXtended Burrows-Wheeler Transform (XBW), is able to index a tree labelled with alphabet symbols. A link between a multi-string BWT and the Aho-Corasick automaton has already been found and led to a way to build a XBW from a multi-string BWT. We exhibit a stronger link between a multi-string BWT and a XBW by using the order of the concatenation in the multi-string. This bijective link has several applications: first, it allows one to build one data structure from the other; second, it enables one to compute an ordering of the input strings that optimises a Run-Length measure (i.e., the compressibility) of the BWT or of the XBW.

**2012 ACM Subject Classification** Mathematics of computing → Discrete mathematics; Theory of computation → Randomness, geometry and discrete structures; Theory of computation → Data structures and algorithms for data management

**Keywords and phrases** Data Structure, Algorithm, Aho-Corasick Tree, compression, RLE

**Digital Object Identifier** 10.4230/LIPIcs.CPM.2019.24

**Related Version** A full version of the paper is available at <https://arxiv.org/abs/1805.10070>.

**Supplement Material** Source code at: <https://framagit.org/bcazaux/compressbwt>

**Funding** Support from the Institut de Biologie Computationnelle (ANR-11-BINF-0002).

*Eric Rivals*: thanks the GEM Flagship project funded from Labex NUMEV (ANR-10-LABX-0020).

## 1 Introduction

A seminal, key data structure, which was used for searching a set of words in a text, is the Aho-Corasick (AC) automaton [1]. Its states form a tree that indexes all the prefixes of the words, and each node in the tree is equipped with another kind of arc, called a Failure Link. A failure link of a node/prefix  $v$  points to the node representing the largest proper suffix of  $v$  in the tree. In a way, the Aho-Corasick automaton can be viewed as a multi-string indexing data structure.

In the early 90's, the Burrows-Wheeler Transform (BWT) of a text  $T$ , which is a reversible permutation of  $T$ , was introduced for the sake of compressing a text. Indeed, the BWT permutation tends to group identical symbols in runs, which favours compression [5]. However, the BWT can also be used as an index for searching in  $T$ , using the Backward Search



© Bastien Cazaux and Eric Rivals;  
licensed under Creative Commons License CC-BY  
30th Annual Symposium on Combinatorial Pattern Matching (CPM 2019).

Editors: Nadia Pisanti and Solon P. Pissis; Article No. 24; pp. 24:1–24:20

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

procedure [10]. In fact, the BWT of  $T$  is the last column of a matrix containing all cyclic-shifts of  $T$  sorted in lexicographical order. As sorting the cyclic shifts of  $T$  is equivalent to sorting its suffixes, there exists a natural link between the suffix array of  $T$  and the BWT of  $T$ . Starting in 2005, the radical increase in textual data and in biological sequences raised the need for multi-string indexes. In the multi-string case, a simple approach is to concatenate all input strings and separate them using the same termination symbol (which does not belong to the alphabet), and then to index the result with a traditional indexing data structure (*e.g.*, a suffix array or a FM-index). Another approach is to add a unique termination symbol for each string without concatenation [3, 21]. In both approaches, an initial order on the strings is given by the permutation of the strings in the first approach, or by the order between the termination symbols in the second. All our results are made with the first approach (which minimises the alphabet size) in mind, and can be adapted to the second. Such multi-string indexes are heavily exploited in bioinformatics: first, to index all chromosomes of a genome [17], or a large collection of similar genomes, which allows aligning sequence reads simultaneously to several reference genomes [19], or second, to store and mine whole sets of raw DNA/RNA sequence reads for the purpose of comparing biological conditions or of identifying splice junctions in RNA [7, 14]. In fact, managing compressed and searchable read data sets is now crucial for bioinformatic analyses.

Initially viewed as a simple extension of single-string BWT construction, the efficient construction of multi-string BWT is not trivial and has been investigated *per se*. Bauer *et al.* proposed, among others, a lightweight incremental algorithm for their construction [3]. Then, Holt *et al.* devised algorithms for directly merging several, already built multi-string BWTs efficiently [14, 15], which has recently been improved to build simultaneously the companion *Longest Common Prefix* (LCP) table [8], or to scale up to terabyte datasets [23].

The notion of BWT has been extended into the XBW to index trees whose arcs are labelled by alphabet symbols [9]. The XBW comprises two arrays, which compactly represent the tree and offer navigational operations.

Recently, Gagie *et al.* propose the notion of Wheeler graphs to subsume several variants of the BWT, including the XBW of a trie for a set of strings [12]. The relation between a multi-string BWT and the XBW representation of the Aho-Corasick automaton has already been studied and exploited. Hon *et al.* first use the XBW representation of the Aho-Corasick trie to speed-up dictionary matching [16] building up on [4]. Manzini gave an algorithm that computes the failure links for the trie using the Suffix Array and LCP tables, and an algorithm to build the XBW of the trie with failure links from the BWT of a multi-string [22]. However, none of these established a bijective link between a multi-string BWT and the XBW of the Aho-Corasick automaton. To generalise these results, one needs to consider the order in which the strings are concatenated to form the multi-string. This idea enables us to exhibit a bijection between a multi-string BWT and the XBW of the Aho-Corasick automaton, which allows building one structure from the other in either direction (from BWT to XBW or from XBW to BWT). Finally, we exploit this bijection to find an optimal string order that minimises a Run-Length Encoding (*i.e.*, maximises the compressibility) of these two data structures. Numerous proofs are in appendix (A preprint version of this work is available as [6]).

## 2 Preliminaries

In this section, we introduce basic concepts and notation, present the Burrows-Wheeler Transform and an extension of it for a set of strings. Finally, we provide a notation for the Aho-Corasick automaton and the eXtended Burrows-Wheeler Transform.

## 2.1 General notation

**Set and Permutation.** Let  $i$  and  $j$  be two integers such that  $i \leq j$ . The *interval*  $\llbracket i, j \rrbracket$  is the set of all integers between  $i$  and  $j$ . An *integer interval partition* of  $\llbracket i, j \rrbracket$  is a set of intervals  $\{\llbracket i_1, j_1 \rrbracket, \dots, \llbracket i_n, j_n \rrbracket\}$  such that  $i_1 = i$ ,  $j_n = j$ , and for all  $k \in \llbracket 1, n-1 \rrbracket$ ,  $j_k + 1 = i_{k+1}$ . We also define the order  $<$  on intervals such that for two intervals  $u := \llbracket i, j \rrbracket$  and  $v := \llbracket i', j' \rrbracket$ ,  $u < v$  iff  $j < i'$ . Let  $E$  be a finite set and let  $\#E$  (or  $\#(E)$ ) denote its cardinality. A *permutation* of  $E$  is an automorphism of  $E$ . A permutation  $\sigma$  of  $E$  is said to be *circular* iff for all  $i$  and  $j \in E$ , there exists a positive integer  $k$  such that  $\sigma^k(i) = j$  (where  $\sigma^1(i) = \sigma(i)$  and  $\sigma^k(i) = \sigma(\sigma^{k-1}(i))$ ). For a circular permutation  $\sigma$  of  $E$  and an element  $e \in E$ , we denote by  $\overline{\sigma}_e$  the function from  $\llbracket 1, \#E \rrbracket$  to  $E$  such that for all  $i \in \llbracket 1, \#E \rrbracket$ ,  $\overline{\sigma}_e(i) := \sigma^{(i+\#E-1)}(e)$ . Given a total order  $<$  for  $E$ , we define  $<_\sigma$  as the order on  $E$  such that  $e <_\sigma f$  iff  $\sigma(e) < \sigma(f)$  for any  $e, f \in E$ .

**String.** Let  $\Sigma$  be a finite alphabet. A *string*  $w$  of length  $n$  over  $\Sigma$  is a sequence of symbols  $w[1] \dots w[n]$  where  $w[i] \in \Sigma$  for all  $i \in \llbracket 1, n \rrbracket$ .  $\Sigma^*$  is the set of all finite strings over  $\Sigma$ . The *length* of a string  $w$  is denoted by  $|w|$ . A *substring* of  $w$  is written as  $w[i, j] := w[i] \dots w[j]$  for some  $i \leq j$ . A *prefix* of  $w$  is a substring of  $w$  starting at position 1, and a *suffix* of  $w$  is a substring of  $w$  ending at position  $|w|$ . A prefix  $x$  (or a suffix) of a string  $y$  is said *proper* if  $x$  is different from  $y$ . The *reverse* of a string  $w$ , denoted by  $\overleftarrow{w}$ , is the string  $w[n] \dots w[2]w[1]$ . We define the *lexicographic order*  $<$  on strings as usual.

**Ordered Set of Strings.** Let  $S = \{s_1, \dots, s_n\}$  be a set of strings. The norm of  $S$ , denoted  $\|S\|$ , is the sum of the length of strings of  $S$ , i.e.  $\|S\| := \sum_{s_i \in S} |s_i|$ . Let  $\text{Prefix}(S)$ , (respectively  $\text{Suffix}(S)$ ) denote the set of all prefixes (resp. all suffixes) of strings of  $S$ . We denote by  $\overleftarrow{S}$  the set of all reverse strings of strings of  $S$ , i.e.  $\overleftarrow{S} := \{\overleftarrow{s_1}, \dots, \overleftarrow{s_n}\}$ . An *ordered set of strings*  $P$  is a pair  $(S, \sigma)$  where  $S$  is a set of strings in lexicographic order, and  $\sigma$  a circular permutation of  $S$ . We denote by  $P.S$  the set of strings  $S$  and by  $P.\sigma$  the circular permutation  $\sigma$ . We denote by  $\overleftarrow{P}$  the pair  $(\overleftarrow{P.S}, P.\sigma)$ .

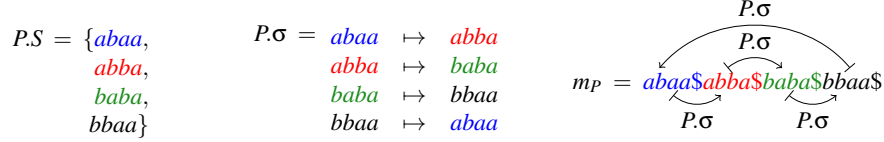
## 2.2 Burrows-Wheeler Transform

**BWT of a string.** Let  $w$  be a string and  $i$  be an integer satisfying  $1 \leq i \leq |w|$ . The *Suffix Array* (SA) of  $w$  [20], denoted  $\text{SA}(w)$ , is the array of integers that stores the starting positions of the  $|w|$  suffixes of  $w$  sorted in lexicographic order. The *Burrows-Wheeler Transform* (BWT) [5] of  $w$ , denoted  $\text{BWT}(w)$ , is the array containing a permutation of the symbols of  $w$  which satisfies  $\text{BWT}(w)[i] := w[\text{SA}(w)[i] - 1]$  if  $\text{SA}(w)[i] > 1$ , and  $\text{BWT}(w)[i] := w[|w|]$  otherwise. The *Longest Common Prefix* table (LCP) [5] of  $w$ , denoted  $\text{LCP}(w)$ , is the array of integers such that  $\text{LCP}(w)[i]$  equals the length of the longest common prefix between the suffixes of  $w$  starting at positions  $\text{SA}(w)[i]$  and  $\text{SA}(w)[i-1]$  if  $i > 1$ , and 0 otherwise.

For any string  $s \in \Sigma^*$  and any symbol  $c \in \Sigma$ , one defines the functions denoted **rank** and **select** as follows:  $\text{rank}_c(s, i)$  is the number of occurrences of  $c$  in  $s[1, i]$ , and  $\text{select}_c(s, j)$  is the position of the  $j^{\text{th}}$  occurrence of  $c$  in  $s$ . The arrays  $\text{BWT}(w)$  and  $\text{LCP}(w)$  can be computed in  $O(|w|)$  time [5, 10]. Simultaneously, one can compute **rank** and **select** for  $\text{BWT}(w)$  at no additional cost and implement them such that any **rank** or **select** query takes constant time [13, 11, 18]. We use such state-of-the-art data structure to store a BWT.

Let  $C$  denote the array of length  $\#\Sigma$  such that  $C[c]$  equals the number of symbols of  $w$  that are alphabetically strictly smaller than  $c$ . The *Last-to-First column mapping* (LF) [10] of  $w$  is the function such that for any  $1 \leq i \leq |w|$  one has  $\text{LF}(w)[i] := C[\text{BWT}(w)[i]] + \text{rank}_{\text{BWT}(w)[i]}(\text{BWT}(w), i)$ . It is proved that  $\text{SA}(w)[\text{LF}(w)[i]] = \text{SA}(w)[i] - 1$  if  $\text{SA}(w)[i] \geq 2$  and  $\text{SA}(w)[\text{LF}(w)[i]] = |w|$  otherwise [5, 10].

**BWT of a set of strings.** From now on, let  $P := (S, \sigma)$  be an ordered set of strings. We assume the symbol  $\$$  is not in  $\Sigma$  and is alphabetically smaller than all other symbols. We denote by  $m_P$  the string obtained by concatenating the strings of  $P.S := \{s_1, \dots, s_{\#(P.S)}\}$  separated by a  $\$$  and following the order  $\overline{P.\sigma}$ . i.e.,  $m_P := s_{\overline{P.\sigma}(1)}\$s_{\overline{P.\sigma}(2)}\$ \dots \$s_{\overline{P.\sigma}(n)}\$$  (See Figure 1). We extend the notion of BWT of a string to an ordered set of strings  $P$ : the BWT of  $P$  is the BWT of the string  $m_P$ , i.e.,  $\text{BWT}(P) := \text{BWT}(m_P)$ . We extend similarly the LF function by setting that  $\text{LF}(P) := \text{LF}(m_P)$  (by using **rank** on  $\text{BWT}(P)$ ).



■ **Figure 1** Running example with  $P.S := \{abaa, abba, baba, bbaa\}$  and the corresponding  $m_P$ .

We define the *Longest Representative Suffix* table (LRS) of  $P$  as the array of  $|m_P|$  integers satisfying: for any  $i \in \llbracket 1, |m_P| \rrbracket$  one has  $\text{LRS}(P)[i] := \text{select}_{\$}(m_P[\text{SA}(m_P)[i] : |m_P|], 1) - 1$ . The entry  $\text{LRS}(P)[i]$  gives the length of the substring of  $m_P$  starting at position  $\text{SA}(m_P)[i]$  up to, but not including the next  $\$$ . Using the LRS table, we extend the notion of LCP table to an ordered set of strings. For  $P$ , we set  $\text{LCP}(P)[i] := \min(\text{LCP}(m_P)[i], \text{LRS}(P)[i])$  for any  $i \in \llbracket 1, |m_P| \rrbracket$ . Using tables  $\text{BWT}(P)$  and  $\text{LCP}(m_P)$ , we can compute the tables  $\text{LRS}(P)$  and  $\text{LCP}(P)$  in linear time in  $|m_P|$  (see Appendix). Figure 2 illustrates the LCP and LRS tables.

### 2.3 Aho-Corasick automaton and eXtended Burrows-Wheeler Transform

**Tree.** Let  $\mathcal{T}$  be a tree and  $u$  be a node of  $\mathcal{T}$ . Let  $\perp$  denote the root of  $\mathcal{T}$ . We denote by  $\text{Parent}_{\mathcal{T}}(u)$  the parent of  $u$  in  $\mathcal{T}$ , by  $\text{Children}_{\mathcal{T}}(u)$  the set of children of  $u$  in  $\mathcal{T}$ , and by  $\text{Leaves}_{\mathcal{T}}(u)$  the set of leaves in the subtree of  $u$  in  $\mathcal{T}$ . Let  $v$  be a leaf of  $\mathcal{T}$ ; we denote by  $\mathcal{B}_{\mathcal{T}}(v)$  the subtree of  $\mathcal{T}$  containing all nodes comprised between  $\perp$  and  $v$  included. As for a leaf  $v$  in the subtree of  $u$  in  $\mathcal{T}$ ,  $\#\text{Children}_{\mathcal{B}_{\mathcal{T}}(v)}(u) = 1$ , we denote by  $\text{Child}_{\mathcal{T}(v)}(u)$  the unique element of  $\text{Children}_{\mathcal{B}_{\mathcal{T}}(v)}(u)$ . Given a total order  $\prec$  on  $\text{Leaves}_{\mathcal{T}}(\perp)$ , we extend this order on the set  $\text{Children}_{\mathcal{T}}(v)$  for any node  $v$  of  $\mathcal{T}$  by: for any  $x, y$  in  $\text{Children}_{\mathcal{T}}(v)$ ,  $x \prec y$  iff  $\min_{x' \in \text{Leaves}_{\mathcal{T}}(x)} x' \prec \min_{y' \in \text{Leaves}_{\mathcal{T}}(y)} y'$ .

**Aho-Corasick tree.** The *Aho-Corasick automaton* (AC) [1] of a set of strings  $S$  is a digraph whose set of nodes is the set of all prefixes of the strings of  $S$ . This graph is composed of two trees on the same node set. The first tree, which we called the *Aho-Corasick Tree* (ACT), has an arc from a prefix  $u$  to a different prefix  $v$  iff  $u$  is the longest proper prefix of  $v$  in  $\text{Prefix}(S)$  (see Figure 4). The second tree, termed *Aho-Corasick Failure link* (ACFL), has an arc from a prefix  $u$  to a different prefix  $v$  iff  $v$  is the longest proper **suffix** of  $u$  in  $\text{Prefix}(S)$ .

**eXtended Burrows-Wheeler Transform.** Let  $\mathcal{T}$  be an ordered tree (i.e., with a total order on the set of children of each node) such that every node of  $\mathcal{T}$  is labelled with a symbol from an alphabet  $\Sigma$ . We define the functions  $\delta$  and  $\pi$  on the set of nodes of  $\mathcal{T}$  except for the root such that for a node  $v$  of  $\mathcal{T}$ ,  $\delta(v)$  is the label of  $v$ , and  $\pi(v)$  is the string obtained by concatenating the labels from  $v$ 's parent up to the root of  $\mathcal{T}$ . Let  $\prec$  be the total order between the nodes of  $\mathcal{T}$  such that for  $u$  and  $v$  two nodes of  $\mathcal{T}$ ,  $u \prec v$  iff  $\pi(u)$  is strictly lexicographically smaller than  $\pi(v)$  or  $u$  is before  $v$  in the order of  $\mathcal{T}$ .

LCP(P)	LRS(P)	SA(m <sub>P</sub> )	BWT(P)
0	0	20	\$ a b a a \$ a b b a \$ b a b a \$ b b a a
0	0	5	\$ a b b a \$ b a b a \$ b b a a \$ a b a a
0	0	10	\$ b a b a \$ b b a a \$ a b a a \$ a b b a \$ b a b a
0	0	15	\$ b b a a \$ a b a a \$ a b b a \$ b a b a \$ b b a a
0	1	19	a \$ a b a a \$ a b b a \$ b a b a \$ b b a a
1	1	4	a \$ a b b a \$ b a b a \$ b b a a \$ a b a a
1	1	9	a \$ b a b a \$ b b a a \$ a b a a \$ a b b a
1	1	14	a \$ b b a a \$ a b a a \$ a b b a \$ b a b a
1	2	18	aa \$ a b a a \$ a b b a \$ b a b a \$ b b a a
2	2	3	aa \$ a b b a \$ b a b a \$ b b a a \$ a b a a
1	3	12	a b a a \$ b b a a \$ a b a a \$ a b b a \$ b a b a
3	4	1	a b a a \$ a b b a \$ b a b a \$ b b a a \$
2	4	6	a b b a \$ b a b a \$ b b a a \$ a b a a \$ a b a a
0	2	8	b a \$ b a b a \$ b b a a \$ a b a a \$ a b a a
2	2	13	b a \$ b b a a \$ a b a a \$ a b b a \$ b a b a
2	3	17	b a a \$ a b a a \$ a b b a \$ b a b a \$ b a b a
3	3	2	b a a \$ a b b a \$ b a b a \$ b b a a \$ a
2	4	11	b a b a \$ b b a a \$ a b a a \$ a b b a \$
1	3	7	b b a \$ b a b a \$ b b a a \$ a b a a \$ a
3	4	16	b b a a \$ a b a a \$ a b b a \$ b a b a \$

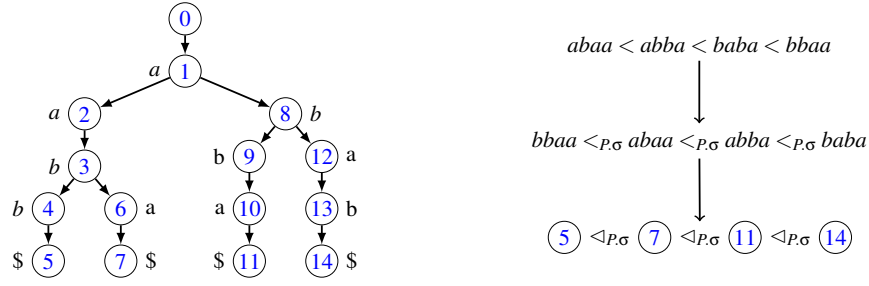
Figure 2 Tables LCP(P), LRS(P), SA(m<sub>P</sub>) and BWT(P) for the running example. One has LCP(P)[10] = LRS(P)[10] = 2 although the longest common prefix between suffixes of rank 9 and 10 is aa\$ab of length 5 (i.e., although LCP(m<sub>P</sub>)[10] = 5).

Example 1. With the tree of Figure 3, we have  $\delta^{(13)} = b$  and  $\pi^{(13)} = aba = \delta^{(12)}\delta^{(8)}\delta^{(1)}$ , and also  $\delta^{(4)} = b$  and  $\pi^{(4)} = baa = \delta^{(3)}\delta^{(2)}\delta^{(1)}$ . Thus,  $^{(13)} \prec ^{(4)}$  since  $\pi^{(13)} < \pi^{(4)}$ .

The Prefix Array (PA) of an ordered tree  $\mathcal{T}$  is the array of pointers to the nodes of  $\mathcal{T}$  (except the root of  $\mathcal{T}$ ) sorted in  $\prec$  order. The eXtended Burrows-Wheeler Transform (XBWT) [9]<sup>1</sup> of a tree  $\mathcal{T}$  is an array of symbols of  $\Sigma$  of length  $\#\text{PA}(\mathcal{T})$ , such that the entry at position  $i$  gives the label of the node  $\text{PA}(\mathcal{T})[i]$ . The eXtended Burrows-Wheeler Last (XBWL) [9]<sup>1</sup> of a tree  $\mathcal{T}$  is the bit array of length  $\#\text{PA}(\mathcal{T})$  such that  $\text{XBWL}(\mathcal{T})[i]$  equals 1 if the node  $\text{PA}(\mathcal{T})[i]$  is the last child of its parent, and 0 otherwise. Figure 5a illustrates XBWT and XBWL.

**eXtended Burrows-Wheeler Transform of Aho-Corasick automaton.** For a set of strings  $S := \{s_1, \dots, s_n\}$ , we denote by  $S^\$$  the set  $\{s_1\$, \dots, s_n\ \$\}$ . Let  $P$  be an ordered set of strings. We define  $\text{ACT}(P)$  as the variant of the Aho-Corasick tree of  $P.S^\$$  where the label of an arc is shifted on the deepest node of this arc (i.e., the label of a node  $v$  of  $\text{ACT}(P)$  is the label of  $(u, v)$  in  $\text{ACT}(P.S)$  where  $u$  is the parent of  $v$  in  $\text{ACT}(P)$ ) and equipped with the order  $\triangleleft_{P,\sigma}$ . Indeed,  $\triangleleft_{P,\sigma}$  is the order on the leaves satisfying: for  $u$  and  $v$  two leaves of  $\text{ACT}(P)$ ,  $u \triangleleft_{P,\sigma} v$

<sup>1</sup> In [9], the XBW-transform is defined as  $\text{XBW}(\mathcal{T})[i] := \langle \text{XBWT}[i], \text{XBWL}[i] \rangle$ , for any position  $i$ .



■ **Figure 3** Tree  $\text{ACT}(\overleftarrow{P})$  for the running example, and links between the orders  $<$ ,  $<_{P,\sigma}$ , and  $<_{P,\sigma}$ .

iff  $\pi(u) <_{P,\sigma} \pi(v)$ . We extend this order to the set of children of all nodes (see Figure 3). Note that  $\text{ACT}(P)$  differs from  $\text{ACT}(P.S)$  (defined above) to correspond to an input of the eXtended Burrows-Wheeler Transform.

### 3 Link between BWT and XBW of Aho-Corasick automaton

Here, we introduce a decomposition of a multi-string BWT that leads us to exhibit a bijection with the Aho-Corasick automaton (see Proposition 2). This builds on and extends Manzini's work [22]. We extend this bijection to the XBWT of the Aho-Corasick automaton (see Figure 4). For space reasons, many proofs are given in the Appendix (or the preprint [6]).

**Link between BWT and AC.** Given an ordered set of strings  $P$ ,  $\text{Decomp\_BWT}(P)$  is the integer interval partition of  $\llbracket 1, |m_P| \rrbracket$  such that

$$\llbracket i, j \rrbracket \in \text{Decomp\_BWT}(P) \quad \text{iff} \quad \begin{cases} \text{LCP}(P)[k] \neq \text{LRS}(P)[k], & \text{for } k \in \{i, j+1\} \\ \text{LCP}(P)[k] = \text{LRS}(P)[k], & \text{for } k \in \llbracket i+1, j \rrbracket. \end{cases}$$

We define  $\text{Dec\_Pre}$  as the function from  $\text{Decomp\_BWT}(P)$  to  $\Sigma^*$  such that

$$\text{Dec\_Pre}[u] := \overleftarrow{m_P[\text{SA}(m_P)[i] : \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1]},$$

for all  $u := \llbracket i, j \rrbracket \in \text{Decomp\_BWT}(P)$ .

► **Proposition 2.** *Dec\_Pre is a bijection between  $\text{Decomp\_BWT}(P)$  and  $\text{Prefix}(\overleftarrow{P.S})$ .*

Let us start the proof of Proposition 2 with the following Lemma.

► **Lemma 3.** *Let  $u = \llbracket i, j \rrbracket \in \text{Decomp\_BWT}(P)$ . For all  $k \in \llbracket i, j \rrbracket$ ,*

$$\text{Dec\_Pre}[u] = \overleftarrow{m_P[\text{SA}(m_P)[k] : \text{SA}(m_P)[k] + \text{LRS}(P)[k] - 1]}.$$

**Proof.** Let us show by contraposition that  $\text{LRS}(P)[k-1] = \text{LCP}(P)[k]$  for all  $k \in \llbracket i+1, j \rrbracket$ . Assume that there exists  $k \in \llbracket i+1, j \rrbracket$  such that  $\text{LRS}(P)[k-1] \neq \text{LCP}(P)[k]$ . Whenever  $\text{LRS}(P)[k-1] < \text{LCP}(P)[k]$ , we get by definition that  $\text{LRS}(P)[k] \geq \text{LCP}(P)[k]$ , and thus  $\text{LRS}(P)[k-1] < \text{LRS}(P)[k]$ . By the definition of  $\text{LRS}(P)$ ,  $m_P[\text{SA}(m_P)[k-1] + \text{LRS}(P)[k-1]] = \$$ . By the definition of  $\text{LCP}(P)$ , for all  $j \in \llbracket 1, \text{LCP}(P)[k-1] - 1 \rrbracket$ ,  $m_P[\text{SA}(m_P)[k] + j] = m_P[\text{SA}(m_P)[k-1] + j]$ . As  $\text{LRS}(P)[k-1] < \text{LCP}(P)[k]$ , we have  $m_P[\text{SA}(m_P)[k] + \text{LRS}(P)[k-1]] = m_P[\text{SA}(m_P)[k-1] + \text{LRS}(P)[k-1]] = \$$ , which is impossible since  $\text{LRS}(P)[k-1] < \text{LRS}(P)[k]$ . Whenever  $\text{LRS}(P)[k-1] > \text{LCP}(P)[k]$ , as  $\text{LCP}(P)[k] = \text{LRS}(P)[k]$ , the string  $m_P[\text{SA}(m_P)[k] : |m_P|]$  is lexicographically strictly smaller than  $m_P[\text{SA}(m_P)[k-1] : |m_P|]$ , which is impossible. This concludes the proof. ◀

**Proof.** Let  $u = \llbracket i, j \rrbracket$  be an interval of  $\text{Decomp\_BWT}(P)$ . First, we prove that  $\text{Dec\_Pre}[u] \in \text{Prefix}(\overleftarrow{P.S})$ , and then to prove the bijection, we show  $\text{Dec\_Pre}$  is injective and surjective. By definition,  $\text{LRS}(P)[i] = \text{select}_{\$}(m_P[\text{SA}(m_P)[i] : |m_P|], 1) - 1$ , we get that  $m_P[\text{SA}(m_P)[i] + \text{LRS}(P)[i]] = \$$  and for any  $j \in \llbracket \text{SA}(m_P)[i], \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1 \rrbracket$ ,  $m_P[j] \neq \$$ . Hence,  $m_P[\text{SA}(m_P)[i] : \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1]$  is a suffix of a string  $w$  of  $P.S$ , and thus  $\text{Dec\_Pre}[u]$  is a prefix of  $\overleftarrow{w}$  in  $\overleftarrow{P.S}$ .

**Injectiveness.** Let  $u_1 = \llbracket i_1, j_1 \rrbracket$  and  $u_2 = \llbracket i_2, j_2 \rrbracket$  be two elements of  $\text{Decomp\_BWT}(P)$ .

Without loss of generality, we take  $i_1 \leq i_2$ . Assume that  $\text{Dec\_Pre}[u_1] = \text{Dec\_Pre}[u_2]$ , we have that  $m_P[\text{SA}(m_P)[i_1] : \text{SA}(m_P)[i_1] + \text{LRS}(P)[i_1] - 1] = m_P[\text{SA}(m_P)[i_2] : \text{SA}(m_P)[i_2] + \text{LRS}(P)[i_2] - 1]$  and thus for all  $k \in \llbracket i_1, i_2 \rrbracket$ ,  $m_P[\text{SA}(m_P)[i_1] : \text{SA}(m_P)[i_1] + \text{LRS}(P)[i_1] - 1] = m_P[\text{SA}(m_P)[k] : \text{SA}(m_P)[k] + \text{LRS}(P)[k] - 1]$ . Hence, we have  $\text{LCP}(P)[k] = \text{LRS}(P)[i_1]$  and  $\text{LRS}(P)[k] = \text{LRS}(P)[i_1]$ . Therefore by the definition of  $\text{Decomp\_BWT}(P)$ , we get  $u_1 = u_2$ .

**Surjectiveness.** Let  $v$  be a prefix of a string of  $\overleftarrow{P.S}$ . The string  $\overleftarrow{v}$  is a suffix of a string of  $P.S$ . By the definition of  $m_P$ ,  $\overleftarrow{v}$  is a prefix of a suffix  $s$  of  $m_P$  such that  $s[\overleftarrow{v} + 1] = \$$ . By the definition of  $\text{BWT}(P)$ , the table  $\text{SA}(m_P)$  gives, for a position  $i$ , the starting position of the  $i^{\text{th}}$  suffix of  $m_P$  in lexicographic order. Hence, there is a bijection between  $\text{Suffix}(m_P)$  and the set of positions in  $\text{SA}(m_P)$ . Let  $k \in \llbracket 1, |m_P| \rrbracket$  such that  $s = m_P[\text{SA}(m_P)[k] : |m_P|]$ . As  $\overleftarrow{v}$  is a suffix of a string of  $P.S$  and a prefix of  $s$ , we have  $\overleftarrow{v} = m_P[\text{SA}(m_P)[k] : \text{SA}(m_P)[k] + \text{LRS}(P)[k] - 1]$ . We take  $u = \llbracket i, j \rrbracket \in \text{Decomp\_BWT}(P)$  such that  $k \in \llbracket i, j \rrbracket$ . By Lemma 3,  $\overleftarrow{v} = \text{Dec\_Pre}[u]$ . ◀

By Proposition 2, there exists an integer interval partition of  $\llbracket 1, |m_P| \rrbracket$  (i.e.,  $\text{Decomp\_BWT}(P)$ ) that is in bijection with the set of nodes of  $\text{AC}(\overleftarrow{P.S})$ . The following theorem extends this bijection by building an isomorphic graph of  $\text{AC}(P.S)$  whose set of nodes is  $\text{Decomp\_BWT}(\overleftarrow{P})$ . This states how to simulate an Aho-Corasick automaton using the BWT (similarly to [22]).

► **Theorem 4.** *Using tables  $\text{BWT}(\overleftarrow{P})$ ,  $\text{LCP}(\overleftarrow{P})$ ,  $\text{LRS}(\overleftarrow{P})$  and the function  $\text{LF}(\overleftarrow{P})$ , we can build the graph  $(\text{Decomp\_BWT}(\overleftarrow{P}), A_T(\overleftarrow{P}) \cup A_F(\overleftarrow{P}))$  that is isomorphic to  $\text{AC}(P.S)$ .*

$$A_T(P) := \{(u, v) \in \text{Decomp\_BWT}(P)^2 \mid \exists x \in u \text{ such that } \text{LF}(P)[x] \neq 0 \text{ and } \text{LF}(P)[x] \in v\},$$

$$A_F(P) := \{(u, v) \in \text{Decomp\_BWT}(P)^2 \mid \left( \max_{\substack{k < i \\ \text{LRS}(P)[k] = \min_{k < l \leq i} (\text{LCP}(P)[l])}} (k) \in v \text{ with } u := \llbracket i, j \rrbracket \right)\}.$$

More precisely, the graph  $(\text{Decomp\_BWT}(\overleftarrow{P}), A_T(\overleftarrow{P}))$  is isomorphic to the tree  $\text{ACT}(P.S)$  and the graph  $(\text{Decomp\_BWT}(\overleftarrow{P}), A_F(\overleftarrow{P}))$  is isomorphic to the tree  $\text{ACFL}(P.S)$  (see Appendix).

**Link between BWT and XBW.** After giving a new proof of the relation between Aho-Corasick automaton and BWT, we exhibit a new (bijective) link between the BWT and the XBW, which takes into account the order in the multi-string (Theorem 5). This leads to both, another construction algorithm of the XBW from the BWT, and to a construction of the BWT from the XBW, and thereby extends Manzini's results (Corollary 6).

Alike the definition of  $\text{Decomp\_BWT}(P)$  for an ordered set of strings  $P$ , we define  $\text{Decomp\_XBW}(\mathcal{T})$  of a tree  $\mathcal{T}$  of  $t$  nodes as the integer interval partition of  $\llbracket 1, t \rrbracket$  such that

$$\llbracket i, j \rrbracket \in \text{Decomp\_XBW}(\mathcal{T}) \text{ iff } \begin{cases} \text{XBWL}(\mathcal{T})[k] = 1, & \text{for } k \in \{i - 1, j\} \\ \text{XBWL}(\mathcal{T})[k] = 0, & \text{for } k \in \llbracket i, j - 1 \rrbracket. \end{cases}$$



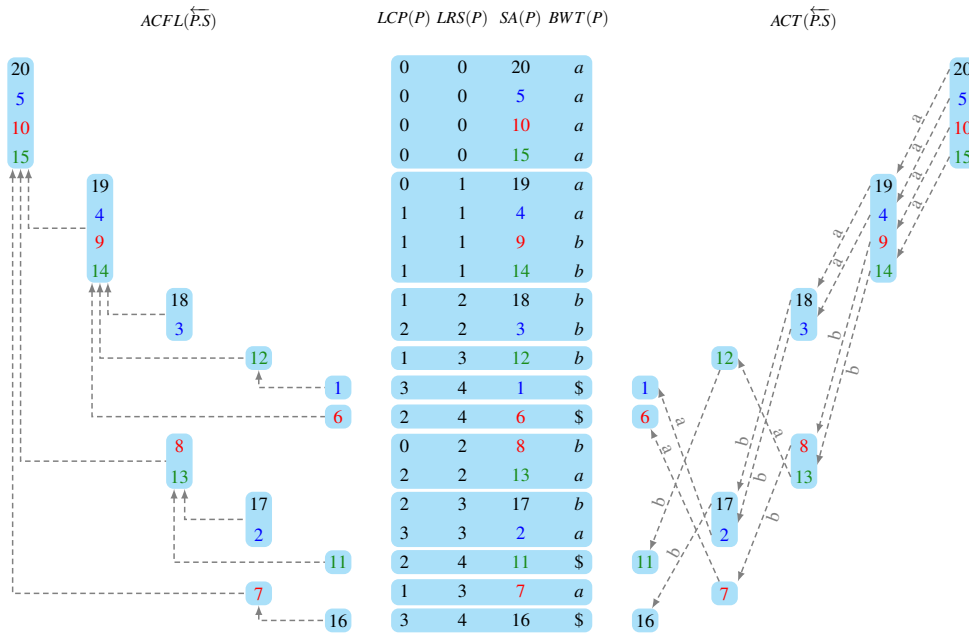


Figure 4 Link between  $BWT(P)$ ,  $ACT(\overleftarrow{P.S})$  and  $ACFL(\overleftarrow{P.S})$  for the running example. The blocks in the tables SA, LCP, LRS and BWT show the decomposition  $Decomp\_BWT(P)$ .

► **Theorem 5** (See Figure 5b). *There exists a bijection, denoted  $BWT\_XBW$ , between  $Decomp\_BWT(P)$  and  $Decomp\_XBW(ACT(\overleftarrow{P}))$  such that for all  $u \in Decomp\_BWT(P)$  with  $u := \llbracket i, j \rrbracket$ ,  $BWT\_XBW(u) := \llbracket i', j' \rrbracket$  and  $z := Parent_{ACT(\overleftarrow{P})}(PA(ACT(\overleftarrow{P}))[i'])$ :*

(1) *Let  $\{y_1, y_2, \dots, y_{\#u}\} := Leaves_{ACT(\overleftarrow{P})}(z)$  such that  $y_k \triangleleft_{P,\sigma} y_l \Rightarrow k < l$ ; then*

$$BWT(P)[i, j] = \delta[Child_{ACT(\overleftarrow{P})}(y_1)(z)]\delta[Child_{ACT(\overleftarrow{P})}(y_2)(z)] \dots \delta[Child_{ACT(\overleftarrow{P})}(y_{\#u})(z)].$$

(2) *Let  $\{x_1, x_2, \dots, x_{\#(BWT\_XBW(u))}\} := Children_{ACT(\overleftarrow{P})}(z)$  such that  $x_k \triangleleft_{P,\sigma} x_l \Rightarrow k < l$ ; then*

$$XBWT(ACT(\overleftarrow{P}))[i', j'] = \delta[x_1]\delta[x_2] \dots \delta[x_{\#(BWT\_XBW(u))}].$$

Theorem 5 provides us with a strong link between  $BWT(P)$  and  $XBWT(ACT(\overleftarrow{P}))$ , which allows transforming one structure into the other. This leads to the following corollary.

► **Corollary 6.** *Let  $P$  be an ordered set of strings.*

(1) *Using tables  $BWT(P)$ ,  $LCP(P)$  and  $LRS(P)$  of an ordered set of strings  $P$ , we can build the tables  $XBWT(ACT(\overleftarrow{P}))$  and  $XBWL(ACT(\overleftarrow{P}))$  in linear time of  $\|P.S\| \times \#\Sigma$ .*

(2) *Using tables  $XBWT(ACT(S))$  and  $XBWL(ACT(S))$  of a set of strings  $S$ , we can build the tables  $BWT(\overleftarrow{P})$ ,  $LCP(\overleftarrow{P})$  and  $LRS(\overleftarrow{P})$  in linear time of  $\|S\| \times \#\Sigma$  where  $P$  is an ordered set of strings such that  $P.S = S$ .*

The idea behind the algorithms is to exploit the link of Theorem 5 to compute each sub-string of the BWT or of the XBWT associated to each element of  $Decomp\_BWT$  or of  $Decomp\_XBW$ .

#### 4 Optimal ordering of strings for maximising compression

In this section, we show how to use Theorem 5 to compute the permutation that leads to a BWT and a XBWT having the minimum Run Length Encoding.

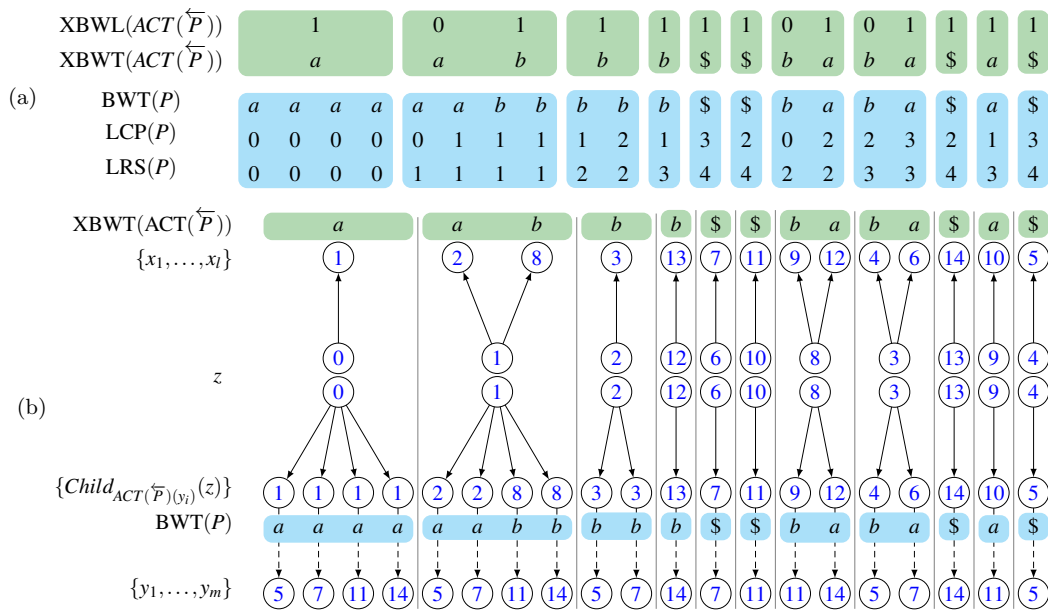


Figure 5 (a) Link between  $\text{Decomp\_BWT}(P)$  (in blue) and  $\text{Decomp\_XBWT}(\text{ACT}(\overleftarrow{P}))$  (in green). (b) Illustration of Theorem 5 for the running example.

### 4.1 Minimum permutation problem for BWT and XBWT

Run-Length Encoding [24] is a widely used method to compress strings. For a string  $w$ , the Run-Length Encoding splits  $w$  into the minimum number of substrings containing a single symbol. The size of the Run-Length Encoding of  $w$  is the cardinality of the minimum decomposition. For example for  $abbaaaccabb = a^1b^2a^3c^2a^1b^3$  (where the power notation  $\alpha^n$  means  $n$  copies of symbol  $\alpha$ ), the size of the Run-Length Encoding is 6 (for the decomposition has 6 blocks).

We define Run-Length measures:  $d_B$  for a BWT, and  $d_X$  for a XBW (similar to those of [15]). For  $P$  an ordered set of strings, let  $d_B(P)$  be the cardinality of the set  $\{i \in [1, \#\text{BWT}(P) - 1] \mid \text{BWT}(P)[i] \neq \text{BWT}(P)[i + 1]\}$ . Similarly, let  $d_X(P)$  be the cardinality of the set  $\{i \in [1, \#\text{XBWT}(\text{ACT}(P)) - 1] \mid \text{XBWT}(\text{ACT}(P))[i] \neq \text{XBWT}(\text{ACT}(P))[i + 1]\}$ .

Given two ordered sets of strings,  $P_1$  and  $P_2$  such that  $P_1.S = P_2.S$  (i.e., they contain the same set of strings), Theorem 5 implies that their BWT may differ, and thus  $d_B(P_1)$  and  $d_B(P_2)$  may also differ. We define the following minimisation problems. As the Run-Length Encoding of  $\text{BWT}(P)$  has size  $d_B(P) + 1$ , finding an optimal solution of  $\text{Min-Permutation-BWT}$  can help compressing  $\text{BWT}(P)$ .

► **Definition 7** ( $\text{Min-Permutation-BWT}$  and  $\text{Min-Permutation-XBWT}$ ). *Let  $S$  be a set of strings. The problem  $\text{Min-Permutation-BWT}$  asks for an ordered set of strings  $P$  that minimises  $d_B(P)$  and such that  $P.S = S$ . The problem  $\text{Min-Permutation-XBWT}$  asks for an ordered set of strings  $P$  that minimises  $d_X(P)$  and such that  $P.S = S$ .*

To simplify  $\text{Min-Permutation-BWT}$ , we consider specific ordered sets of strings. Let  $P$  be an ordered set of strings and let  $\perp$  denote the root of  $\text{ACT}(\overleftarrow{P})$ . We say that  $P$  is *topologically planar* if for each node  $u$  in  $\text{ACT}(\overleftarrow{P})$  and  $v \in \text{Leaves}_{\text{ACT}(\overleftarrow{P})}(\perp) \setminus \text{Leaves}_{\text{ACT}(\overleftarrow{P})}(u)$ , there does not exist  $u_1$  and  $u_2$  in  $\text{Leaves}_{\text{ACT}(\overleftarrow{P})}(u)$  such that  $u_1 \triangleleft_{P,\sigma} v \triangleleft_{P,\sigma} u_2$ . In other words,  $P$  is *topologically planar* if we can draw the tree  $\text{ACT}(\overleftarrow{P})$  by ordering the leaves with  $\triangleleft_{P,\sigma}$  without arcs crossing each other.

Let  $P$  be an ordered set of strings, which is not necessarily topologically planar. We denote by  $P_{tp}$  the ordered set of strings such that  $P_{tp}.S = P.S$ , and such that  $P_{tp}.\sigma$  satisfies for any  $u$  in  $\text{ACT}(\overleftarrow{P})$ , any  $v \in \text{Leaves}_{\text{ACT}(\overleftarrow{P})}(\perp) \setminus \text{Leaves}_{\text{ACT}(\overleftarrow{P})}(u)$ , and any  $u_1, u_2$  in  $\text{Leaves}_{\text{ACT}(\overleftarrow{P})}(u)$  such that  $u_1 \triangleleft_{P.\sigma} v \triangleleft_{P.\sigma} u_2$ , we have  $u_1 \triangleleft_{P_{tp}.\sigma} u_2 \triangleleft_{P_{tp}.\sigma} v$ . As we have a bijection between the set of circular permutations of  $P.S$  and the set of leaves of  $\text{ACT}(\overleftarrow{P})$ , we can unambiguously define the ordered set of strings  $P_{tp}$  that is topologically planar.

► **Proposition 8.** *Let  $P$  be an ordered set of strings. We have  $d_B(P_{tp}) \leq d_B(P)$  and  $d_B(P_{tp}) = d_X(\overleftarrow{P_{tp}})$ .*

**Proof.** For the inequality, let us prove that any modification of the order used to create  $P_{tp}$  decreases the value of  $d_B$ . Let  $P$  be an ordered set of strings that is not topologically planar. Let  $u$  in  $\text{ACT}(\overleftarrow{P})$ ,  $v \in \text{Leaves}_{\text{ACT}(\overleftarrow{P})}(\perp) \setminus \text{Leaves}_{\text{ACT}(\overleftarrow{P})}(u)$ , and let  $u_1$  and  $u_2$  in  $\text{Leaves}_{\text{ACT}(\overleftarrow{P})}(u)$  such that  $u_1 \triangleleft_{P.\sigma} v \triangleleft_{P.\sigma} u_2$ . Let  $P'$  be the copy of  $P$  where the only difference is  $u_1 \triangleleft_{P_{tp}.\sigma} u_2 \triangleleft_{P_{tp}.\sigma} v$ . Let  $x \in \text{Decomp\_BWT}(P)$  such that  $\text{Dec\_Pre}[x] = u$ . By Theorem 5, for all  $y \in \text{Decomp\_BWT}(P) \setminus \{x\}$ , we have  $\text{BWT}(P)[y] = \text{BWT}(P')[y]$ ,  $\text{BWT}(P)[x] = \dots \delta[u_1] \delta[v_1] \delta[u_2] \dots$  and  $\text{BWT}(P')[x] = \dots \delta[u_1] \delta[u_2] \dots \delta[v_1] \dots$  with  $v_1 = \text{Child}_{\text{ACT}(\overleftarrow{P})(v)}(u)$ . As  $\delta[u_1] = \delta[u_2]$  and  $\delta[u_1] \neq \delta[v_1]$ , we have  $d_B(P') \leq d_B(P)$ .

For the equality, it is enough to see that for any element  $u$  in  $\text{Decomp\_BWT}(P_{tp})$ , the numbers of distinct successive symbols are identical in  $\text{BWT}(P_{tp})[u]$  and in  $\text{XBWT}(\text{ACT}(\overleftarrow{P_{tp}}))[\text{BWT\_XBW}[u]]$ . Thus, for two successive elements  $u = \llbracket i_1, j_1 \rrbracket$  and  $v = \llbracket i_2, j_2 \rrbracket$  of  $\text{Decomp\_BWT}(P_{tp})$ , we obtain an equivalence between  $\text{BWT\_XBW}[u] = \llbracket i'_1, j'_1 \rrbracket$  and  $\text{BWT\_XBW}[v] = \llbracket i'_2, j'_2 \rrbracket$ :

$$\text{BWT}(P_{tp})[j_1] = \text{BWT}(P_{tp})[i_2] \quad \text{iff} \quad \text{BWT}(\text{ACT}(\overleftarrow{P_{tp}}))[j'_1] = \text{BWT}(\text{ACT}(\overleftarrow{P_{tp}}))[i'_2]. \quad \blacktriangleleft$$

Thanks to Proposition 8, we can restrict the search to ordered sets of strings that are topologically planar when solving **Min-Permutation-BWT** or **Min-Permutation-XBWT**. Furthermore, an optimal solution of **Min-Permutation-BWT** for  $S$  is also an optimal solution of **Min-Permutation-XBWT** for  $\overleftarrow{S}$ , and vice versa. This yields the following theorem.

► **Theorem 9.** *Let  $S$  be a set of strings. We can find an optimal solution for **Min-Permutation-BWT** and for **Min-Permutation-XBWT** in  $O(\|S\| \times \#\Sigma)$  time.*

## 4.2 Proof of Theorem 9

As a reminder, Proposition 8 states that an optimal solution of **Min-Permutation-BWT** is also an optimal solution of **Min-Permutation-XBWT**, and vice versa. In the following of this proof, we only prove the result regarding **Min-Permutation-XBWT**.

To start, let us give an overview of algorithm:

1. we take a random permutation  $\sigma$  of  $S$  and define  $P$  such that  $P.S = S$  and  $P.\sigma = \sigma$ ,
2. we build  $\text{ACT}(P)$ ,  $\text{XBWT}(\text{ACT}(P))$  and  $\text{Decomp\_XBW}(\text{ACT}(P))$ ,
3. we find  $P'$  which is an optimal solution of **Min-Permutation-XBWT**.

In the following, we define the problem **Min-Permutation-Table** and explicit its link to **Min-Permutation-XBWT** (Lemma 10). Lemma 10 gives us a linear time algorithm for finding an optimal solution of **Min-Permutation-Table**, and thus we can apply this algorithm to obtain an optimal solution for **Min-Permutation-XBWT**.

Given  $A$  an array of symbols of  $\Sigma$ , we define  $\text{Char}(A)$  as the set of (different) symbols in  $A$ . Given  $T$  an array of  $n$  symbols of  $\Sigma$  and  $D$  an integer interval partition of  $\llbracket 1, n \rrbracket$  such for each interval  $\llbracket i, j \rrbracket$  of  $D$ ,  $\#(\text{Char}(T[i, j])) = j - i + 1$  (i.e., all the symbols of

$T[i, j]$  are different), the problem **Min-Permutation-Table** is to find a  $T'$  such that for all  $\llbracket i, j \rrbracket \in D$ ,  $\#(\text{Char}(T[i, j])) = \#(\text{Char}(T'[i, j]))$  and which minimises  $d_A(T', D) := \#\{i \in \llbracket 1, n-1 \rrbracket \mid T'[i] \neq T'[i+1]\}$ . A proof of the following lemma is given in Appendix.

► **Lemma 10.** *Let  $S$  be a set of strings and let  $P$  be an ordered set of strings such that  $P.S = S$ . For an optimal solution  $T'$  of **Min-Permutation-Table** for  $\text{XBWT}(\text{ACT}(P))$  and for  $\text{Decomp\_XBW}(\text{ACT}(P))$ , there exists an optimal solution  $P'$  of **Min-Permutation-XBWT** for  $S$  such that  $\text{XBWT}(\text{ACT}(P')) = T'$ .*

Let  $T$  be an array of  $n$  symbols of  $\Sigma$  and let  $D$  be an integer interval partition of  $\llbracket 1, n \rrbracket$  such for each interval  $\llbracket i, j \rrbracket$  of  $D$ ,  $\#(\text{Char}(T[i, j])) = j - i + 1$ . Let  $A(D)$  be the array of all intervals in  $D$  in the order  $<$  and  $B(T, D)$  the array of size  $\#A(D) - 1$  such that the position  $i$  of  $B(T, D)$  is  $B(T, D)[i] = \text{Char}(A(D)[i]) \cap \text{Char}(A(D)[i+1])$ . For any set of symbols  $C := \{c_1, \dots, c_m\}$  where  $c_1 < \dots < c_m$ , we define  $\text{word}(C)$  as the string  $c_1 \dots c_m$ .

► **Lemma 11.** *Let  $T$  be an array of  $n$  symbols of  $\Sigma$  and let  $D$  be an integer interval partition of  $\llbracket 1, n \rrbracket$  such for any interval  $\llbracket i, j \rrbracket$  of  $D$ ,  $\#(\text{Char}(T[i, j])) = j - i + 1$ .*

- *If there exists  $i$  in  $\llbracket 1, n \rrbracket$  such that  $\llbracket i, i \rrbracket \in D$ , then  $T'_1[1, i-1]T'_2$  is an optimal solution of **Min-Permutation-Table** for  $T$  and for  $D$ , where  $T'_1$  is an optimal solution of **Min-Permutation-Table** for  $T[1, i]$  and for  $\{\llbracket i', j' \rrbracket \in D \mid j' \leq i\}$ , and  $T'_2$  is an optimal solution of **Min-Permutation-Table** for  $T[i, n]$  and for  $\{\llbracket i', j' \rrbracket \in D \mid i' \geq i\}$ .*
- *If there exists  $i$  in  $\llbracket 1, \#B(T, D) \rrbracket$  such that  $\#(B(T, D)[i]) = 0$ , then  $T'_1T'_2$  is an optimal solution of **Min-Permutation-Table** for  $T$  and for  $D$ , where  $T'_1$  is an optimal solution of **Min-Permutation-Table** for  $T[1, A(D)[i][1]]$  and for  $\{\llbracket i', j' \rrbracket \in D \mid j' \leq A(D)[i][1]\}$ , and  $T'_2$  is an optimal solution of **Min-Permutation-Table** for  $T[A(D)[i+1][0], n]$  and for  $\{\llbracket i', j' \rrbracket \in D \mid i' \geq A(D)[i+1][0]\}$ .*
- *If there exists  $i$  in  $\llbracket 1, \#B(T, D) \rrbracket$  such that  $\#(B(T, D)[i]) = 1$ , then  $T'_1aT'_2$  is an optimal solution of **Min-Permutation-Table** for  $T$  and for  $D$ , where  $B(T, D)[i] = \{a\}$ ,  $T'_1$  is an optimal solution of **Min-Permutation-Table** for  $T[1, A(D)[i-1][1]]\text{word}(\text{Char}(A(D)[i] \setminus \{a\}))$  and for  $\{\llbracket i', j' \rrbracket \in D \mid j' < A(D)[i][1] \cup \llbracket A(D)[i][0], A(D)[i][1] - 1 \rrbracket\}$ , and  $T'_2$  is an optimal solution of **Min-Permutation-Table** for  $\text{word}(\text{Char}(A(D)[i+1] \setminus \{a\}))T[A(D)[i+2][0], n]$  and for  $\{\llbracket i', j' \rrbracket \in D \mid i' > A(D)[i+1][0] \cup \llbracket A(D)[i+1][0] + 1, A(D)[i][1] \rrbracket\}$ .*

**Proof.** All the proofs are derived from the equality  $d_A(T[1, n], D) = d_A(T[1, i], \{\llbracket i', j' \rrbracket \in D \mid j' \leq i\}) + d_A(T[i, n], \{\llbracket i', j' \rrbracket \in D \mid i' \geq i\})$  for all  $i \in \llbracket 1, n \rrbracket$ . ◀

► **Lemma 12.** *Let  $T$  be an array of  $n$  symbols of  $\Sigma$  and let  $D$  be an integer interval partition of  $\llbracket 1, n \rrbracket$  such that for any interval  $\llbracket i, j \rrbracket$  of  $D$ ,  $\#(\text{Char}(T[i, j])) = j - i + 1$ . Whenever for any  $i \in \llbracket 1, \#B(T, D) \rrbracket$  one has  $\#(B(T, D)[i]) \geq 2$ , Algorithm 1 gives an optimal solution for **Min-Permutation-Table** in  $O(\#T \times \#\Sigma)$  time.*

**Proof.**

**Complexity.** To build the tables  $A(D)$  and  $B(T, D)$ , we need  $O(\#T \times \#\Sigma)$  time. As both tables are smaller than  $T$ , the **for** loop of Algorithm 1 also takes  $O(\#T \times \#\Sigma)$  time.

**Optimality.** As for any  $i$  in  $\llbracket 1, \#B(T, D) \rrbracket$ , one has  $\#(B(T, D)[i]) \geq 2$ , we get  $B(T, D)[i] \setminus \{\text{last}\} \neq \emptyset$  (within the **for** loop of Algorithm 1). Let  $T^*$  be a string such that for all  $\llbracket i, j \rrbracket \in D$ ,  $\text{Char}(T[i, j]) = \text{Char}(T^*[i, j])$ . The size of  $T^*$  is  $\#T$  and the number of intervals of  $D$  is  $\#D$ , i.e., the maximum number of positions where two consecutive letters can be identical. Hence, we have  $d_A(T^*, D) \leq \#T - \#D + 1$ . Let  $T'$  be the string given by Algorithm 1. We have  $d_A(T', D) = \sum_{u \in D} (\#u - 1) + 1 = \#T - \#D + 1$ . This concludes the proof. ◀

---

**Algorithm 1:** Computation of an array  $T'$  of symbols satisfying for any  $\llbracket i, j \rrbracket \in D$ ,  $\text{Char}(T[i, j]) = \text{Char}(T'[i, j])$ .

---

**Input** : An instance of Min-Permutation-Table  $T$  and  $D$

**Output** : A string  $T'$

$last \leftarrow \$$  such that  $\$ \notin \Sigma$ ;

$T' \leftarrow$  empty string;

**for**  $i \in \llbracket 1, \#B(T, D) \rrbracket$  **do**

$lettre \leftarrow$  random element of  $B(T, D)[i] \setminus \{last\}$ ;

$T' \leftarrow T'$  word( $\text{Char}(A(D)[i]) \setminus \{lettre, last\}$ );

$T' \leftarrow T'$   $lettre$   $lettre$ ;

$last \leftarrow lettre$ ;

$T' \leftarrow T'$  word( $\text{Char}(A(D)[\#A(D)]) \setminus \{last\}$ );

**return**  $T'$ ;

---

► **Lemma 13.** Let  $T$  be an array of  $n$  symbols of  $\Sigma$  and let  $D$  be an integer interval partition of  $\llbracket 1, n \rrbracket$  such that for each interval  $\llbracket i, j \rrbracket$  of  $D$ ,  $\#(\text{Char}(T[i, j])) = j - i + 1$ . The problem **Min-Permutation-Table** can be solved in  $O(\#T \times \#\Sigma)$  time.

**Proof.** By Lemma 12 and Lemma 11, we can compute an optimal solution of **Min-Permutation-Table** by splitting the interval, applying Algorithm 1 on each part, and then merging the strings output by Algorithm 1. ◀

## 5 Conclusion and Perspectives

Here, we present a new view of the Burrows-Wheeler Transform: as the text representation of an Aho-Corasick automaton that depends on the concatenation order. This induces a link between the Burrows-Wheeler Transform and the eXtended Burrows-Wheeler Transform, via the Aho-Corasick automaton. This link allows one to transform one structure into the other (for which we provide algorithms). We also exploit this link to find in linear time an ordering of input strings that optimises the compression of the concatenated strings.

Of course, it would be interesting to evaluate even empirically the gain of this compression on real life data. In bioinformatics, one wishes to index a collection of genomes from individual of the same species. For instance, it can be all variants of a virus genome within a host (*e.g.*, HIV or Ebola virus in human); their number and relative frequency may change along time [2]. An application is then to compare the sequencing reads obtained from any new infected individual with this index to determine whether some variants are becoming more frequent or resistant to some treatment. Some viruses evolve rapidly to circumvent immune response, the number of potential variants can be large. Hence, it is practically relevant to reduce the index space by finding an optimal permutation of the genomes that would maximised the Run-Length Encoding of the index.

---

## References

- 1 Alfred V. Aho and Margaret J. Corasick. Efficient String Matching: An Aid to Bibliographic Search. *Communications of the ACM*, 18(6):333–340, 1975.
- 2 Jasmijn A. Baaijens, Amal Zine El Aabidine, Eric Rivals, and Alexander Schönhuth. De novo assembly of viral quasispecies using overlap graphs. *Genome Research*, 27(5):835–848, 2017.

- 3 Markus J. Bauer, Anthony J. Cox, and Giovanna Rosone. Lightweight algorithms for constructing and inverting the BWT of string collections. *Theoretical Computer Science*, 483:134–148, 2013.
- 4 Djamal Belazzougui. Succinct Dictionary Matching with No Slowdown. In *Combinatorial Pattern Matching*, pages 88–100, 2010.
- 5 Michael Burrows and David J. Wheeler. A block-sorting lossless data compression algorithm. Technical report, Digital Equipment Corporation, 1994.
- 6 Bastien Cazaux and Eric Rivals. Strong link between BWT and XBW via Aho-Corasick automaton and applications to Run-Length Encoding. *CoRR*, abs/1805.10070, 2018. [arXiv:1805.10070](#).
- 7 Anthony J. Cox, Markus J. Bauer, Tobias Jakobi, and Giovanna Rosone. Large-scale compression of genomic sequence databases with the Burrows–Wheeler Transform. *Bioinformatics*, 28(11):1415–1419, 2012.
- 8 Lavinia Egidi and Giovanni Manzini. Lightweight BWT and LCP Merging via the Gap Algorithm. In *String Processing and Information Retrieval*, pages 176–190, 2017.
- 9 Paolo Ferragina, Fabrizio Luccio, Giovanni Manzini, and S. Muthukrishnan. Compressing and indexing labeled trees, with applications. *Journal of the ACM*, 57(1):1–33, 2009.
- 10 Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *Journal of the ACM*, 52(4):552–581, 2005.
- 11 Luca Foschini, Roberto Grossi, Ankur Gupta, and Jeffrey Scott Vitter. When indexing equals compression: Experiments with compressing suffix arrays and applications. *ACM Transactions on Algorithms*, 2(4):611–639, 2006.
- 12 Travis Gagie, Giovanni Manzini, and Jouni Sirén. Wheeler graphs: A framework for BWT-based data structures. *Theoretical Computer Science*, 698:67–78, 2017.
- 13 Roberto Grossi, Ankur Gupta, and Jeffrey Scott Vitter. High-order entropy-compressed text indexes. In *ACM-SIAM Symp. on Discrete Algorithms*, pages 841–850, 2003.
- 14 James Holt and Leonard McMillan. Constructing Burrows-Wheeler Transforms of large string collections via merging. In *ACM Conf. on Bioinformatics, Computational Biology, and Health Informatics*, pages 464–471, 2014.
- 15 James Holt and Leonard McMillan. Merging of multi-string BWTs with applications. *Bioinformatics*, 30(24):3524–3531, 2014.
- 16 Wing-Kai Hon, Tsung-Han Ku, Rahul Shah, Sharma V. Thankachan, and Jeffrey Scott Vitter. Faster compressed dictionary matching. *Theoretical Computer Science*, 475:113–119, 2013.
- 17 Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*, 25(14):1754–1760, 2009.
- 18 Veli Mäkinen and Gonzalo Navarro. Rank and select revisited and extended. *Theoretical Computer Science*, 387(3):332–347, 2007.
- 19 Veli Mäkinen, Gonzalo Navarro, Jouni Sirén, and Niko Välimäki. Storage and Retrieval of Individual Genomes. In *RECOMB, Tucson, AZ, USA*, pages 121–137, 2009.
- 20 Udi Manber and Eugene W. Myers. Suffix Arrays: A New Method for On-Line String Searches. *SIAM Journal on Computing*, 22(5):935–948, 1993.
- 21 Sabrina Mantaci, Antonio Restivo, Giovanna Rosone, and Marinella Sciortino. An Extension of the Burrows Wheeler Transform and Applications to Sequence Comparison and Data Compression. In *Combinatorial Pattern Matching*, pages 178–189, 2005.
- 22 Giovanni Manzini. XBWT tricks. In *String Processing and Information Retrieval*, pages 80–92, 2016.
- 23 Jouni Sirén. Burrows-Wheeler Transform for Terabases. In *Data Compression Conf.*, pages 211–220, 2016.
- 24 Jouni Sirén, Niko Välimäki, Veli Mäkinen, and Gonzalo Navarro. Run-Length Compressed Indexes Are Superior for Highly Repetitive Sequence Collections. In *String Processing and Information Retrieval*, pages 164–175, 2009.

## A Appendix

### A.1 Linear time construction of tables $\text{LRS}(P)$ and $\text{LCP}(P)$

To build in linear time the tables  $\text{LRS}(P)$  and  $\text{LCP}(P)$ , we use Algorithm 2, which is proved in Lemma 14. We summarise this property in Proposition 15. The tables for the running example are illustrated in Figure 2.

► **Lemma 14.** *Let  $i \in \llbracket 1, |m_P| \rrbracket$ .*

$$\text{LRS}(P)[i] = \begin{cases} 0 & \text{if } i \in \llbracket 1, \#(P.S) \rrbracket, \\ \text{LRS}(P)[\text{LF}(P)[i]] - 1 & \text{otherwise.} \end{cases}$$

**Proof.** For any  $i \in \llbracket 1, |w| \rrbracket$ , we know that  $\text{SA}(w)[\text{LF}(w)[i]] = \text{SA}(w)[i] - 1$  if  $\text{SA}(w)[i] \geq 2$  and  $\text{SA}(w)[\text{LF}(w)[i]] = |w|$  otherwise; hence, we get

■ if  $\text{SA}(w)[i] \geq 2$ ,

$$\begin{aligned} \text{LRS}(P)[\text{LF}(P)[i]] &= \text{select}_{\$}(m_P[\text{SA}(m_P)[\text{LF}(P)[i]] : |m_P|], 1) - 1 \\ &= \text{select}_{\$}(m_P[\text{SA}(w)[i] - 1 : |m_P|], 1) - 1 \\ &= \begin{cases} 0 & \text{if } m_P[\text{SA}(w)[i] - 1] = \$, \\ \text{LRS}(P)[i] + 1 & \text{otherwise.} \end{cases} \end{aligned}$$

■ If  $\text{SA}(w)[i] = 1$ ,

$$\begin{aligned} \text{LRS}(P)[\text{LF}(P)[i]] &= \text{select}_{\$}(m_P[\text{SA}(m_P)[\text{LF}(P)[i]] : |m_P|], 1) - 1 \\ &= \text{select}_{\$}(m_P[|m_P| : |m_P|], 1) - 1 \\ &= 0. \end{aligned}$$

We define the function **Letter** from  $\llbracket 1, |m_P| \rrbracket$  to  $\Sigma$  such that  $C[\text{Letter}[i]] < i \leq C[\text{Letter}[i+1]]$  (see the definition of **LF** on page 3) where  $\Sigma = \{c_1 \dots, c_{\#\Sigma}\}$  with  $c_1 < \dots < c_{\#\Sigma}$ . Thus, we define **RLF**( $P$ ) from  $\llbracket 1, |m_P| \rrbracket$  to  $\llbracket 1, |m_P| \rrbracket$  such that

$$\text{RLF}(P)[i] = \text{select}_{\text{Letter}[i]}(\text{BWT}(P), i - C[\text{Letter}[i]]).$$

Let  $i$  be an integer between 1 and  $|m_P|$ . We have

$$\begin{aligned} \text{RLF}(P)[\text{LF}(P)[i]] &= \text{select}_{\text{Letter}[\text{LF}(P)[i]]}(\text{BWT}(P), \text{LF}(P)[i] - C[\text{Letter}[\text{LF}(P)[i]]]) \\ &= \text{select}_{\text{BWT}(P)[i]}(\text{BWT}(P), \text{LF}(P)[i] - C[\text{BWT}(P)[i]]) \\ &= \text{select}_{\text{BWT}(P)[i]}(\text{BWT}(P), \text{rank}_{\text{BWT}(w)[i]}(\text{BWT}(w), i)) \\ &= i. \end{aligned}$$

Hence, the function **RLF**( $P$ ) is the reverse bijection of **LF**( $P$ ), and as  $m_P[\text{SA}(m_P)[i] - 1] = \text{BWT}(P)[i]$ , one gets

$$\begin{aligned} \text{LRS}(P)[i] = 0 &\Leftrightarrow \text{BWT}(P)[\text{RLF}(P)[i]] = \$ \\ &\Leftrightarrow \text{BWT}(P)[\text{select}_{\text{Letter}[i]}(\text{BWT}(P), i - C[\text{Letter}[i]])] = \$ \\ &\Leftrightarrow \text{Letter}[i] = \$ \\ &\Leftrightarrow i \in \llbracket 1, \#(P.S) \rrbracket. \end{aligned}$$

Therefore, we derive the desired equality, which concludes the proof. ◀

► **Proposition 15.** *Let  $P$  be an ordered set of strings. Using tables  $\text{BWT}(P)$  and  $\text{LCP}(m_P)$ , we can compute the tables  $\text{LRS}(P)$  and  $\text{LCP}(P)$  in a time that is linear in  $|m_P|$ .*

**Proof.** As the value of  $\text{LCP}(P)$  at each position corresponds to the minimum between the values at same positions in  $\text{LCP}(m_P)$  and in  $\text{LRS}(P)$ , we only need to prove that the table  $\text{LRS}(P)$  can be computed from  $\text{BWT}(P)$  in linear time.

Using Algorithm 2 and Lemma 14, we can compute table  $\text{LRS}(P)$  in  $O(|m_P|)$  time.

---

**Algorithm 2:** Computation of table  $\text{LRS}(P)$ .

---

**Input** : The ordered set of strings  $P$   
**Output** : The table  $\text{LRS}(P)$   
 Compute the table  $\text{BWT}(P)$ ;  
 Initialise  $\text{LRS}$  as an empty table of length  $|m_P|$ ;  
**for**  $i \in \llbracket 1, \#(P.S) \rrbracket$  **do**  
    $position \leftarrow i$ ;  
    $nb \leftarrow 0$ ;  
    $\text{LRS}[position] \leftarrow nb$ ;  
   **while**  $\text{BWT}(P)[position] \neq \$$  **do**  
      $\text{LRS}[position] \leftarrow nb$ ;  
      $nb \leftarrow nb + 1$ ;  
      $position \leftarrow \text{LF}[position]$ ;  
**return**  $\text{LRS}$ ;

---

## A.2 Proof of Theorem 4

The next propositions exhibit the links between  $\text{Decomp\_BWT}(P)$  and first, the arcs of  $\text{ACT}(\overleftarrow{P.S})$  (Proposition 16), and second, the arcs of  $\text{ACFL}(\overleftarrow{P.S})$  (Proposition 17).

► **Proposition 16** (See Figure 4). *The graph  $G_T(P) = (\text{Decomp\_BWT}(P), A_T(P))$  is isomorphic to the tree  $\text{ACT}(\overleftarrow{P.S})$ , where*

$$A_T(P) := \{(u, v) \in \text{Decomp\_BWT}(P)^2 \mid \exists x \in u \text{ such that } \text{LF}(P)[x] \neq 0 \text{ and } \text{LF}(P)[x] \in v\}.$$

**Proof.** First, we show that there exists a bijection between the node set of  $\text{ACT}(\overleftarrow{P.S})$  and that of  $G_T(P)$ . We reuse the bijection  $\text{Dec\_Pre}$ , which served in Proposition 2. Let us show that for each arc  $(u, v)$  of  $A_T(P)$ ,  $(\text{Dec\_Pre}[u], \text{Dec\_Pre}[v])$  is an arc of  $\text{ACT}(\overleftarrow{P.S})$ , and vice versa. Let  $(u, v) \in A_T(P)$ , i.e.,  $(u, v) \in \text{Decomp\_BWT}(P)^2$  such that there exists  $x \in u$  satisfying  $\text{LF}(P)[x] \in v$ .

According to [5], we know that  $\text{SA}(m_P)[\text{LF}(m_P)[i]] = \text{SA}(m_P)[i] - 1$ . By Lemma 14, we have  $\text{LRS}(P)[\text{LF}(P)[x]] = \text{LRS}(P)[x] + 1$  for all  $x$  such that  $\text{LF}(P)[x] \neq 0$ . With both equalities and Lemma 3, we obtain

$$\begin{aligned} \text{Dec\_Pre}[v] &= \overleftarrow{m_P[\text{SA}(m_P)[\text{LF}(P)[x]] : \text{SA}(m_P)[\text{LF}(P)[x]] + \text{LRS}(P)[\text{LF}(P)[x]] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[x] - 1 : \text{SA}(m_P)[x] - 1 + \text{LRS}(P)[\text{LF}(P)[x]] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[x] - 1 : \text{SA}(m_P)[x] - 1 + \text{LRS}(P)[x]]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[x] : \text{SA}(m_P)[x] + \text{LRS}(P)[x] - 1] m_P[\text{SA}(m_P)[x] - 1]} \\ &= \text{Dec\_Pre}[u] m_P[\text{SA}(m_P)[x] - 1]. \end{aligned}$$



The string  $\text{Dec\_Pre}[u]$  is thus the longest prefix of  $\text{Dec\_Pre}[v]$ .

Let  $(x, y)$  be an arc of  $\text{ACT}(\overleftarrow{P.S})$ . We take  $z$  a leaf in the subtree of  $\text{ACT}(\overleftarrow{P.S})$  in  $y$ . As  $x$  is the parent of  $y$  in  $\text{ACT}(\overleftarrow{P.S})$ ,  $z$  is also a leaf in the subtree of  $\text{ACT}(\overleftarrow{P.S})$  in  $x$ . We take  $i \in \llbracket 1, |m_P| \rrbracket$  such that  $\overleftarrow{z}$  is a prefix of  $m_P[i : |m_P|]$ . Hence  $\overleftarrow{y}$  is a prefix of  $m_P[i + |z| - |y| : |m_P|]$  and  $\overleftarrow{x}$  is a prefix of  $m_P[i + |z| - |x| : |m_P|]$ . As  $(x, y)$  is an arc of  $\text{ACT}(\overleftarrow{P.S})$ ,  $|y| - |x| = 1$ . Thus, choosing  $k$  such that  $\text{SA}(m_P)[k] = i + |z| - |y| + 1$ , and  $u, v$  in  $\text{Decomp\_BWT}(P)^2$  such that  $k \in u$  and  $\text{LF}(P)[k] \in v$ , we get  $\overleftarrow{x} = \text{Dec\_Pre}[u]$  and  $\overleftarrow{y} = \text{Dec\_Pre}[v]$ . This concludes the proof.  $\blacktriangleleft$

► **Proposition 17** (See Figure 4). *The graph  $G_F(P) = (\text{Decomp\_BWT}(P), A_F(P))$  is isomorphic to the tree  $\text{ACFL}(\overleftarrow{P.S})$ , where*

$$A_F(P) := \{(u, v) \in \text{Decomp\_BWT}(P)^2 \mid \left( \max_{\substack{k < i \\ \text{LRS}(P)[k] = \min_{k-1 \leq l \leq i} (\text{LCP}(P)[l])}} (k) \right) \in v \text{ with } u = \llbracket i, j \rrbracket\}.$$

**Proof.** First, let us show the following equivalence. Let  $u = \llbracket i, j \rrbracket \in \text{Decomp\_BWT}(P)$ .

$$\begin{aligned} w \in \text{Decomp\_BWT}(P) \text{ such that } \exists k \in w \text{ with } k < i \text{ and } \text{LRS}(P)[k] = \min_{k-1 \leq l \leq i} (\text{LCP}(P)[l]) \\ \Leftrightarrow \\ \text{Dec\_Pre}[w] \text{ is a suffix of } \text{Dec\_Pre}[u]. \end{aligned}$$

Let  $w \in \text{Decomp\_BWT}(P)$  such that  $\exists k \in w$  with  $k < i$  and  $\text{LRS}(P)[k] = \min_{k-1 \leq l \leq i} (\text{LCP}(P)[l])$ . Hence, we have for all  $l \in \llbracket k-1, i \rrbracket$ ,  $\text{LRS}(P)[k] \leq \text{LCP}(P)[l]$ , and thus

$$\begin{aligned} \text{Dec\_Pre}[u] &= \overleftarrow{m_P[\text{SA}(m_P)[i] : \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[i] + \text{LRS}(P)[k] : \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[i] : \text{SA}(m_P)[i] + \text{LRS}(P)[k] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[i] + \text{LRS}(P)[k] : \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[k] : \text{SA}(m_P)[k] + \text{LRS}(P)[k] - 1]} \\ &= \overleftarrow{m_P[\text{SA}(m_P)[i] + \text{LRS}(P)[k] : \text{SA}(m_P)[i] + \text{LRS}(P)[i] - 1]} \text{Dec\_Pre}[w]. \end{aligned}$$

Let  $w = \llbracket i_1, j_1 \rrbracket$  and  $u = \llbracket i_2, j_2 \rrbracket$  be two elements of  $\text{Decomp\_BWT}(P)$  such that  $\text{Dec\_Pre}[w]$  is a suffix of  $\text{Dec\_Pre}[u]$ . Hence, we have that  $m_P[\text{SA}(m_P)[i_1] : \text{SA}(m_P)[i_1] + \text{LRS}(P)[i_1] - 1]$  is a prefix of  $m_P[\text{SA}(m_P)[i_2] : \text{SA}(m_P)[i_2] + \text{LRS}(P)[i_2] - 1]$ . By the definition of  $\text{BWT}(P)$ , for all  $l \in \llbracket i_1, i_2 \rrbracket$ ,  $\text{LRS}(P)[i_1] \leq \text{LCP}(P)[l]$ . This concludes the proof of the equivalence. By the equivalence, given  $u$  and  $v$  in  $\text{Decomp\_BWT}(P)$  such that  $\text{Dec\_Pre}[u]$  is a suffix of  $\text{Dec\_Pre}[v]$ , for all  $k_1 \in u$  and  $k_2 \in v$ , we have  $k_1 \leq k_2$ . Hence, by taking the largest  $w$  satisfying the first step of the inequality, we obtain the longest suffix and vice versa.  $\blacktriangleleft$

### A.3 Proof of Theorem 5

**Proof.** We define  $T_B$  (resp.  $T_X$ ) as the array of intervals of  $\text{Decomp\_BWT}(P)$  (resp.  $\text{Decomp\_XBW}(\text{ACT}(\overleftarrow{P}))$ ) sorted in the interval order. Let us prove that  $T_B$  and  $T_X$  have the same length, and that at the same position  $i$ ,  $T_B[i]$  and  $T_X[i]$  represent the same prefix of  $\overleftarrow{P}$ . By Proposition 2, the length of  $T_B$  is  $\#\text{Prefix}(\overleftarrow{P.S})$ . By the definition of  $\text{Decomp\_XBW}$ , the length of  $T_X$  is the number of 1 in  $\text{XBWL}(\text{ACT}(\overleftarrow{P}))$ , i.e., the number of internal nodes of  $\text{ACT}(\overleftarrow{P})$ , and thus is equal to  $\#\text{Prefix}(\overleftarrow{P.S})$ . Hence,  $T_B$  and  $T_X$  have the same cardinalities.

Let  $i \in \llbracket 1, \#\text{Prefix}(\overleftarrow{P.S}) \rrbracket$ . By Proposition 2,  $T_B[i]$  represents the  $i^{\text{th}}$  suffix of strings of  $P.S$  in lexicographic order. By the definition of  $\text{XBWT}(\text{ACT}(\overleftarrow{P}))$ , all the nodes  $\text{PA}(\text{ACT}(\overleftarrow{P}))[k]$  for  $k \in T_X[i]$  have the same parent  $z$  in  $\text{ACT}(\overleftarrow{P})$ , and  $z$  represents the  $i^{\text{th}}$  node in  $\prec$  order, i.e., the  $i^{\text{th}}$  suffix of strings of  $P.S$  in lexicographic order.

As for a given position  $i$ ,  $T_B[i]$  and  $T_X[i]$  represent the same prefix of  $\overleftarrow{P}$ , we define the bijection  $\text{BWT\_XBW}$  such that for any  $i$  in  $\llbracket 1, \#\text{Prefix}(\overleftarrow{P}.S) \rrbracket$ , one has  $\text{BWT\_XBW}[T_B[i]] = T_X[i]$ . As the tree  $\text{ACT}(\overleftarrow{P})$  represents the Aho-Corasick tree of  $\overleftarrow{P}.S^\$$ , we have a bijection  $b_1$  from the node set of  $\text{ACT}(\overleftarrow{P})$  onto the set of prefixes of  $\overleftarrow{P}.S^\$$ . By the definition of functions  $\pi$  and  $\delta$  (see page 4), for any node  $v$  of  $\text{ACT}(\overleftarrow{P})$ , we have  $b_1(v) = \overleftarrow{\pi}(v) \delta(v)$ .

By the definition of  $\text{XBWT}(\text{ACT}(\overleftarrow{P}))$ , we have that for  $T_X[i] = \llbracket i', j' \rrbracket$

$$\begin{aligned} \#T_X[i] &= \#\{v \text{ node of } \text{ACT}(\overleftarrow{P}) \mid \pi(v) = \overleftarrow{\text{Dec\_Pre}}[\text{BWT\_XBW}^{-1}[T_X[i]]]\} \\ &= \#\{v \text{ node of } \text{ACT}(\overleftarrow{P}) \mid v \text{ is a child of } b_1^{-1}(\text{Dec\_Pre}[\text{BWT\_XBW}^{-1}[T_X[i]])]\}. \end{aligned}$$

As  $i' \in T_X[i]$  and  $b_1^{-1}(\text{Dec\_Pre}[\text{BWT\_XBW}^{-1}[T_X[i]])] = \text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])$ ,

$$\begin{aligned} \#T_X[i] &= \#\{v \text{ node of } \text{ACT}(\overleftarrow{P}) \mid v \text{ is a child of } \text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])\} \\ &= \#\text{Children}_{\text{ACT}(\overleftarrow{P})}(\text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])). \end{aligned}$$

Let  $\{x_1, x_2, \dots, x_{\#T_X[i]}\}$  be the set of children of  $\text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])$  sorted such that  $x_1 \triangleleft_{P,\sigma} \dots \triangleleft_{P,\sigma} x_{\#T_X[i]}$ . By the definition of  $\text{XBWT}(\text{ACT}(\overleftarrow{P}))$ , for any  $k \in \llbracket 1, \#T_X[i] \rrbracket$ , we have  $\text{XBWT}(\text{ACT}(\overleftarrow{P}))[i' + k - 1] = \delta[x_k]$ . By the definition of  $\text{BWT}(P)$ , for  $T_B[i] = \llbracket i'', j'' \rrbracket$ , we get

$$\begin{aligned} \#T_B[i] &= \#\{w \in \overleftarrow{P}.S^\$ \mid \text{Dec\_Pre}[T_B[i]] \text{ is a prefix of } w\} \\ &= \#\{v \text{ leaf of } \text{ACT}(\overleftarrow{P}) \mid b_1^{-1}(\text{Dec\_Pre}[T_B[i]]) \text{ is an ancestor of } v \text{ in } \text{ACT}(\overleftarrow{P})\}. \end{aligned}$$

As  $b_1^{-1}(\text{Dec\_Pre}[T_B[i]]) = b_1^{-1}(\text{Dec\_Pre}[\text{BWT\_XBW}^{-1}[T_X[i]])] = \text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])$ ,

$$\begin{aligned} \#T_B[i] &= \#\{v \text{ a leaf of } \text{ACT}(\overleftarrow{P}) \mid \text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i']) \\ &\quad \text{is an ancestor of } v \text{ in } \text{ACT}(\overleftarrow{P})\} \\ &= \#\text{Leaves}_{\text{ACT}(\overleftarrow{P})}(\text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])). \end{aligned}$$

Let  $\{y_1, y_2, \dots, y_{\#T_B[i]}\}$  be the set of leaves of the subtree of  $\text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i'])$  in  $\text{ACT}(\overleftarrow{P})$  sorted such that  $y_1 \triangleleft_{P,\sigma} \dots \triangleleft_{P,\sigma} y_{\#T_B[i]}$ . Given  $k \in \llbracket i'', j'' \rrbracket$ , we define  $w_P[k]$  as the string  $m_P[\text{select}_\$(m_P, \text{rank}_\$(m_P, \text{SA}(P)[k] - 1)) + 1 : \text{SA}(P)[k] + \text{LRS}(P)[k] - 1]$ . For all  $k \in \llbracket i'', j'' \rrbracket$ , the string  $w_P[k]$  is a string of  $P.S$ . Moreover, the definition of  $\triangleleft_{P,\sigma}$  implies  $w_P[i''] \triangleleft_{P,\sigma} \dots \triangleleft_{P,\sigma} w_P[j'']$ . As for  $l \in \llbracket 1, \#T_B[i] \rrbracket$ , the string  $\pi(y_l)$  is also a string of  $P.S$ , we get  $\pi(y_l) = w_P[i'' + l - 1]$ . Given  $x$  and  $y$  in  $\llbracket i'', j'' \rrbracket \in \text{Decomp\_BWT}(P)$ , we obtain the following equivalences between orders:

$$x < y \Leftrightarrow w_P[x] \triangleleft_{P,\sigma} w_P[y] \Leftrightarrow \pi(y_{x-i''+1}) \triangleleft_{P,\sigma} \pi(y_{y-i''+1}) \Leftrightarrow y_{x-i''+1} \triangleleft_{P,\sigma} y_{y-i''+1}.$$

By the definition of  $\text{BWT}(P)$ , for all  $k \in \llbracket 1, \#T_B[i] \rrbracket$ , we get that

$$\begin{aligned} \text{BWT}(P)[i'' + k - 1] &= m_P[\text{SA}(m_P)[i'' + k - 1] - 1] \\ &= w_P[i'' + k - 1][|w_P[i'' + k - 1]| - \text{LRS}[i'' + k - 1] + 1] \\ &= \pi(y_k)[|w_P[i'' + k - 1]| - \text{LRS}[i'' + k - 1] + 1] \\ &= \delta[\text{Child}_{\text{ACT}(\overleftarrow{P})(y_k)}(\text{Parent}_{\text{ACT}(\overleftarrow{P})}(\text{PA}(\text{ACT}(\overleftarrow{P}))[i']))]. \end{aligned} \quad \blacktriangleleft$$

#### A.4 Proof of Corollary 6

Let  $P$  be an ordered set of strings. To compute tables  $\text{XBWT}(\overleftarrow{\text{ACT}}(\overleftarrow{P}))$  and  $\text{XBWL}(\overleftarrow{\text{ACT}}(\overleftarrow{P}))$  using only  $\text{BWT}(P)$ ,  $\text{LCP}(P)$  and  $\text{LRS}(P)$ , we first define a new table  $\text{BWD}(P)$ .

The *Burrows-Wheeler Decomposition* of  $P$ , denoted by  $\text{BWD}(P)$ , is the array of length  $\#\text{Decomp\_BWT}(P)$  such that for each position  $i$ ,  $\text{BWD}(P)[i]$  is the cardinality of the  $i^{\text{th}}$  element of  $\text{Decomp\_BWT}(P)$  in interval order.

► **Lemma 18.** *Using tables  $\text{LCP}(P)$  and  $\text{LRS}(P)$ , Algorithm 3 computes  $\text{BWD}(P)$  in linear time in  $\|P.S\|$  and the table  $\text{BWD}(P)$  can be stored with  $\|P.S\| \times \log(\#(P.S))$  bits.*

**Proof of Lemma 18.** For each  $i$  in  $\llbracket 1, |m_P| \rrbracket$ , at the beginning of the loop **for**, we have that  $\text{LCP}(P)[i-j] \neq \text{LRS}(P)[i-j]$  and for all  $k \in \llbracket i-j+1, i-1 \rrbracket$ ,  $\text{LCP}(P)[k] \neq \text{LRS}(P)[k]$ . Hence, if  $\text{LCP}(P)[i] \neq \text{LRS}(P)[i]$ , the interval  $\llbracket i-j, i-1 \rrbracket$  is an element of  $\text{Decomp\_BWT}(P)$  and the cardinality of  $\llbracket i-j, i-1 \rrbracket$  is  $j$ . Otherwise, we increase  $j$  by 1 because the position  $i$  does not correspond to a new interval of  $\text{Decomp\_BWT}(P)$ . For the complexity, as each step of the loop can be computed in constant time, Algorithm 3 computes  $\text{BWD}(P)$  in linear time in  $\|P.S\|$ . As for each position  $i$  of  $\text{BWD}(P)$ ,  $\text{BWD}(P)[i]$  represents the number of strings of  $P.S$  having as suffix  $\text{Dec\_Pre}[u]$ , where  $u$  is the  $i^{\text{th}}$  element of  $\text{Decomp\_BWT}(P)$  sorted in interval order, it follows that  $\text{BWD}(P)[i] \leq \#(P.S)$ . This concludes the proof. ◀

---

#### Algorithm 3: Computation of table $\text{BWD}(P)$ .

---

**Input** : The tables  $\text{LCP}(P)$  and  $\text{LRS}(P)$

**Output**: The table  $\text{BWD}(P)$

$\text{BWD} \leftarrow$  empty list;

$j \leftarrow 0$ ;

**for**  $i \in \llbracket 1, |m_P| \rrbracket$  **do**

**if**  $\text{LCP}(P)[i] \neq \text{LRS}(P)[i]$  **then**

        insert  $j$  at the end of  $\text{BWD}$ ;

$j \leftarrow 1$ ;

**else**  $j \leftarrow j + 1$  ;

**return**  $\text{BWD}$ ;

---

► **Lemma 19.** *Using tables  $\text{BWT}(P)$  and  $\text{BWD}(P)$ , Algorithm 4 computes the tables  $\text{XBWT}(\overleftarrow{\text{ACT}}(\overleftarrow{P}))$  and  $\text{XBWL}(\overleftarrow{\text{ACT}}(\overleftarrow{P}))$  in linear time of  $\|P.S\| \times \#\Sigma$ .*

**Proof of Lemma 19.** Algorithm 4 is an application of Theorem 5. ◀

Using Algorithm 3 to build  $\text{BWD}(P)$  and Algorithm 4, we can compute tables  $\text{XBWT}(\overleftarrow{\text{ACT}}(\overleftarrow{P}))$  and  $\text{XBWL}(\overleftarrow{\text{ACT}}(\overleftarrow{P}))$  in linear time of  $\|P.S\| \times \#\Sigma$ .

We define the equivalent of  $\text{BWD}$  for  $\text{Decomp\_XBW}(P)$ . The *eXtended Burrows Wheeler Decomposition* ( $\text{XBWD}$ ) of a tree  $\mathcal{T}$  is the array of length  $\#\text{XBWL}(\mathcal{T})$  such that for any position  $i$ ,  $\text{XBWD}(P)[i]$  equals the cardinality of the  $i^{\text{th}}$  element of  $\text{Decomp\_XBW}(\mathcal{T})$  sorted in interval order.

► **Lemma 20.** *From the table  $\text{XBWL}(\text{ACT}(S))$ , Algorithm 5 computes  $\text{XBWD}(\text{ACT}(S))$  in linear time in  $\|P.S\|$  and the table  $\text{XBWD}(\text{ACT}(S))$  can be stored in  $\|P.S\| \times \log(\#\Sigma)$  bits.*

---

**Algorithm 4:** Computation of tables  $\text{XBWT}(\overleftarrow{P})$  and  $\text{XBWL}(\overleftarrow{P})$ .

---

**Input** : The tables  $\text{BWT}(P)$  and  $\text{BWD}(P)$   
**Output** : The tables  $\text{XBWT}(\overleftarrow{P})$  and  $\text{XBWL}(\overleftarrow{P})$   
 Initialise both  $\text{XBWT}$  and  $\text{XBWL}$  as empty lists;  
 $nb \leftarrow 1$ ;  
**for**  $i \in \llbracket 1, \#\text{BWD}(P) \rrbracket$  **do**  
    $D \leftarrow$  Dictionary such that for all  $l \in \Sigma$ ,  $D[l] \leftarrow \text{true}$ ;  
    $begin \leftarrow nb$ ;  
    $last \leftarrow begin + \text{BWD}(P)[i] - 1$ ;  
    $nb \leftarrow last + 1$ ;  
   **for**  $j \in \llbracket begin, last \rrbracket$  **do**  
     **if**  $D[\text{BWT}(P)[j]]$  **then**  
       insert  $\text{BWT}(P)[j]$  at the end of  $\text{XBWT}$  ;  
       insert 0 at the end of  $\text{XBWL}$  ;  
        $D[\text{BWT}(P)[j]] \leftarrow \text{false}$ ;  
    $\text{XBWL}[\#\text{XBWL}] \leftarrow 1$ ;  
**return**  $\text{XBWT}$  and  $\text{XBWL}$ ;

---

**Proof.** The proof of Lemma 20 is similar to that of Lemma 18. As for each position  $i$  of  $\text{XBWD}(\text{ACT}(S))$ ,  $\text{XBWD}(\text{ACT}(S))[i]$  represents the number of the right extensions of the strings  $\text{PA}(\text{ACT}(S))[i]$  in  $S$  (i.e., the number of different strings  $\text{PA}(\text{ACT}(S))[i]a$  which are substrings of a string of  $S$  for some  $a \in \Sigma$ ), we have  $\text{XBWD}(\text{ACT}(S))[i] \leq \#\Sigma$ . ◀

---

**Algorithm 5:** Computation of table  $\text{XBWD}(\text{ACT}(S))$ .

---

**Input** : The table  $\text{XBWL}(\text{ACT}(S))$ ;   **Output** : The table  $\text{XBWD}(\text{ACT}(S))$   
 $\text{XBWD} \leftarrow$  empty list;  
 $j \leftarrow 1$ ;  
**for**  $i \in \llbracket 1, \#(\text{XBWL}(\text{ACT}(S))) \rrbracket$  **do**  
   **if**  $\text{XBWL}(\text{ACT}(S)) = 1$  **then**  
     insert  $j$  at the end of  $\text{XBWD}$  ;  
      $j \leftarrow 1$ ;  
   **else**  $j \leftarrow j + 1$  ;  
**return**  $\text{XBWD}$ ;

---

► **Lemma 21.** *Using tables  $\text{XBWT}(\text{ACT}(S))$  and  $\text{XBWD}(\text{ACT}(S))$ , we can build the tables  $\text{BWT}(\overleftarrow{P})$ ,  $\text{LCP}(\overleftarrow{P})$  and  $\text{LRS}(\overleftarrow{P})$  in linear time of  $\|P.S\| \times \#\Sigma$ , where  $P$  is a topologically planar, ordered set of strings such that  $P.S = S$ .*

**Proof.** In [9], Ferragina *et al.* prove that with both tables  $\text{XBWT}(\text{ACT}(S))$  and  $\text{XBWD}(\text{ACT}(S))$  one can access in constant time the children and the parents in  $\text{ACT}(S)$ . Hence, we can compute in linear time in  $\|P.S\|$ , the table  $\text{TL}(\text{ACT}(S))$ , where in each position of  $i$  we store the number of leaves in the subtree of the node  $\text{PA}(\text{ACT}(S))[i]$ . We finish the proof using the results of Theorem 5 and an algorithm similar to Algorithm 4. ◀

### A.5 Proof of Lemma 10

**Proof.** Let  $T'$  be an optimal solution of `Min-Permutation-Table` for  $\text{XBWT}(\text{ACT}(P))$  and for  $\text{Decomp\_XBW}(\text{ACT}(P))$ . By Theorem 5, for each  $\llbracket i, j \rrbracket \in \text{Decomp\_XBW}(\text{ACT}(P))$ , the order of the symbols in  $\text{XBWT}(\text{ACT}(P))[i, j]$  depends on the order on the children of the parent of  $\text{PA}(\text{ACT}(P))[i]$ . Hence, the choice of  $T'$  corresponds to the choice of an order for each internal node of  $\text{ACT}(P)$  over all its children. As we can extend this order to a total order on the leaves of  $\text{ACT}(P)$ , we can build  $P'$  the ordered set of strings satisfying  $P'.S = S^\$$  and  $P'.\sigma(\pi(f_i)) = s_i$ , where the order on the leaves of  $\text{ACT}(P)$  is  $f_1 < \dots < f_{\#S}$  and  $s_1 < \dots < s_{\#S}$  are the strings of  $S$  in lexicographic order. ◀