

# Lightweight merging of compressed indices based on BWT variants<sup>☆</sup>

Lavinia Egidi<sup>a</sup>, Giovanni Manzini<sup>b</sup>

<sup>a</sup>University of Eastern Piedmont, Italy

<sup>b</sup>University of Eastern Piedmont and IIT-CNR Pisa, Italy

---

## Abstract

In this paper we propose a flexible and lightweight technique for merging compressed indices based on variants of Burrows-Wheeler transform (BWT), thus addressing the need for algorithms that compute compressed indices over large collections using a limited amount of working memory. Merge procedures make it possible to use an incremental strategy for building large indices based on merging indices for progressively larger subcollections.

Starting with a known lightweight algorithm for merging BWTs [Holt and McMillan, Bionformatics 2014], we show how to modify it in order to merge, or compute from scratch, also the Longest Common Prefix (LCP) array. We then expand our technique for merging compressed tries and circular/permuterm compressed indices, two compressed data structures for which there were hitherto no known merging algorithms.

*Keywords:* multi-string BWT, Longest Common Prefix array, XBWT, trie compression, compressed permuterm index, circular patterns.

---

## 1. Introduction

The Burrows Wheeler transform (BWT), originally introduced as a tool for data compression [1], has found application in the compact representation of many different data structures. After the seminal work [2] showing that the BWT can be used as a compressed full text index for a single string, many researchers have proposed variants of this transformation for string collections [3, 4, 5, 6, 7], trees [8, 9], graphs [10, 11, 12], and alignments [13, 14]. See [15] for an attempt to provide a unified view of these variants.

---

<sup>☆</sup>Partially supported by INdAM-GNCS project 2019 *Innovative methods for the solution of medical and biological big data*, MIUR-PRIN project *Multicriteria Data Structures and Algorithms: from compressed to learned indexes, and beyond* grant n. 2017WR7SHH, and the projects *LSBC\_19-21* and *HySecEn* from the University of Eastern Piedmont. Postprint version. The final publication is available on ScienceDirect <https://doi.org/10.1016/j.tcs.2019.11.001> © 2019. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0>.

In this paper we consider the problem of constructing compressed indices for *string collections* based on BWT variants. A compressed index is obviously most useful when working with very large amounts of data. Therefore, a fundamental requirement for construction algorithms, in order to be of practical use, is that they are *lightweight* in the sense that they use a limited amount of working space, i.e. space in addition to the space used for the input and the output. Indeed, the construction of compressed indices in linear time and small working space is an active and promising area of research, see [16, 17, 18] and references therein.

A natural approach when working with string collections is to build the indexing data structure *incrementally*, that is, for progressively larger subcollections [11, 19, 20, 21]. For example, when additional data should be added to an already large index, the incremental construction appears much more reasonable, and often works better in practice, than rebuilding the complete index from scratch, even when the from-scratch option has better theoretical bounds.

Along this path, Holt and McMillan [22, 23] proposed a simple and elegant algorithm, that we call the H&M algorithm from now on, for merging BWTs of collections of sequences. For collections of total size  $n$ , their fastest version takes  $\mathcal{O}(n \text{aveLcp}_{01})$  time where  $\text{aveLcp}_{01}$  is the average length of the longest common prefix between suffixes in the collection. The average length of the longest common prefix is  $\mathcal{O}(n)$  in the worst case but  $\mathcal{O}(\log n)$  for random strings and for many real world datasets [24]. However, even when  $\text{aveLcp}_{01} = \mathcal{O}(\log n)$  the H&M algorithm is not theoretically optimal since computing the BWT from scratch takes  $\mathcal{O}(n)$  time. Despite its theoretical shortcomings, because of its simplicity and small space usage, the H&M algorithm is competitive in practice for collections with relatively small average LCP. In addition, since the H&M algorithm accesses all data by sequential scans, it has been adapted to work on very large collections in external memory [23].

In this paper we revisit the H&M algorithm and we show that its main technique can be adapted to solve the merging problem for three different compressed indices based on the BWT.

First, in Section 4 we describe a procedure to merge, in addition to the BWTs, the Longest Common Prefix (LCP) arrays of string collections. The LCP array is often used to provide additional functionalities to indices based on the BWT [25], and the issue of efficiently computing and storing LCP values has recently received much attention [26, 27, 28, 29]. Our algorithm has the same  $\mathcal{O}(n \text{aveLcp})$  complexity as the H&M algorithm.

Next, in Section 5 we describe a procedure for merging *compressed labelled trees* (tries) as produced by the eXtended BWT transform (XBWT) [8, 9]. This result is particularly interesting since at the moment there are no time and space optimal algorithms for the computation from scratch of the XBWT. Our algorithm takes time proportional to the number of nodes in the output tree times the *average* node height.

Finally, in Section 6 we describe algorithms for merging *compressed indices for circular patterns* [30], and *compressed permuterm indices* [7]. The time complexity of these algorithms is proportional to the total collection size times the

average *circular LCP*, a notion that naturally extends the LCP to the modified lexicographic order used for circular strings.

Our algorithms are based on the H&M technique specialized to the particular features of the different compressed indices given as input. They all make use of techniques to recognize blocks of the input that become irrelevant for the computation and skip them in successive iterations. Because of the skipping of irrelevant blocks we call our merging procedures **Gap** algorithms. Our algorithms are all lightweight in the sense that, in addition to the input and the output, they use only a few bitarrays of working space and the space for handling the irrelevant blocks. The latter amount of space can be significant for pathological inputs, but in practice we found it takes between 2% and 9% of the overall space, depending on the alphabet size.

The **Gap** algorithms share with the H&M algorithm the feature of accessing all data by sequential scans and are therefore suitable for implementation in external memory. In [31] an external memory version of the **Gap** algorithm for merging BWT and LCP arrays is engineered, analyzed, and extensively tested on collections of DNA sequences. The results reported there show that the external memory version of **Gap** outperforms the known external memory algorithms for BWT/LCP computation when the average LCP of the collection is relatively small or when the strings of the input collection have widely different lengths.

To the best of our knowledge, the problem of incrementally building compressed indices via merging has been previously addressed only in [11, 19, 20, 21]. Sirén presents in [20, 21] an algorithm that maintains a BWT-based compressed index in RAM and incrementally merges new collections to it. The algorithm is the first that makes it possible to build indices for Terabytes of data without using a specialized machine with a lot of RAM. However, Sirén’s algorithm is specific for a particular compressed index (which doesn’t use the LCP array), while ours can be more easily adapted to build different flavors of compressed indices as shown in this paper. In [19], Li proposes an algorithm for the incremental computation of the BWT of string collections specialized to DNA sequences. The algorithm makes use of B+ trees, has a complexity of  $\mathcal{O}(n \log n)$ , and is probably the current fastest algorithm for the construction of an FM-index for a collection of long sequence reads. In [11] the authors present a merge algorithm for colored de Bruijn graphs. Their algorithm is also inspired by the H&M algorithm and the authors report a threefold reduction in working space compared to the state of the art methods for from scratch de Bruijn graphs. Inspired by the techniques introduced in this paper, in [32] we proposed an improved de Bruijn graph merging algorithm that also supports the construction of succinct Variable Order de Bruijn graph representations [33].

## 2. Background

Let  $t[1, n]$  denote a string of length  $n$  over an alphabet  $\Sigma$  of constant size  $\sigma$ . We write  $t[i, j]$  to denote the substring  $t[i]t[i + 1] \cdots t[j]$ . If  $j \geq n$  we assume  $t[i, j] = t[i, n]$ . If  $i > j$  or  $i > n$  then  $t[i, j]$  is the empty string. Given two strings  $t$  and  $s$  we write  $t \preceq s$  ( $t \prec s$ ) to denote that  $t$  is lexicographically (strictly)

lcp	bwt	context	lcp	bwt	context	id	lcp <sub>01</sub>	bwt <sub>01</sub>	context
-1	b	\$ <sub>0</sub>	-1	c	\$ <sub>1</sub>	0	-1	b	\$ <sub>0</sub>
0	c	ab\$ <sub>0</sub>	0	\$ <sub>1</sub>	aabcabc\$ <sub>1</sub>	1	0	c	\$ <sub>1</sub>
2	\$ <sub>0</sub>	abcab\$ <sub>0</sub>	1	c	abc\$ <sub>1</sub>	1	0	\$ <sub>1</sub>	aabcabc\$ <sub>1</sub>
0	a	b\$ <sub>0</sub>	3	a	abcabc\$ <sub>1</sub>	0	1	c	ab\$ <sub>0</sub>
1	a	bcab\$ <sub>0</sub>	0	a	bc\$ <sub>1</sub>	1	2	c	abc\$ <sub>1</sub>
0	b	cab\$ <sub>0</sub>	2	a	bcabc\$ <sub>1</sub>	0	3	\$ <sub>0</sub>	abcab\$ <sub>0</sub>
-1			0	b	c\$ <sub>1</sub>	1	5	a	abcabc\$ <sub>1</sub>
			1	b	cab\$ <sub>1</sub>	0	0	a	b\$ <sub>0</sub>
			-1			1	1	a	bc\$ <sub>1</sub>
						0	2	a	bcab\$ <sub>0</sub>
						1	4	a	bcabc\$ <sub>1</sub>
						1	0	b	c\$ <sub>1</sub>
						0	1	b	cab\$ <sub>0</sub>
						1	3	b	cab\$ <sub>1</sub>
							-1		

Figure 1: LCP array and BWT for  $t_0 = \text{abcab}\$0$  and  $t_1 = \text{aabcabc}\$1$ , and multi-string BWT and corresponding LCP array for the same strings. Column *id* shows, for each entry of  $\text{bwt}_{01} = \text{bc}\$1\text{cc}\$0\text{aaaabbb}$  whether it comes from  $t_0$  or  $t_1$ .

smaller than  $s$ . We denote by  $\text{LCP}(t, s)$  the length of the longest common prefix between  $t$  and  $s$ .

The *suffix array*  $\text{sa}[1, n]$  associated to  $t$  is the permutation of  $[1, n]$  giving the lexicographic order of  $t$ 's suffixes, that is, for  $i = 1, \dots, n-1$ ,  $t[\text{sa}[i], n] < t[\text{sa}[i+1], n]$ . The *longest common prefix array*  $\text{lcp}[1, n+1]$  is defined for  $i = 2, \dots, n$  by

$$\text{lcp}[i] = \text{LCP}(t[\text{sa}[i-1], n], t[\text{sa}[i], n]); \quad (1)$$

the *lcp* array stores the length of the longest common prefix between lexicographically consecutive suffixes. For convenience we define  $\text{lcp}[1] = \text{lcp}[n+1] = -1$ . We also define the maximum and average LCP as:

$$\text{maxLcp} = \max_{1 < i \leq n} \text{lcp}[i], \quad \text{aveLcp} = \left( \sum_{1 < i \leq n} \text{lcp}[i] \right) / n. \quad (2)$$

The *Burrows-Wheeler transform*  $\text{bwt}[1, n]$  of  $t$  is defined by

$$\text{bwt}[i] = \begin{cases} t[n] & \text{if } \text{sa}[i] = 1 \\ t[\text{sa}[i] - 1] & \text{if } \text{sa}[i] > 1. \end{cases}$$

$\text{bwt}$  is best seen as the permutation of  $t$  in which the position of  $t[j]$  coincides with the lexicographic rank of  $t[j+1, n]$  (or of  $t[1, n]$  if  $j = n$ ) in the suffix array. We call the string  $t[j+1, n]$  *context* of  $t[j]$ . See Figure 1 for an example.

The longest common prefix (LCP) array, and Burrows-Wheeler transform (BWT) can be generalized to the case of multiple strings. Historically, the

first of such generalizations is the circular BWT [4] considered in Section 6. Here we consider the generalization proposed in [3, 5] which is the one most used in applications. Let  $\mathbf{t}_0[1, n_0]$  and  $\mathbf{t}_1[1, n_1]$  be such that  $\mathbf{t}_0[n_0] = \$_0$  and  $\mathbf{t}_1[n_1] = \$_1$  where  $\$_0 < \$_1$  are two symbols not appearing elsewhere in  $\mathbf{t}_0$  and  $\mathbf{t}_1$  and smaller than any other symbol. Let  $\mathbf{sa}_{01}[1, n_0 + n_1]$  denote the suffix array of the concatenation  $\mathbf{t}_0\mathbf{t}_1$ . The *multi-string* BWT of  $\mathbf{t}_0$  and  $\mathbf{t}_1$ , denoted by  $\mathbf{bwt}_{01}[1, n_0 + n_1]$ , is defined by

$$\mathbf{bwt}_{01}[i] = \begin{cases} \mathbf{t}_0[n_0] & \text{if } \mathbf{sa}_{01}[i] = 1 \\ \mathbf{t}_0[\mathbf{sa}_{01}[i] - 1] & \text{if } 1 < \mathbf{sa}_{01}[i] \leq n_0 \\ \mathbf{t}_1[n_1] & \text{if } \mathbf{sa}_{01}[i] = n_0 + 1 \\ \mathbf{t}_1[\mathbf{sa}_{01}[i] - n_0 - 1] & \text{if } n_0 + 1 < \mathbf{sa}_{01}[i]. \end{cases}$$

In other words,  $\mathbf{bwt}_{01}[i]$  is the symbol preceding the  $i$ -th lexicographically larger suffix, with the exception that if  $\mathbf{sa}_{01}[i] = 1$  then  $\mathbf{bwt}_{01}[i] = \$_0$  and if  $\mathbf{sa}_{01}[i] = n_0 + 1$  then  $\mathbf{bwt}_{01}[i] = \$_1$ . Hence,  $\mathbf{bwt}_{01}[i]$  is always a character of the string ( $\mathbf{t}_0$  or  $\mathbf{t}_1$ ) containing the  $i$ -th largest suffix (see again Figure 1). The above notion of multi-string BWT can be immediately generalized to define  $\mathbf{bwt}_{1\dots k}$  for a family of distinct strings  $\mathbf{t}_1, \dots, \mathbf{t}_k$ . Essentially  $\mathbf{bwt}_{1\dots k}$  is a permutation of the symbols in  $\mathbf{t}_1, \dots, \mathbf{t}_k$  such that the position in  $\mathbf{bwt}_{1\dots k}$  of  $\mathbf{t}_i[j]$  is given by the lexicographic rank of its context  $\mathbf{t}_i[j + 1, n_i]$  (or  $\mathbf{t}_i[1, n_i]$  if  $j = n_i$ ).

Given the concatenation  $\mathbf{t}_0\mathbf{t}_1$  and its suffix array  $\mathbf{sa}_{01}[1, n_0 + n_1]$ , we consider the corresponding LCP array  $\mathbf{lcp}_{01}[1, n_0 + n_1 + 1]$  defined as in (1) (see again Figure 1). Note that, for  $i = 2, \dots, n_0 + n_1$ ,  $\mathbf{lcp}_{01}[i]$  gives the length of the longest common prefix between the contexts of  $\mathbf{bwt}_{01}[i]$  and  $\mathbf{bwt}_{01}[i - 1]$ . This definition can be immediately generalized to a family of  $k$  strings to define the LCP array  $\mathbf{lcp}_{1\dots k}$  associated to the multi-string BWT  $\mathbf{bwt}_{1\dots k}$ .

### 2.1. The H&M Algorithm

In [22] Holt and McMillan introduced a simple and elegant algorithm, we call it the H&M algorithm, to merge multi-string BWTs<sup>1</sup>. Because it is the starting point for our results, we now briefly recall its main properties.

Given  $\mathbf{bwt}_{1\dots k}$  and  $\mathbf{bwt}_{k+1\dots h}$  the H&M algorithm computes  $\mathbf{bwt}_{1\dots h}$ . The computation doesn't explicitly need  $\mathbf{t}_1, \dots, \mathbf{t}_h$  but only the (multi-string) BWTs to be merged. For simplicity of notation we describe the algorithm assuming we are merging two single-string BWTs  $\mathbf{bwt}_0 = \mathbf{bwt}(\mathbf{t}_0)$  and  $\mathbf{bwt}_1 = \mathbf{bwt}(\mathbf{t}_1)$ ; the same algorithm works in the general case with multi-string BWTs in input. Note also that the algorithm can be easily adapted to merge more than two (multi-string) BWTs at the same time.

Computing  $\mathbf{bwt}_{01}$  amounts to sorting the symbols of  $\mathbf{bwt}_0$  and  $\mathbf{bwt}_1$  according to the lexicographic order of their contexts, where the context of symbol  $\mathbf{bwt}_0[i]$  (resp.  $\mathbf{bwt}_1[i]$ ) is  $\mathbf{t}_0[\mathbf{sa}_0[i], n_0]$  (resp.  $\mathbf{t}_1[\mathbf{sa}_1[i], n_1]$ ). By construction, the symbols

<sup>1</sup>Unless explicitly stated otherwise, in the following we use H&M to refer to the algorithm from [22], and not to its variant proposed in [23].

in  $\mathbf{bwt}_0$  and  $\mathbf{bwt}_1$  are already sorted by context, hence to compute  $\mathbf{bwt}_{01}$  we only need to merge  $\mathbf{bwt}_0$  and  $\mathbf{bwt}_1$  without changing the relative order of the symbols within the two input sequences.

The H&M algorithm works in successive iterations. After the  $h$ -th iteration the entries of  $\mathbf{bwt}_0$  and  $\mathbf{bwt}_1$  are sorted on the basis of the first  $h$  symbols of their context. More formally, the output of the  $h$ -th iteration is a binary vector  $Z^{(h)}$  containing  $n_0 = |\mathbf{t}_0|$   $\mathbf{0}$ 's and  $n_1 = |\mathbf{t}_1|$   $\mathbf{1}$ 's and such that the following property holds.

**Property 1.** *For  $i = 1, \dots, n_0$  and  $j = 1, \dots, n_1$  the  $i$ -th  $\mathbf{0}$  precedes the  $j$ -th  $\mathbf{1}$  in  $Z^{(h)}$  if and only if*

$$\mathbf{t}_0[\mathbf{sa}_0[i], \mathbf{sa}_0[i] + h - 1] \preceq \mathbf{t}_1[\mathbf{sa}_1[j], \mathbf{sa}_1[j] + h - 1] \quad (3)$$

(recall that according to our notation if  $\mathbf{sa}_0[i] + h - 1 > n_0$  then  $\mathbf{t}_0[\mathbf{sa}_0[i], \mathbf{sa}_0[i] + h - 1]$  coincides with  $\mathbf{t}_0[\mathbf{sa}_0[i], n_0]$ , and similarly for  $\mathbf{t}_1$ ).  $\square$

Following Property 1 we identify the  $i$ -th  $\mathbf{0}$  in  $Z^{(h)}$  with  $\mathbf{bwt}_0[i]$  and the  $j$ -th  $\mathbf{1}$  in  $Z^{(h)}$  with  $\mathbf{bwt}_1[j]$  so that  $Z^{(h)}$  encodes a permutation of  $\mathbf{bwt}_{01}$ . Property 1 is equivalent to stating that we can logically partition  $Z^{(h)}$  into  $b(h) + 1$  blocks

$$Z^{(h)}[1, \ell_1], Z^{(h)}[\ell_1 + 1, \ell_2], \dots, Z^{(h)}[\ell_{b(h)} + 1, n_0 + n_1] \quad (4)$$

such that each block corresponds to a set of  $\mathbf{bwt}_{01}$  symbols whose contexts are prefixed by the same length- $h$  string (the symbols with a context shorter than  $h$  are contained in singleton blocks). Within each block the symbols of  $\mathbf{bwt}_0$  precede those of  $\mathbf{bwt}_1$ , and the context of any symbol in block  $Z^{(h)}[\ell_j + 1, \ell_{j+1}]$  is lexicographically smaller than the context of any symbol in block  $Z^{(h)}[\ell_k + 1, \ell_{k+1}]$  with  $k > j$ .

The H&M algorithm initially sets  $Z^{(0)} = \mathbf{0}^{n_0} \mathbf{1}^{n_1}$ : since the context of every  $\mathbf{bwt}_{01}$  symbol is prefixed by the same length-0 string (the empty string), there is a single block containing all  $\mathbf{bwt}_{01}$  symbols. At iteration  $h$  the algorithm computes  $Z^{(h+1)}$  from  $Z^{(h)}$  using the procedure in Figure 2. The following lemma is a restatement of Lemma 3.2 in [22] using our notation (see [34] for a proof in our notation).

**Lemma 1.** *For  $h = 0, 1, 2, \dots$  the bit vector  $Z^{(h)}$  satisfies Property 1.  $\square$*

### 3. Computing LCP values with the H&M algorithm

Our first result is to show that with a simple modification to the H&M algorithm it is possible to compute the LCP array  $\mathbf{lcp}_{01}$ , in addition to merging  $\mathbf{bwt}_0$  and  $\mathbf{bwt}_1$ . Our strategy consists in keeping explicit track of the logical blocks we have defined for  $Z^{(h)}$  and represented in (4). We maintain an integer array  $B[1, n_0 + n_1 + 1]$  such that at the end of iteration  $h$  it is  $B[i] \neq 0$  if and only if a block of  $Z^{(h)}$  starts at position  $i$ . The use of such integer array is shown

---

```

1: Initialize array  $F[1, \sigma]$ 
2:  $k_0 \leftarrow 1; k_1 \leftarrow 1$  ▷ Init counters for  $\text{bwt}_0$  and  $\text{bwt}_1$ 
3: for  $k \leftarrow 1$  to  $n_0 + n_1$  do
4:    $b \leftarrow Z^{(h-1)}[k]$  ▷ Read bit  $b$  from  $Z^{(h-1)}$ 
5:    $c \leftarrow \text{bwt}_b[k_b++]$  ▷ Get symbol from  $\text{bwt}_0$  or  $\text{bwt}_1$  according to  $b$ 
6:   if  $c \neq \$$  then
7:      $j \leftarrow F[c]++$  ▷ Get destination for  $b$  according to symbol  $c$ 
8:   else
9:      $j \leftarrow b$  ▷ Symbol  $\$b$  goes to position  $b$ 
10:  end if
11:   $Z^{(h)}[j] \leftarrow b$  ▷ Copy bit  $b$  to  $Z^{(h)}$ 
12: end for

```

---

Figure 2: Main loop of algorithm H&M for computing  $Z^{(h)}$  given  $Z^{(h-1)}$ . Array  $F$  is initialized so that  $F[c]$  contains the number of occurrences of symbols smaller than  $c$  in  $\text{bwt}_0$  and  $\text{bwt}_1$  plus one. Note that the bits stored in  $Z^{(h)}$  immediately after reading symbol  $c \neq \$$  are stored in positions from  $F[c]$  to  $F[c+1] - 1$  of  $Z^{(h)}$ .

---

```

1: Initialize arrays  $F[1, \sigma]$  and  $\text{Block\_id}[1, \sigma]$ 
2:  $k_0 \leftarrow 1; k_1 \leftarrow 1$  ▷ Init counters for  $\text{bwt}_0$  and  $\text{bwt}_1$ 
3: for  $k \leftarrow 1$  to  $n_0 + n_1$  do
4:   if  $B[k] \neq 0$  and  $B[k] \neq h$  then
5:      $\text{id} \leftarrow k$  ▷ A new block of  $Z^{(h-1)}$  is starting
6:   end if
7:    $b \leftarrow Z^{(h-1)}[k]$  ▷ Read bit  $b$  from  $Z^{(h-1)}$ 
8:    $c \leftarrow \text{bwt}_b[k_b++]$  ▷ Get symbol from  $\text{bwt}_0$  or  $\text{bwt}_1$  according to  $b$ 
9:   if  $c \neq \$$  then
10:     $j \leftarrow F[c]++$  ▷ Get destination for  $b$  according to symbol  $c$ 
11:  else
12:     $j \leftarrow b$  ▷ Symbol  $\$b$  goes to position  $b$ 
13:  end if
14:   $Z^{(h)}[j] \leftarrow b$  ▷ Copy bit  $b$  to  $Z^{(h)}$ 
15:  if  $\text{Block\_id}[c] \neq \text{id}$  then
16:     $\text{Block\_id}[c] \leftarrow \text{id}$  ▷ Update block id for symbol  $c$ 
17:    if  $B[j] = 0$  then ▷ Check if already marked
18:       $B[j] = h$  ▷ A new block of  $Z^{(h)}$  will start here
19:    end if
20:  end if
21: end for

```

---

Figure 3: Main loop of the H&M algorithm modified for the computation of the lcp values. At Line 1 for each symbol  $c$  we set  $\text{Block\_id}[c] = -1$  and  $F[c]$  as in Figure 2. At the beginning of the algorithm we initialize the array  $B[1, n_0 + n_1 + 1]$  as  $B = 1 0^{n_0 + n_1 - 1} 1$ .

in Figure 3. Note that: (i) initially we set  $B = 1 0^{n_0+n_1-1} 1$  and once an entry in  $B$  becomes nonzero it is never changed, (ii) during iteration  $h$  we only write to  $B$  the value  $h$ , (iii) because of the test at Line 4 the values written during iteration  $h$  influence the algorithm only in subsequent iterations. In order to identify new blocks, we maintain an array  $\text{Block\_id}[1, \sigma]$  such that  $\text{Block\_id}[c]$  is the id of the block of  $Z^{(h-1)}$  to which the last seen occurrence of symbol  $c$  belonged.

The following lemma shows that the nonzero values of  $B$  at the end of iteration  $h$  mark the boundaries of  $Z^{(h)}$ 's logical blocks.

**Lemma 2.** *For any  $h \geq 0$ , let  $\ell, m$  be such that  $1 \leq \ell \leq m \leq n_0 + n_1$  and*

$$\text{lcp}_{01}[\ell] < h, \quad \min(\text{lcp}_{01}[\ell + 1], \dots, \text{lcp}_{01}[m]) \geq h, \quad \text{lcp}_{01}[m + 1] < h. \quad (5)$$

*Then, at the end of iteration  $h$  the array  $B$  is such that*

$$B[\ell] \neq 0, \quad B[\ell + 1] = \dots = B[m] = 0, \quad B[m + 1] \neq 0 \quad (6)$$

*and  $Z^{(h)}[\ell, m]$  is one of the blocks in (4).*

**Proof:** We prove the result by induction on  $h$ . For  $h = 0$ , hence before the execution of the first iteration, (5) is only valid for  $\ell = 1$  and  $m = n_0 + n_1$  (recall that we defined  $\text{lcp}_{01}[1] = \text{lcp}_{01}[n_0 + n_1 + 1] = -1$ ). Since initially  $B = 1 0^{n_0+n_1-1} 1$  our claim holds.

Suppose now that (5) holds for some  $h > 0$ . Let  $s = \mathbf{t}_{01}[\mathbf{sa}_{01}[\ell], \mathbf{sa}_{01}[\ell] + h - 1]$ ; by (5)  $s$  is a common prefix of the suffixes starting at positions  $\mathbf{sa}_{01}[\ell], \mathbf{sa}_{01}[\ell + 1], \dots, \mathbf{sa}_{01}[m]$ , and no other suffix of  $\mathbf{t}_{01}$  is prefixed by  $s$ . By Property 1 the  $\mathbf{0}$ s and  $\mathbf{1}$ s in  $Z^{(h)}[\ell, m]$  correspond to the same set of suffixes, that is, if  $\ell \leq v \leq m$  and  $Z^{(h)}[v]$  is the  $i$ th  $\mathbf{0}$  (resp.  $j$ th  $\mathbf{1}$ ) of  $Z^{(h)}$  then the suffix starting at  $\mathbf{t}_0[\mathbf{sa}_0[i]]$  (resp.  $\mathbf{t}_1[\mathbf{sa}_1[j]]$ ) is prefixed by  $s$ .

To prove (6) we start by showing that, if  $\ell < m$ , then at the end of iteration  $h - 1$  it is  $B[\ell + 1] = \dots = B[m] = 0$ . To see this observe that the range  $\mathbf{sa}_{01}[\ell, m]$  is part of a (possibly) larger range  $\mathbf{sa}_{01}[\ell', m']$  containing all suffixes prefixed by the length  $h - 1$  prefix of  $s$ . By inductive hypothesis, at the end of iteration  $h - 1$  it is  $B[\ell' + 1] = \dots = B[m'] = 0$  which proves our claim since  $\ell' \leq \ell$  and  $m \leq m'$ .

To complete the proof, we need to show that during iteration  $h$ : (i) we do not modify  $B[\ell + 1, m]$  and (ii) we write a nonzero to  $B[\ell]$  and  $B[m + 1]$  if they do not already contain a nonzero. Let  $c = s[0]$  and  $s' = s[1, h - 1]$  so that  $s = cs'$ . Consider now the range  $\mathbf{sa}_{01}[e, f]$  containing the suffixes prefixed by  $s'$ . By inductive hypothesis at the end of iteration  $h - 1$  it is

$$B[e] \neq 0, \quad B[e + 1] = \dots = B[f] = 0, \quad B[f + 1] \neq 0. \quad (7)$$

During iteration  $h$ , the bits in  $Z^{(h)}[\ell, m]$  are possibly changed only when we are scanning the region  $Z^{(h-1)}[e, f]$  and we find an entry  $b = Z^{(h-1)}[k]$ ,  $e \leq k \leq f$ , such that the corresponding value in  $\text{bwt}_b$  is  $c$ . Note that by (7) as soon as  $k$



reaches  $e$  the variable `id` changes and becomes different from all values stored in `Block_id`. Hence, at the first occurrence of symbol  $c$  the value  $h$  will be stored in  $B[\ell]$  (Line 18) unless a nonzero is already there. Again, because of (7), during the scanning of  $Z^{(h-1)}[e, f]$  the variable `id` does not change so subsequent occurrences of  $c$  will not cause a nonzero value to be written to  $B[\ell + 1, m]$ . Finally, as soon as we leave region  $Z^{(h-1)}[e, f]$  and  $k$  reaches  $f + 1$ , the variable `id` changes again and at the next occurrence of  $c$  a nonzero value will be stored in  $B[m + 1]$ . If there are no more occurrences of  $c$  after we leave region  $Z^{(h-1)}[e, f]$  then either  $\text{sa}_{01}[m + 1]$  is the first suffix array entry prefixed by symbol  $c + 1$  or  $m + 1 = n_0 + n_1 + 1$ . In the former case  $B[m + 1]$  gets a nonzero value at iteration 1, in the latter case  $B[m + 1]$  gets a nonzero value when we initialize array  $B$ .  $\square$

**Corollary 1.** *For  $i = 2, \dots, n_0 + n_1$ , if  $\text{lcp}_{01}[i] = \ell$ , then starting from the end of iteration  $\ell + 1$  it is  $B[i] = \ell + 1$ .*

**Proof:** By Lemma 2 we know that  $B[i]$  becomes nonzero only after iteration  $\ell + 1$ . Since at the end of iteration  $\ell$  it is still  $B[i] = 0$  during iteration  $\ell + 1$   $B[i]$  gets the value  $\ell + 1$  which is never changed in successive iterations.  $\square$

The above corollary suggests the following algorithm to compute  $\text{bwt}_{01}$  and  $\text{lcp}_{01}$ : repeat the procedure of Figure 3 until the iteration  $h$  in which all entries in  $B$  become nonzero. At that point  $Z^{(h)}$  describes how  $\text{bwt}_0$  and  $\text{bwt}_1$  should be merged to get  $\text{bwt}_{01}$  and for  $i = 2, \dots, n_0 + n_1$   $\text{lcp}_{01}[i] = B[i] - 1$ . The above strategy requires a number of iterations, each one taking  $\mathcal{O}(n_0 + n_1)$  time, equal to the maximum of the `lcp` values, for an overall complexity of  $\mathcal{O}((n_0 + n_1) \max \text{lcp}_{01})$ , where  $\max \text{lcp}_{01} = \max_i \text{lcp}_{01}[i]$ . Note that in addition to the space for the input and the output the algorithm only uses two bit arrays (one for the current and the next  $Z^{(\cdot)}$ ) and a constant number of counters (the arrays  $F$  and `Block_id`). Summing up we have the following result.

**Lemma 3.** *Given  $\text{bwt}_0$  and  $\text{bwt}_1$ , the algorithm in Figure 3 computes  $\text{bwt}_{01}$  and  $\text{lcp}_{01}$  in  $\mathcal{O}(n \max \text{Lcp})$  time and  $2n + \mathcal{O}(\log n)$  bits of working space, where  $n = |\mathbf{t}_{01}|$  and  $\max \text{Lcp} = \max_i \text{lcp}_{01}[i]$  is the maximum LCP of  $\mathbf{t}_{01}$ .  $\square$*

#### 4. The Gap BWT/LCP merging Algorithm

The `Gap` algorithm, as well as its variants described in the following sections, are based on the notion of *monochrome blocks*.

**Definition 1.** *If  $B[\ell] \neq 0$ ,  $B[m + 1] \neq 0$  and  $B[\ell + 1] = \dots = B[m] = 0$ , we say that block  $Z^{(h)}[\ell, m]$  is monochrome if it contains only  $\mathbf{0}$ 's or only  $\mathbf{1}$ 's.  $\square$*

Since a monochrome block only contains suffixes from either  $\mathbf{t}_0$  or  $\mathbf{t}_1$ , whose relative order is known, it does not need to be further modified. If in addition, the LCP arrays of  $\mathbf{t}_0$  and  $\mathbf{t}_1$  are given in input, then also LCP values inside monochrome blocks are known without further processing. This intuition is formalized by the following lemmas.

**Lemma 4.** *If at the end of iteration  $h$  bit vector  $Z^{(h)}$  contains only monochrome blocks we can compute  $\mathbf{bwt}_{01}$  and  $\mathbf{lcp}_{01}$  in  $\mathcal{O}(n_0 + n_1)$  time from  $\mathbf{bwt}_0$ ,  $\mathbf{bwt}_1$ ,  $\mathbf{lcp}_0$  and  $\mathbf{lcp}_1$ .*

**Proof:** By Property 1, if we identify the  $i$ -th  $\mathbf{0}$  in  $Z^{(h)}$  with  $\mathbf{bwt}_0[i]$  and the  $j$ -th  $\mathbf{1}$  with  $\mathbf{bwt}_1[j]$  the only elements which could be not correctly sorted by context are those within the same block. However, if the blocks are monochrome all elements belong to either  $\mathbf{bwt}_0$  or  $\mathbf{bwt}_1$  so their relative order is correct.

To compute  $\mathbf{lcp}_{01}$  we observe that if  $B[i] \neq 0$  then by (the proof of) Corollary 1 it is  $\mathbf{lcp}_{01}[i] = B[i] - 1$ . If instead  $B[i] = 0$  we are inside a block hence  $\mathbf{sa}_{01}[i-1]$  and  $\mathbf{sa}_{01}[i]$  belong to the same string  $\mathbf{t}_0$  or  $\mathbf{t}_1$  and their LCP is directly available in  $\mathbf{lcp}_0$  or  $\mathbf{lcp}_1$ .  $\square$

Notice that a lazy strategy of not completely processing monochrome blocks, makes it impossible to compute LCP values from scratch. In this case, in order to compute  $\mathbf{lcp}_{01}$  it is necessary that the algorithm also takes  $\mathbf{lcp}_1$  and  $\mathbf{lcp}_0$  in input.

**Lemma 5.** *Suppose that, at the end of iteration  $h$ ,  $Z^{(h)}[\ell, m]$  is a monochrome block. Then (i) for  $g > h$ ,  $Z^{(g)}[\ell, m] = Z^{(h)}[\ell, m]$ , and (ii) processing  $Z^{(h)}[\ell, m]$  during iteration  $h + 1$  creates a set of monochrome blocks in  $Z^{(h+1)}$ .*

**Proof:** The first part of the Lemma follows from the observation that subsequent iterations of the algorithm will only reorder the values within a block (and possibly create new sub-blocks); but if a block is monochrome the reordering will not change its actual content.

For the second part, we observe that during iteration  $h + 1$  as  $k$  goes from  $\ell$  to  $m$  the algorithm writes to  $Z^{(h+1)}$  the same value which is in  $Z^{(h)}[\ell, m]$ . Hence, a new monochrome block will be created for each distinct symbol encountered (in  $\mathbf{bwt}_0$  or  $\mathbf{bwt}_1$ ) as  $k$  goes through the range  $[\ell, m]$ .  $\square$

The lemma implies that, if block  $Z^{(h)}[\ell, m]$  is monochrome at the end of iteration  $h$ , starting from iteration  $g = h + 2$  processing the range  $[\ell, m]$  will not change  $Z^{(g)}$  with respect to  $Z^{(g-1)}$ . Indeed, by the lemma the monochrome blocks created in iteration  $h + 1$  do not change in subsequent iterations (in a subsequent iteration a monochrome block can be split in sub-blocks, but the actual content of the bit vector does not change). The above observation suggests that, after we have processed block  $Z^{(h+1)}[\ell, m]$  in iteration  $h + 1$ , we can mark it as *irrelevant* and avoid to process it again. As the computation goes on, more and more blocks become irrelevant. Hence, at the generic iteration  $h$  instead of processing the whole  $Z^{(h-1)}$  we process only the blocks which are still “active” and skip irrelevant blocks. Adjacent irrelevant blocks are merged so that among two active blocks there is at most one irrelevant block (the *gap* after which the algorithm is named). The overall structure of a single iteration is shown in Figure 4. The algorithm terminates when there are no more active blocks since this implies that all blocks have become monochrome and by Lemma 4 we are able to compute  $\mathbf{bwt}_{01}$  and  $\mathbf{lcp}_{01}$ .

---

```

1: if (next block is irrelevant) then
2:   skip it
3: else
4:   process block
5:   if (processed block is monochrome) then
6:     mark it irrelevant
7:   end if
8: end if
9: if (last two blocks are irrelevant) then
10:  merge them
11: end if

```

---

Figure 4: Main loop of the **Gap** algorithm. The processing of active blocks at Line 4 is done as in Lines 7–20 of Figure 3.

We point out that at Line 2 of the **Gap** algorithm we cannot simply skip an irrelevant block ignoring its content. To keep the algorithm consistent we must correctly update the global variables of the main loop, i.e. the array  $F$  and the pointers  $k_0$  and  $k_1$  in Figure 3. To this end a simple approach is to store for each irrelevant block the number of occurrences  $o_c$  of each symbol  $c \in \Sigma$  in it and the pair  $(r_0, r_1)$  providing the number of **0**'s and **1**'s in the block (recall that an irrelevant block may consist of adjacent monochrome blocks coming from different strings). When the algorithm reaches an irrelevant block,  $F$ ,  $k_0$ ,  $k_1$  are updated setting  $k_0 \leftarrow k_0 + r_0$ ,  $k_1 \leftarrow k_1 + r_1$  and  $\forall c F[c] \leftarrow F[c] + o_c$ . The above scheme for handling irrelevant blocks is simple and effective for most applications. However, for a large non-constant alphabet it would imply a multiplicative  $\mathcal{O}(\sigma)$  slowdown. In [34, Sect. 4] we present a different scheme for large alphabets with a slowdown reduced to  $\mathcal{O}(\log \sigma)$ .

We point out that our **Gap** algorithm is related to the **H&M** variant with  $\mathcal{O}(n \text{ aveLcp})$  time complexity described in [23, Sect. 2.1]: Indeed, the sorting operations are essentially the same in the two algorithms. The main difference is that **Gap** keeps explicit track of the irrelevant blocks while **H&M** keeps explicit track of the active blocks (called buckets in [23]): this difference makes the non-sorting operations completely different. An advantage of working with irrelevant blocks is that they can be easily merged, while this is not the case for the active blocks in **H&M**. Of course, the main difference is that **Gap** merges simultaneously **BWT** and **LCP** values.

**Theorem 6.** *Given  $\text{bwt}_0, \text{lcp}_0$  and  $\text{bwt}_1, \text{lcp}_1$  let  $n = |\text{bwt}_0| + |\text{bwt}_1|$ . The **Gap** algorithm computes  $\text{bwt}_{01}$  and  $\text{lcp}_{01}$  in  $\mathcal{O}(n \text{ aveLcp}_{01})$  time, where  $\text{aveLcp}_{01} = (\sum_i \text{lcp}_{01}[i])/n$  is the average **LCP** of the string  $\text{t}_{01}$ . The working space is  $2n + \mathcal{O}(\log n)$  bits, plus the space used for handling irrelevant blocks.*

**Proof:** For the running time we reason as in [23] and observe that the sum, over all iterations, of the length of all active blocks is bounded by  $\mathcal{O}(\sum_i \text{lcp}_{01}[i]) = \mathcal{O}(n \text{ aveLcp}_{01})$ . The time bound follows observing that at any iteration the cost of processing an active block of length  $\ell$  is bounded by  $\mathcal{O}(\ell)$  time.

Name	Size GB	$\sigma$	Max Len	Ave Len	Max LCP	Ave LCP
Pacbio	6.24	5	40212	9567.43	1055	17.99
Illumina	7.60	6	103	102.00	102	27.53
Wiki-it	4.01	210	553975	4302.84	93537	61.02
Proteins	6.11	26	35991	410.22	25065	100.60

Table 1: Collections used in our experiments sorted by average LCP. Columns 4 and 5 refer to the lengths of the single documents. Pacbio are NGS reads from a *D.melanogaster* dataset. Illumina are NGS reads from Human ERA015743 dataset. Wiki-it are pages from Italian Wikipedia. Proteins are protein sequences from Uniprot. Collections and source files are available on <https://people.unipmn.it/manzini/gap>.

Name	$k$	gSACA-K + $\Phi$	$\tau = 50$		$\tau = 100$		$\tau = 200$	
			time	space	time	space	time	space
Pacbio	7	0.46	0.41	4.35	0.46	4.18	0.51	4.09
Illumina	4	0.48	0.93	3.31	1.02	3.16	1.09	3.08
Wiki-it	5	0.41	—	—	—	—	3.07	6.55
Proteins	4	0.59	3.90	4.55	5.18	4.29	7.05	4.15

Table 2: For each collection we report the number  $k$  of subcollections, the average running time of gSACA-K+ $\Phi$  in  $\mu$ secs per symbol, and the running time ( $\mu$ secs) and space usage (bytes) per symbol for Gap for different values of the  $\tau$  parameter. All tests were executed on a desktop with 32GB RAM and eight Intel-I7 3.40GHz CPUs, using a single CPU in each experiment.

For the analysis of the working space we observe for the array  $B$  we can use the space for the output LCP, hence the working space consists only in  $2n$  bits for two instances of the arrays  $Z^{(\cdot)}$  and a constant number of counters (the arrays  $F$  and  $\text{Block\_id}$ ).  $\square$

It is unfortunately impossible to give a clean bound for the space needed for keeping track of irrelevant blocks. Our scheme uses  $\mathcal{O}(1)$  words per block, but in the worst case we can have  $\Theta(n)$  blocks. Although such worst case is rather unlikely, it is important to have some form of control on this additional space. We use the following simple heuristic: we choose a threshold  $\tau$  and we keep track of an irrelevant block only if its size is at least  $\tau$ . This strategy introduces a  $\mathcal{O}(\tau)$  time slowdown but ensures that there are at most  $n/(\tau + 1)$  irrelevant blocks simultaneously. The experiments in the next section show that in practice the space used to keep track of irrelevant blocks is less than 10% of the total.

Note that also in [23] the authors faced the problem of limiting the memory used to keep track of the active blocks. They suggested the heuristic of keeping track of active blocks only after the  $h$ -th iteration ( $h = 20$  for their dataset).

#### 4.1. Experimental Results

We have implemented the Gap algorithm in C and tested it on the collections shown in Table 1 which have documents of different size, LCP, and alphabet

size. We represented LCP values with the minimum possible number of bytes for each collection: 1 byte for **Illumina**, 2 bytes for **Pacbio** and **Proteins**, and 4 bytes for **Wiki-it**. We always used 1 byte for each BWT value and  $n$  bytes to represent a pair of  $Z^{(h)}$  arrays using 4 bits for each entry so that the tested implementation can merge simultaneously up to 16 BWTs.

Referring to Table 2, we split each collection into  $k$  subcollections of size less than 2GB and we computed the multi-string SA of each subcollection using **gSACA-K** [35]. From the SA we computed the multi-string BWT and LCP arrays using the  $\Phi$  algorithm [36] (implemented in **gSACA-K**). This computation used 13 bytes per input symbol. Then, we merged the subcollections BWTs and LCPs using **Gap** with different values of the parameter  $\tau$  which determines the size of the smallest irrelevant block we keep track of. Since skipping a block takes time proportional to  $\sigma + k$ , regardless of  $\tau$  **Gap** never keeps track of blocks smaller than that threshold; therefore for **Wiki-it** we performed a single experiment where the smallest irrelevant block size was  $\sigma + k = 215$ .

From the results in Table 2 we see that **Gap**'s running time is indeed roughly proportional to the average LCP. For example, **Pacbio** and **Illumina** collections both consist of DNA reads but, despite **Pacbio** reads being longer and having a larger maximum LCP, **Gap** is twice as fast on them because of the smaller average LCP. Similarly, **Gap** is faster on **Wiki-it** than on **Proteins** despite the latter collection having a smaller alphabet and shorter documents.

As expected, the parameter  $\tau$  offers a time-space tradeoff for the **Gap** algorithm. In the space reported in Table 2, the fractional part is the peak space usage for irrelevant blocks, while the integral value is the space used by the arrays  $\text{bwt}_i$ ,  $B$  and  $Z^{(h)}$ . For example, for **Wiki-it** we use  $n$  bytes for the BWTs,  $4n$  bytes for the LCP values (the  $B$  array),  $n$  bytes for  $Z^{(h)}$ , and the remaining  $0.55n$  bytes are mainly used for keeping track of irrelevant blocks. This is a relatively high value, about 9% of the total space, since in our current implementation the storage of a block grows linearly with the alphabet size. For DNA sequences and  $\tau = 200$  the cost of storing blocks is less than 3% of the total without a significant slowdown in the running time.

For completeness, we tested the **H&M** implementation from [23] on the **Pacbio** collection. The running time was  $14.57 \mu\text{secs}$  per symbol and the space usage 2.28 bytes per symbol. These values are only partially significant for several reasons: (i) **H&M** computes the BWT from scratch, hence doing also the work of **gSACA-K**, (ii) **H&M** doesn't compute the LCP array, hence the lower space usage, (iii) the algorithm is implemented in Cython which makes it easier to use in a Python environment but is not as fast and space efficient as C.

#### 4.2. Merging only BWTs

If we are not interested in LCP values but we only need to merge BWTs, we can still use **Gap** instead of **H&M** to do the computation in  $\mathcal{O}(n \text{aveLcp})$  time. In that case however, the use of the integer array  $B$  recording LCP values is wasteful. We can save space replacing it with an array  $B_2[1, n_0 + n_1 + 1]$  containing two bits per entry representing four possible states called  $\{0, 1, 2, 3\}$ . The rationale for this is that, if we are not interested in LCP values, the entries

---

```

4: if  $B_2[k] \neq 0$  and  $B_2[k] \neq 2$  then
5:    $\text{id} \leftarrow k$                                 ▷ A new block of  $Z^{(h-1)}$  is starting
6: end if
7: if  $B_2[k] = 1$  then
8:    $B_2 \leftarrow 3$                                 ▷ Mark the block as old
9: end if
    $\vdots$ 
13: if  $B_2[j] = 0$  then                            ▷ Check if already marked
14:    $B_2[j] \leftarrow 2$                             ▷ A new block of  $Z^{(h)}$  will start here
15: end if

```

---

Figure 5: Modification of the H&M algorithm to use a two-bit array  $B_2$  instead of the integer array  $B$ . The code shows the case for  $h$  even; if  $h$  is odd, the value  $2$  is replaced by  $1$  and viceversa.

of  $B$  are only used in Line 4 of Figure 3 where it is tested whether they are different from 0 or  $h$ .

During iteration  $h$ , the values in  $B_2$  are used instead of the ones in  $B$  as follows: An entry  $B_2[i] = 0$  corresponds to  $B[i] = 0$ , an entry  $B_2[i] = 3$  corresponds to  $0 < B[i] < h - 1$ . If  $h$  is even, an entry  $B_2[i] = 2$  corresponds to  $B[i] = h$  and an entry  $B_2[i] = 1$  corresponds to  $B[i] = h - 1$ ; while if  $h$  is odd the correspondence is  $2 \rightarrow h - 1$ ,  $1 \rightarrow h$ . The array  $B_2$  is initialized as  $\mathfrak{z}(0)^{n_0+n_1-1}(\mathfrak{z})$ , and it is updated appropriately in lines 13–14. The reason for this apparently involved scheme is that during iteration  $h$ , an entry in  $B_2$  can be modified either before or after we read it at Line 4. The resulting code is shown in Figure 5. Using the array  $B_2$  we can still define (and skip) monochrome blocks and therefore achieve the  $\mathcal{O}(n \text{ aveLcp})$  complexity.

Notice that, by Corollary 1, the value in  $B_2[i]$  changes from 0 to 2 or 1 during iteration  $h = \text{lcp}_{01}[i] + 1$ . Hence, if every time we do such change we write to an external file the pair  $\langle i, h - 1 \rangle$ , when the merging is complete the file contains all the information required to compute the LCP array  $\text{lcp}_{01}$  even if we do not know  $\text{lcp}_0$  and  $\text{lcp}_1$ . This idea has been introduced and investigated in [31] where the Gap algorithm has been modified to work in external memory with a limited amount of RAM. In [31], BWTs of subcollections are initially computed in RAM, the dimensions of the subcollections chosen according to the available RAM. Then, BWTs are merged, in an appropriate number of rounds. In the very final merging round, the LCP values are written to disk as outlined above (i.e., as pairs  $\langle i, h - 1 \rangle$ ) and, at the end of the BWT merging, the LCP array is reconstructed from the above pairs with an external memory merge sort algorithm.

## 5. Merging compressed tries

Tries [37] are a fundamental data structure for representing a collection of  $k$  distinct strings. A trie consists of a rooted tree in which each edge is labeled

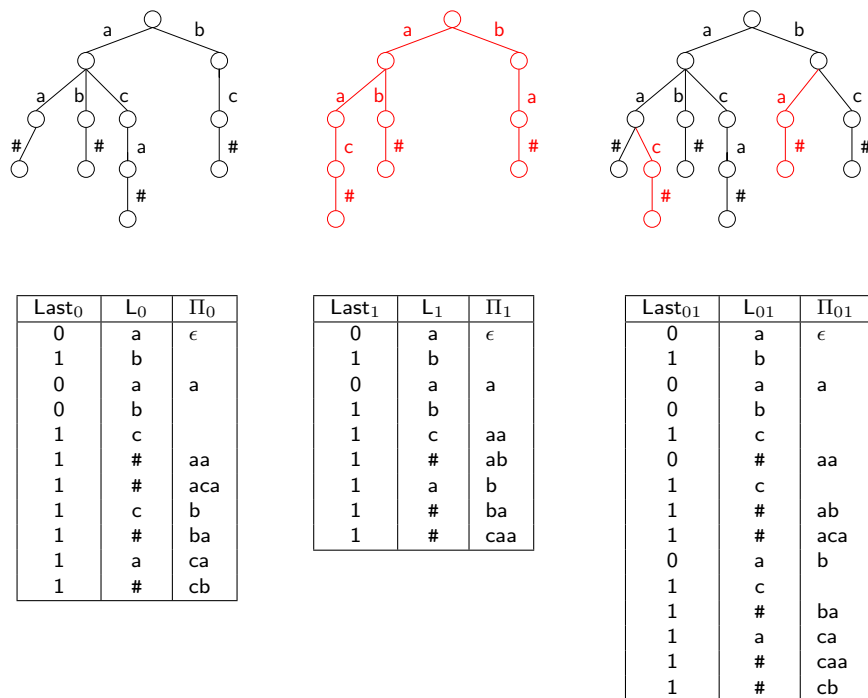


Figure 6: The trie  $T_0$  containing the strings  $aa\#, ab\#, aca\#, bc\#$  (left), the trie  $T_1$  containing  $aac\#, ab\#, ba\#$  (center) and the trie  $T_{01}$  containing the union of the two set of strings (right). Below each trie we show the corresponding XBWT representation.

with a symbol in the input alphabet, and each string is represented by a path from the root to one of the leaves. To simplify the algorithms, and ensure that no string is the prefix of another one, it is customary to add a special symbol  $\# \notin \Sigma$  at the end of each string.<sup>2</sup> Tries for different sets of strings are shown in Figure 6. For any trie node  $u$  we write  $\text{hgt}(u)$  to denote its height, that is the length of the path from the root to  $u$ . We define the height of the trie  $T$  as the maximum node height  $\text{hgt}(T) = \max_u \text{hgt}(u)$ , and the average height  $\text{avehgt}(T) = (\sum_u \text{hgt}(u))/|T|$ , where  $|T|$  denotes the number of trie nodes.

The eXtended Burrows-Wheeler Transform [9, 38, 39] is a generalization of the BWT designed to compactly represent any labeled tree  $T$ . To define  $\text{xbwt}(T)$ , to each *internal* node  $w$  we associate the string  $\lambda_w$  obtained by concatenating the symbols in the edges in the upward path from  $w$  to the root of  $T$ . If  $T$  has  $n$  internal nodes we have  $n$  strings overall; let  $\Pi[1, n]$  denote the array

<sup>2</sup>In this and in the following section we purposely use a special symbol  $\#$  different from  $\$$ . The reason is that  $\$$  is commonly used for sorting purposes, while  $\#$  simply represents a symbol different from the ones in  $\Sigma$ .

containing such strings sorted lexicographically. Note that  $\Pi[1]$  is always the empty string corresponding to the root of  $T$ . For  $i = 1, \dots, n$  let  $L(i)$  denote the set of symbols labeling the edges exiting from the node corresponding to  $\Pi[i]$ . We define the array  $\mathbf{L}$  as the concatenation of the arrays  $L(1), \dots, L(n)$ . If  $T$  has  $m$  edges (and therefore  $m + 1$  nodes), it is  $|\mathbf{L}| = m$  and  $\mathbf{L}$  contains  $n - 1$  symbols from  $\Sigma$  and  $m + 1 - n$  occurrences of  $\#$ . To keep an explicit representation of the intervals  $L(1), \dots, L(n)$  within  $\mathbf{L}$ , we define a binary array  $\mathbf{Last}[1, m]$  such that  $\mathbf{Last}[i] = \mathbf{1}$  iff  $\mathbf{L}[i]$  is the last symbol of some interval  $L(j)$ . See Figure 6 for a complete example.

In [8] it is shown that the two arrays  $\mathbf{xbwt}(T) = \langle \mathbf{Last}, \mathbf{L} \rangle$  are sufficient to represent  $T$ , and that if they are enriched with data structures supporting constant time rank and select operations,  $\mathbf{xbwt}(T)$  can be used for efficient upward and downward navigation and for substring search in  $T$ . The fundamental property for efficient navigation and search is that there is an one-to-one correspondence between the symbols in  $\mathbf{L}$  different from  $\#$  and the strings in  $\Pi$  different from the empty string. The correspondence is order preserving in the sense that the  $i$ -th occurrence of symbol  $c$  corresponds to the  $i$ -th string in  $\Pi$  starting with  $c$ . For example, in Figure 6 (right) the third **a** in  $\mathbf{Last}_{01}$  corresponds to the third string in  $\Pi_{01}$  starting with **a**, namely **ab**. Note that **ab** is the string associated to the node reached by following the edge associated to the third **a** in  $\mathbf{L}_{01}$ .

In this section, we consider the problem of merging two distinct XBWTs. More formally, let  $T_0$  (resp.  $T_1$ ) denote the trie containing the set of strings  $\mathbf{t}_1, \dots, \mathbf{t}_k$  (resp.  $\mathbf{s}_1, \dots, \mathbf{s}_h$ ), and let  $T_{01}$  denote the trie containing the strings in the union  $\mathbf{t}_1, \dots, \mathbf{t}_k, \mathbf{s}_1, \dots, \mathbf{s}_h$  (see Figure 6). Note that  $T_{01}$  might contain less than  $h + k$  strings: if the same string appears in both  $T_0$  and  $T_1$  it will be represented in  $T_{01}$  only once. Given  $\mathbf{xbwt}(T_0) = \langle \mathbf{Last}_0, \mathbf{L}_0 \rangle$  and  $\mathbf{xbwt}(T_1) = \langle \mathbf{Last}_1, \mathbf{L}_1 \rangle$  we want to compute the XBWT representation of the trie  $T_{01}$ .

We observe that if we had at our disposal the sorted string arrays  $\Pi_0$  and  $\Pi_1$ , then the construction of  $\mathbf{xbwt}(T_{01})$  could be done as follows: First, we merge lexicographically the strings in  $\Pi_0$  and  $\Pi_1$ , then we scan the resulting sorted array of strings. During the scan

- if we find a string appearing only once then it corresponds to an internal node belonging to either  $T_0$  or  $T_1$ ; the labels on the outgoing edges can be simply copied from the appropriate range of  $\mathbf{L}_0$  or  $\mathbf{L}_1$ .
- if we find two consecutive equal strings they correspond respectively to an internal node in  $T_0$  and to one in  $T_1$ . The corresponding node in  $T_{01}$  has a set of outgoing edges equal to the union of the edges of those nodes in  $T_0$  and  $T_1$ : thus, the labels in the outgoing edges are the union of the symbols in the appropriate ranges of  $\mathbf{L}_0$  and  $\mathbf{L}_1$ .

Although the arrays  $\Pi_0$  and  $\Pi_1$  are not available, by properly modifying the H&M algorithm we can compute how their elements would be interleaved by the merge operation. Let  $m_0 = |\mathbf{L}_0| = |\mathbf{Last}_0|$ ,  $n_0 = |\Pi_0|$ , and similarly  $m_1 = |\mathbf{L}_1| = |\mathbf{Last}_1|$ ,  $n_1 = |\Pi_1|$ . Figure 7 shows the code for the generic  $h$ -th iteration of the H&M algorithm adapted for the XBWT. Iteration  $h$  computes



---

```

1: Initialize array  $F[1, \sigma]$ 
2:  $k_0 \leftarrow 1; k_1 \leftarrow 1$  ▷ Init counters for  $L_0$  and  $L_1$ 
3:  $Z^{(h)} \leftarrow \mathbf{01}$  ▷ First two entries correspond to  $\Pi_0[1] = \Pi_1[1] = \epsilon$ 
4: for  $k \leftarrow 1$  to  $n_0 + n_1$  do
5:    $b \leftarrow Z^{(h-1)}[k]$  ▷ Read bit  $b$  from  $Z^{(h-1)}$ 
6:   repeat
7:      $c \leftarrow L_b[k_b]$  ▷ Get symbol from  $L_0$  or  $L_1$  according to  $b$ 
8:     if  $c \neq \#$  then ▷  $\#$  is ignored: it is not in  $\Pi_0$  or  $\Pi_1$ 
9:        $j \leftarrow F[c]++$  ▷ Get destination for  $b$  according to symbol  $c$ 
10:       $Z^{(h)}[j] \leftarrow b$  ▷ Copy bit  $b$  to  $Z^{(h)}$ 
11:     end if
12:      $\ell \leftarrow \text{Last}_b[k_b++]$  ▷ Check if  $c$  labels last outgoing edge
13:   until  $\ell \neq 1$ 
14: end for

```

---

Figure 7: Main loop of algorithm **H&M** modified to merge XBWTs. Array  $F$  is initialized so that  $F[c]$  contains the number of occurrences of symbols smaller than  $c$  in  $L_0$  and  $L_1$  plus three, to account for  $\Pi_0[1] = \Pi_1[1] = \epsilon$  which are smaller than any other string.

a binary vector  $Z^{(h)}$  containing  $n_0 = |t_0|$   $\mathbf{0}$ 's and  $n_1 = |t_1|$   $\mathbf{1}$ 's and such that the following property holds (compare with Property 1)

**Property 2.** *At the end of iteration  $h$ , for  $i = 2, \dots, n_0$  and  $j = 2, \dots, n_1$  the  $i$ -th  $\mathbf{0}$  precedes the  $j$ -th  $\mathbf{1}$  in  $Z^{(h)}$  if and only if*

$$\Pi_0[i][1, h] \preceq \Pi_1[j][1, h]. \quad (8)$$

□

In (8)  $\Pi_0[i][1, h]$  denotes the length- $h$  prefix of  $\Pi_0[i]$ . If  $\Pi_0[i]$  has length smaller than  $h$  then  $\Pi_0[i][1, h] = \Pi_0[i]$  (and similarly for  $\Pi_1[j]$ ). Note that Property 2 does not mention the first  $\mathbf{0}$  and the first  $\mathbf{1}$  in  $Z^{(h)}$ : By construction it is  $\Pi_0[1] = \Pi_1[1] = \epsilon$  so we know their lexicographic rank is the smallest possible. Note also that because of Step 3 in Figure 7, the first  $\mathbf{0}$  and the first  $\mathbf{1}$  in  $Z^{(h)}$  are always the first two elements of  $Z^{(h)}$ .

Apart from the first two entries, during iteration  $h$  the array  $Z^{(h)}$  is logically partitioned into  $\sigma$  subarrays, one for each alphabet symbol different from  $\#$ . If  $Occ(c)$  denotes the number of occurrences in  $L_0$  and  $L_1$  of the symbols smaller than  $c$ , then the subarray corresponding to  $c$  starts at position  $Occ(c) + 3$ . Hence, if  $c < c'$  the subarray corresponding to  $c$  precedes the one corresponding to  $c'$ . Because of how the array  $F$  is initialized and updated, we see that every time we read a symbol  $c$  from  $L_0$  and  $L_1$  we write a value in the portion of  $Z^{(h)}$  corresponding to  $c$ , and that each portion is filled sequentially. Armed with these observations, we are ready to establish the correctness of the algorithm in Figure. 7.

**Lemma 7.** *Let  $Z^{(0)} = \mathbf{010}^{n_0-1}\mathbf{1}^{n_1-1}$ , and let  $Z^{(h)}$  be obtained from  $Z^{(h-1)}$  by the algorithm in Figure 7. Then, for  $h = 0, 1, 2, \dots$ , the array  $Z^{(h)}$  satisfies Property 2.*

**Proof:** We prove the result by induction. For  $h = 0$ ,  $\Pi_0[i][1, 0] = \Pi_1[j][1, 0] = \epsilon$  so (8) is always true and  $Z^{(0)}$  satisfies Property 2.

Suppose now  $h > 0$ . To prove the “if” part, let  $3 \leq v < w \leq n_0 + n_1$  denote two indexes such that  $Z^{(h)}[v]$  is the  $i$ -th  $\mathbf{0}$  and  $Z^{(h)}[w]$  is the  $j$ -th  $\mathbf{1}$  in  $Z^{(h)}$  for some  $2 \leq i \leq n_0$  and  $2 \leq j \leq n_1$  (it is  $v \geq 3$  since  $i \geq 2$  and  $Z^{(h)}[1, 2] = \mathbf{01}$ ). We need to show that (8) holds.

Assume first  $\Pi_0[i][1] \neq \Pi_1[j][1]$ . The hypothesis  $v < w$  implies  $\Pi_0[i][1] < \Pi_1[j][1]$  hence (3) certainly holds. Assume now  $\Pi_0[i][1] = \Pi_1[j][1] = c$ . Let  $v'$ ,  $w'$  denote respectively the values of the main loop variable  $k$  in the procedure of Figure 7 when the entries  $Z^{(h)}[v]$  and  $Z^{(h)}[w]$  are written (hence, during the scanning of  $Z^{(h-1)}$ ). The hypotheses  $v < w$  and  $\Pi_0[i][1] = \Pi_1[j][1]$  imply  $v' < w'$ . By construction  $Z^{(h-1)}[v'] = \mathbf{0}$  and  $Z^{(h-1)}[w'] = \mathbf{1}$ . Say  $v'$  is the  $i'$ -th  $\mathbf{0}$  in  $Z^{(h-1)}$  and  $w'$  is the  $j'$ -th  $\mathbf{1}$  in  $Z^{(h-1)}$ . By the inductive hypothesis on  $Z^{(h-1)}$  we have

$$\Pi_0[i'][1, h-1] \preceq \Pi_1[j'][1, h-1] \quad (9)$$

(we could have  $v' = 1$  that would imply  $i' = 1$ ; in that case we cannot apply the inductive hypothesis, but (9) still holds). By the properties of the XBWT we have

$$\Pi_0[i][1, h] = c \Pi_0[i'][1, h-1] \quad \text{and} \quad \Pi_1[j][1, h] = c \Pi_1[j'][1, h-1]$$

which combined with (9) gives us (8).

For the “only if” part assume (8) holds for some  $i \geq 2$  and  $j \geq 2$ . We need to prove that in  $Z^{(h)}$  the  $i$ -th  $\mathbf{0}$  precedes the  $j$ -th  $\mathbf{1}$ . If  $\Pi_0[i][1] \neq \Pi_1[j][1]$  the proof is immediate. If  $c = \Pi_0[i][1] = \Pi_1[j][1]$  then

$$\Pi_0[i][2, h] \preceq \Pi_1[j][2, h].$$

Let  $i'$  and  $j'$  be such that  $\Pi_0[i'][1, h-1] = \Pi_0[i][2, h]$  and  $\Pi_1[j'][1, h-1] = \Pi_1[j][2, h]$ . By induction, in  $Z^{(h-1)}$  the  $i'$ -th  $\mathbf{0}$  precedes the  $j'$ -th  $\mathbf{1}$  (again we could have  $i' = 1$  and in that case we cannot apply the inductive hypothesis, but the claim still holds).

During iteration  $h$ , the  $i$ -th  $\mathbf{0}$  in  $Z^{(h)}$  is written to position  $v$  when processing the  $i'$ -th  $\mathbf{0}$  of  $Z^{(h-1)}$ , and the  $j$ -th  $\mathbf{1}$  in  $Z^{(h)}$  is written to position  $w$  when processing the  $j'$ -th  $\mathbf{1}$  of  $Z^{(h-1)}$ . Since in  $Z^{(h-1)}$  the  $i'$ -th  $\mathbf{0}$  precedes the  $j'$ -th  $\mathbf{1}$  and since  $v$  and  $w$  both belongs to the subarray of  $Z^{(h)}$  corresponding to the symbol  $c$ , their relative order does not change and the  $i$ -th  $\mathbf{0}$  precedes the  $j$ -th  $\mathbf{1}$  as claimed.  $\square$

As in the original H&M algorithm we stop the merge phase after the first iteration  $h$  such that  $Z^{(h)} = Z^{(h-1)}$ . Since in subsequent iterations we would have  $Z^{(g)} = Z^{(h)}$  for any  $g > h$ , we get that by Property 2,  $Z^{(h)}$  gives the correct lexicographic merge of  $\Pi_0$  and  $\Pi_1$ . Note however that the lexicographic order

is not sufficient to establish whether two consecutive nodes, say  $\Pi_0[i]$  and  $\Pi_1[j]$  have the same upward path and therefore should be merged in a single node of  $T_{01}$ . To this end, we consider the integer array  $B$  used in Section 2.1 to mark the starting point of each block. We have shown in Corollary 1 that at the end of the original H&M algorithm  $B$  contains the LCP values plus one. Indeed, at iteration  $h$  the algorithm sets  $B[k] = h$  since it “discovers” that the suffixes in  $\text{sa}_{01}[k-1]$  and  $\text{sa}_{01}[k]$  differ in the  $h$ -th symbol (hence  $\text{lcp}_{01}[k] = h-1$ ). If we maintain the array  $B$  in the XBWT merging algorithm, we get that at the end of the computation if the strings associated to  $\Pi_0[i]$  and  $\Pi_1[j]$  are identical then the entry in  $B$  corresponding to  $\Pi_1[j]$  would be zero, since the two strings do not differ in any position. Hence, at the end of the modified H&M algorithm the array  $Z^{(h)}$  provides the lexicographic order of the nodes, and the array  $B$  the position of the nodes of  $T_0$  and  $T_1$  with the same upward path. We conclude that with a single scan of  $Z^{(h)}$  and  $B$  we can merge all paths and compute  $\text{xbwt}(T_0)$ . Finally, we observe that instead of  $B$  we can use a two-bit array  $B_2$  as in Sect. 4.2, since we are only interested in determining whether a certain entry is zero, and not in its exact value.

**Lemma 8.** *The modified H&M algorithm computes  $\text{xbwt}(T_{01})$  given  $\text{xbwt}(T_0)$  and  $\text{xbwt}(T_1)$  in  $\mathcal{O}(|T_{01}|\text{hgt}(T_{01}))$  time and  $4n + \mathcal{O}(\log n)$  bits of working space, where  $n = n_0 + n_1$ .*

**Proof:** Each iteration of the merging algorithm takes  $\mathcal{O}(m_0 + m_1)$  time since it consists of a scan of the arrays  $Z^{(h-1)}$ ,  $\mathbf{L}_0$ ,  $\mathbf{L}_1$ ,  $\mathbf{Last}_0$  and  $\mathbf{Last}_1$ . After at most  $\text{hgt}(T_{01})$  iterations the strings in  $\Pi_0$  and  $\Pi_1$  are lexicographically sorted and  $Z^{(h)}$  no longer changes. The final scan of  $Z^{(h)}$  and  $B_2$  to compute  $\text{xbwt}(T_0)$  takes  $\mathcal{O}(m_0 + m_1)$  time. Since  $|T_{01}| \geq \max(m_0, m_1)$  the overall cost is  $\mathcal{O}(|T_{01}|\text{hgt}(T_{01}))$  time. The working space of the algorithm, consists of  $B_2$  and of two instances of the  $Z^{(h)}$  array (for the current and the previous iteration), in addition to  $\mathcal{O}(\sigma)$  counters (recall that  $\sigma$  is assumed to be constant).  $\square$

As for BWT/LCP merging, we now show how to reduce the running time by skipping the portions of  $Z^{(h)}$  that no longer change from one iteration to the next. Note that we cannot use monochrome blocks to early terminate XBWT merging. Indeed, from the previous discussion we know that if two strings  $\Pi_0[i]$  and  $\Pi_1[j]$  are equal, they will form a non-monochrome block that will never be split.

For this reason we introduce an array  $C[1, n_0 + n_1]$  that, at the beginning of iteration  $h$ , keeps track of all the strings in  $\Pi_0$  and  $\Pi_1$  that have length less than  $h$ . More precisely, for  $i = 1, \dots, n_0$  (resp.  $j = 1, \dots, n_1$ ) if the  $i$ -th  $\mathbf{0}$  (resp. the  $j$ -th  $\mathbf{1}$ ) is in position  $k$  of  $Z^{(h)}$ , then  $C[k] = \ell > 0$  iff the length  $|\Pi_0[i]|$  (resp.  $|\Pi_1[j]|$ ) is equal to  $\ell - 1$  with  $\ell - 1 < h$ . As a consequence, if  $C[k] = 0$  then the string corresponding to  $C[k]$  has length  $h$  or more. Note that by Property 2 at the beginning of iteration  $h$  the algorithm has already determined the lexicographic rank of all the strings in  $\Pi_0$  and  $\Pi_1$  of length smaller than  $h$ . Hence, the entry in  $Z^{(h)}[k]$  will not change in successive iterations and will remain associated to the same string from  $\Pi_0$  or  $\Pi_1$ .

The array  $C$  is initialized as  $110^{n_0+n_1-2}$  since at the beginning of iteration 1 it is  $Z^{(h)} = \mathbf{010}^{n_0-1}\mathbf{1}^{n_1-1}$  and indeed the only strings of length 0 are  $\Pi_0[0] = \Pi_1[0] = \epsilon$ . During iteration  $h$ , we update  $C$  adding, immediately after Line 10 in Figure 7, the line

**if**  $C[k] = h$  **then**  $C[j] \leftarrow h + 1$

The rationale is that if, during iteration  $h - 1$  we found out that the string  $\alpha$  corresponding to  $Z^{(h-1)}[k]$  has length  $h - 1$  (so we set  $C[k] = h$ ), then the string corresponding to  $Z^{(h)}[j]$  is  $c\alpha$  and has therefore length  $h$ .

By the above discussion we see that if at iteration  $h$  we write  $h + 1$  to position  $C[j]$ , then at iteration  $h + 1$  we can possibly use  $C[j]$  to write  $h + 2$  in some other position in  $C$ , but starting from iteration  $h + 2$  it is no longer necessary to process neither  $C[j]$  nor  $Z^{(h+2)}[j]$  since they will not affect neither  $C$  nor  $Z^{(h+3)}$ . In other words, during iteration  $h$  we can skip all ranges  $Z^{(h)}[\ell, m]$  such that  $C[\ell, m]$  contains only positive values smaller than  $h$ . These ranges grown larger and larger as the algorithm proceeds and are handled in the same way as the irrelevant blocks in **Gap**. Finally, we observe that, using the same techniques as in Section 4.2, we can replace the integer array  $C$  with an array  $C_2$  containing only two bits per entry.

**Theorem 9.** *The modified **Gap** algorithm computes  $\text{xbwt}(T_{01})$  given  $\text{xbwt}(T_0)$  and  $\text{xbwt}(T_1)$  in  $\mathcal{O}(|T_{01}|\text{avehgt}(T_{01}))$  time. The working space is  $6n + \mathcal{O}(\log n)$  bits, where  $n = n_0 + n_1$ , plus the space required for handling irrelevant blocks.*

**Proof:** The analysis is similar to the one in Theorem 6. Here the algorithm executes  $\text{hgt}(T_{01})$  iterations; however, because of irrelevant blocks, iterations have decreasing costs. To bound the overall running time, observe that the cost of each iteration is dominated by the cost of processing the entries in  $L_0$  and  $L_1$ . The generic entry  $L_0[i]$  corresponds to a trie node  $u_i$  with upward path of length  $\text{hgt}(u_i)$ . Entry  $L_0[i]$  is processed when the **Gap** algorithm reaches the entry in  $Z^{(h)}$  corresponding to the string  $\Pi_1[i']$  associated to  $u_i$ 's parent. We know that  $Z^{(h)}$ 's entry corresponding to  $\Pi_1[i']$  becomes irrelevant after iteration  $|\Pi_1[i']| + 1 = \text{hgt}(u_i)$ . Hence, the overall cost of processing  $u_i$  is  $\mathcal{O}(\text{hgt}(u_i))$ . Summing over all entries in  $L_0$  and  $L_1$  the total cost is  $\mathcal{O}(|T_0|\text{avehgt}(T_0) + |T_1|\text{avehgt}(T_1))$ . The thesis follows observing that  $|T_{01}|\text{avehgt}(T_{01}) \geq \max(|T_0|\text{avehgt}(T_0), |T_1|\text{avehgt}(T_1))$ .  $\square$

## 6. Merging indices for circular patterns

Another well known variant of the BWT is the *multistring circular BWT* which is defined by sorting the cyclic rotations of the input strings instead of their suffixes. However, to make the transformation reversible, the cyclic rotations have to be sorted according to an order relation, different from the lexicographic order, that we now quickly review.

For any string  $t$ , we define the *infinite form*  $t^\infty$  of  $t$  as the infinite length string obtained concatenating  $t$  to itself infinitely many times. Given two strings

$t$  and  $s$  we write  $t \preceq^\infty s$  to denote that  $t^\infty \preceq s^\infty$ . For example, for  $t = abaa$  and  $s = aba$ , it is  $t^\infty = abaaabaa \dots$  and  $s^\infty = abaabaaba \dots$  so  $t \preceq^\infty s$ . Notice that  $t^\infty = s^\infty$  does not necessarily imply that  $t = s$ . For example, for  $t = ababab$  and  $s = abab$  it is  $t^\infty = s^\infty$ . The following lemma, which is a consequence of Fine and Wilf Theorem [40] and a restatement of Proposition 5 in [4], provides an upper bound to the number of comparisons required to establish whether  $t^\infty = s^\infty$ .

**Lemma 10.** *If  $t^\infty \neq s^\infty$  then there exists an index  $i \leq |t| + |s| - \gcd(|t|, |s|)$  such that  $t^\infty[i] \neq s^\infty[i]$ .  $\square$*

A string is *primitive* if all its cyclic rotations are distinct. The following Lemma is another well known consequence of the Fine and Wilf Theorem.

**Lemma 11.** *If  $t$  and  $s$  are primitive,  $t^\infty = s^\infty$  implies  $t = s$ .  $\square$*

Let  $t_0[1, n_0], t_1[1, n_1]$  be two primitive strings and  $t_{01}[1, n]$  their concatenation of length  $n = n_0 + n_1$ . For  $i = 1, \dots, n$ , let  $\text{rot}_{01}(i)$  define the rotation of substrings  $t_0$  and  $t_1$  within  $t_{01}$  as follows:

$$\text{rot}_{01}(i) = \begin{cases} t_0[i, n_0]t_0[1, i-1] & \text{if } 0 < i \leq n_0 \\ t_1[i-n_0, n_1]t_1[1, i-n_0-1] & \text{if } n_0 < i \leq n_0 + n_1. \end{cases}$$

For example, if  $t_0 = abc$  and  $t_1 = abbb$ , it is  $\text{rot}_{01}(2) = bca$  and  $\text{rot}_{01}(7) = babb$ . The above definition of rotations of substrings can be obviously generalized to a collection of  $k$  strings.

We are interested in circular BWTs as the core of compressed indices and therefore we assume in the following that  $t_0$  and  $t_1$  are primitive and that  $t_0$  is not a rotation of  $t_1$ . Notice, though, that the algorithms we present are correct and work with the same time and space complexities also in the general case. We define the *circular Suffix Array* of  $t_0$  and  $t_1$ ,  $\text{csa}_{01}$  as the permutation of  $[1, n]$  such that:

$$\text{rot}_{01}(\text{csa}_{01}[i]) \preceq^\infty \text{rot}_{01}(\text{csa}_{01}[i+1]). \quad (10)$$

Note that because of our assumptions and Lemma 11, the inequality in (10) is always strict. Finally, the multistring *circular Burrows-Wheeler Transform* (*cBWT*) is defined as

$$\text{cbwt}_{01}[i] = \begin{cases} t_0[n_0] & \text{if } \text{csa}_{01}[i] = 1 \\ t_0[\text{csa}_{01}[i] - 1] & \text{if } 1 < \text{csa}_{01}[i] \leq n_0 \\ t_1[n_1] & \text{if } \text{csa}_{01}[i] = n_0 + 1 \\ t_1[\text{csa}_{01}[i] - n_0 - 1] & \text{if } \text{csa}_{01}[i] > n_0 + 1. \end{cases}$$

The above definition given for  $t_0$  and  $t_1$  can be generalized to any number of strings. The  $\preceq^\infty$  order and the above multistring circular BWT has been introduced in [4]. In [7] the authors use a data structure equivalent to a circular BWT to design a *compressed permuterm index* for prefix/suffix queries. The

crucial observation is that if we add a unique symbol  $\#$  at the end of each string, the same symbol for every string, then searching  $\beta\#\alpha$  in a circular BWT returns all the strings simultaneously prefixed by  $\alpha$  and suffixed by  $\beta$ . In [30] Hon et al. use the circular BWT to design a succinct index for circular patterns. Note that Hon et al. in addition to  $\text{cbwt}_{01}$  use an additional data structure  $\text{length}_{01}$  such that  $\text{length}_{01}(i)$  provides the length of the string  $t_j$  to which the symbol  $\text{cbwt}_{01}[i]$  belongs. Finally, a lightweight algorithm for the construction of the circular BWT has been described in [41]: for a string of length  $n$  the proposed algorithm takes  $\mathcal{O}(n)$  time and uses  $\mathcal{O}(n \log \sigma)$  bits of space.

To simplify our analysis, we preliminary extend the concept of longest common prefix to the  $\preceq^\infty$  order. For any pair of strings  $t, s$ , we define

$$\text{clcp}(t, s) = \begin{cases} \text{LCP}(t^\infty, s^\infty) & \text{if } t^\infty \neq s^\infty \\ |t| + |s| - \text{gcd}(|t|, |s|) & \text{otherwise.} \end{cases} \quad (11)$$

Because of Lemma 10,  $\text{clcp}(t, s)$  generalizes the standard LCP in that it provides the number of comparisons that are necessary in order to establish the  $\preceq^\infty$  ordering between  $t, s$ . It is then natural to define for  $i = 2, \dots, n$

$$\text{clcp}_{01}[i] = \text{clcp}(\text{rot}_{01}(\text{csa}_{01}[i-1]), \text{rot}_{01}(\text{csa}_{01}[i])) \quad (12)$$

and the values

$$\text{maxLcp} = \max_i \text{clcp}_{01}[i], \quad \text{aveLcp} = \left( \sum_i \text{clcp}_{01}[i] \right) / n. \quad (13)$$

that generalize the standard notions of maximum LCP and average LCP.

Let  $\text{cbwt}_0$  (resp.  $\text{cbwt}_1$ ) denote the circular BWT for the collection of strings  $t_1, \dots, t_k$  (resp.  $s_1, \dots, s_h$ ). In this section we consider the problem of computing the circular BWT  $\text{cbwt}_{01}$  for the union collection  $t_1, \dots, t_k, s_1, \dots, s_h$ . As we previously observed, we assume that all strings are primitive and that within each input collection no string is the rotation of another. However, we cannot rule out the possibility that some  $t_i$  is the rotation of some  $s_j$ . The merging algorithm should therefore recognize this occurrence and eliminate from the union one of the two strings, say  $s_j$ . In practice, this means that all symbols of  $\text{cbwt}_1$  coming from  $s_j$  must not be included in  $\text{cbwt}_{01}$ .

To merge  $\text{cbwt}_0$  and  $\text{cbwt}_1$  we need to merge their symbols according to their context. By construction, the context of  $\text{cbwt}_0[i]$  (resp.  $\text{cbwt}_1[j]$ ) is  $\text{rot}_0(\text{csa}_0[i])$  (resp.  $\text{rot}_1(\text{csa}_1[j])$ ), where  $\text{rot}_0(\text{csa}_0[i])$  is a cyclic rotation of the string  $t_k$  to which the symbol  $\text{cbwt}_0[i]$  belongs (and similarly for  $\text{rot}_1(\text{csa}_1[j])$ ). Note, however, that contexts must be sorted according to the  $\prec^\infty$  order; hence  $\text{cbwt}_0[i]$  should precede  $\text{cbwt}_1[j]$  in  $\text{cbwt}_{01}$  iff  $\text{rot}_0(\text{csa}_0[i]) \preceq^\infty \text{rot}_1(\text{csa}_1[j])$ . The good news is that the H&M algorithm, as described in Figure 2, when applied to  $\text{cbwt}_0$  and  $\text{cbwt}_1$  will sort each symbol according to the  $\preceq^\infty$  order of its context. Notice that the  $\preceq^\infty$  order induces a significant difference with respect to the merging of BWTs: indeed, since there are no  $\$$ 's in  $\text{cbwt}_0$  and  $\text{cbwt}_1$ , Line 9 is never executed and the destination of each symbol is always determined by its predecessor in the cyclic rotation. More formally, reasoning as in Lemma 1, it is possible to prove the following property.

**Property 3.** For  $i = 1, \dots, n_0$  and  $j = 1, \dots, n_1$  the  $i$ -th  $\mathbf{0}$  precedes the  $j$ -th  $\mathbf{1}$  in  $Z^{(h)}$  if and only if

$$\text{rot}_0(\text{csa}_0[i])^\infty[1, h] \preceq \text{rot}_1(\text{csa}_1[j])^\infty[1, h]. \quad (14)$$

□

Property 3 states that after iteration  $h$  the infinite strings  $\text{rot}_0(\text{csa}_0[i])^\infty$  and  $\text{rot}_1(\text{csa}_1[j])^\infty$  have been sorted according to their length  $h$  prefix. As for the original H&M algorithm, as soon as  $Z^{(h+1)} = Z^{(h)}$  the  $Z^{(\cdot)}$  array will not change in any successive iteration and the merging is complete. By Lemma 10 it is  $Z^{(h+1)} = Z^{(h)}$  for some  $h \leq \text{maxLcp}$ .

Since we do not simply need to sort the context, but also recognize if some string  $\mathbf{t}_i$  is a rotation of some  $\mathbf{s}_j$ , we make use of the algorithm in Figure 3 which, in addition to  $Z^{(h)}$ , also computes the integer array  $B$  that marks the boundaries of the groups of all rotations whose infinite form have a common prefix of length  $h$ . We can prove a result analogous to Lemma 2 replacing the LCP between suffixes ( $\text{lcp}_{01}$ ) with the LCP between the infinite strings  $\text{rot}_b(\text{csa}_b[i])^\infty$  (that is  $\text{clcp}_{01}$ ). After iteration  $h = \text{maxLcp}$  all distinct rotations have been sorted according to the  $\preceq^\infty$  order; thus an entry  $B[k] = 0$  denotes two rotations  $\text{rot}_0(\text{csa}_0[i])^\infty$  and  $\text{rot}_1(\text{csa}_1[j])^\infty$  which have a common prefix of length  $\text{maxLcp}$ . By Lemma 10 it is  $\text{rot}_0(\text{csa}_0[i])^\infty = \text{rot}_1(\text{csa}_1[j])^\infty$  and by Lemma 11  $\text{rot}_0(\text{csa}_0[i]) = \text{rot}_1(\text{csa}_1[j])$ . The two rotations are therefore identical and the symbol  $\text{cbwt}_1[j]$  should not be included in  $\text{cbwt}_{01}$ .

Summing up, to merge  $\text{cbwt}_0$  and  $\text{cbwt}_1$  we execute the procedure of Figure 3 until both  $Z^{(h)}$  and  $B$  do not change. Then, we compute  $\text{cbwt}_{01}$  by merging  $\text{cbwt}_0$  and  $\text{cbwt}_1$  according to  $Z^{(h)}$ , discarding those symbols corresponding to zero entries in  $B$ . The number of iterations will be at most  $\text{maxLcp}$ . In addition, since we are only interested in zero/nonzero entries, instead of  $B$  we can use a 2-bit array  $B_2$  as in Section 4.2. Reasoning as for Lemma 3, setting  $n = n_0 + n_1$  we get the following result.

**Lemma 12.** *The modified H&M algorithm computes  $\text{cbwt}_{01}$  given  $\text{cbwt}_0$  and  $\text{cbwt}_1$  in  $\mathcal{O}(n \text{maxLcp})$  time and  $4n + \mathcal{O}(\log n)$  bits of working space.* □

As we have done in the previous sections, we now show how to reduce the running time of the merging algorithm by avoiding to re-process the blocks of  $Z^{(h-1)}$  that have become irrelevant for the computation of the new bitarray  $Z^{(h)}$ . Reasoning as in Section 4 we observe that monochrome blocks, i.e. blocks containing entries only from  $\text{cbwt}_0$  or  $\text{cbwt}_1$ , after having been processed once, become irrelevant and can be skipped in successive iterations. Note however, that whenever  $\text{rot}_0(\text{csa}_0[i])^\infty = \text{rot}_1(\text{csa}_1[j])^\infty$  these two entries will always belong to the same block. To handle this case we first assume  $\text{cbwt}_{01}$  is to be used as a compressed index for circular patterns [30] and we later consider the case in which  $\text{cbwt}_{01}$  is to be used for a compressed permuterm index.

### 6.1. Compressed indices of circular patterns

In this setting,  $\text{cbwt}_{01}$  is to be used as a compressed index for circular patterns and therefore we have access to the  $\text{length}_0$  and  $\text{length}_1$  data structures providing the length of each rotation. Under this assumption we modify the `Gap` algorithm described in Section 4 as follows: in addition to skipping monotone blocks, every time there is a size-2 non monochrome block containing, say  $\text{cbwt}_0[i]$  and  $\text{cbwt}_1[j]$ , we mark it as *quasi-irrelevant* and compute  $\ell_{ij} = |\text{length}_0(i)| + |\text{length}_1(j)| - \gcd(|\text{length}_0(i)|, |\text{length}_1(j)|)$ . As soon as this block is split or we reach iteration  $\ell_{ij}$  the block becomes irrelevant and is skipped in successive iterations. As in the original `Gap` algorithm, the computation stops when all blocks have become irrelevant.

For simplicity, in the next theorem we assume that the access to the data structures  $\text{length}_0$  and  $\text{length}_1$  takes constant time. If not, and random access to the individual lengths takes  $\mathcal{O}(\rho)$  time, the overall cost of the algorithm is increased by  $\mathcal{O}((n_0 + n_1)\rho)$  since each length is computed at most once.

**Theorem 13.** *The modified `Gap` algorithm computes  $\langle \text{cbwt}_{01}, \text{length}_{01} \rangle$  given  $\langle \text{cbwt}_0, \text{length}_0 \rangle$  and  $\langle \text{cbwt}_1, \text{length}_1 \rangle$  in  $\mathcal{O}(n \text{ avecLcp})$  time, where  $n = n_0 + n_1$ . The working space is  $2n + \mathcal{O}(\log n)$  bits plus the space required for handling (quasi-)irrelevant blocks.*

**Proof:** If  $\text{rot}_{01}(\text{csa}_{01}[k])$  is different from any other rotation, by definitions (11) and (12) after at most  $\max(\text{clcp}_{01}[k], \text{clcp}_{01}[k+1])$  iterations it will be in a monochrome (possibly singleton) block. If instead  $\text{rot}_{01}(\text{csa}_{01}[k])$  is identical to another rotation, which can only be either  $\text{rot}_{01}(\text{csa}_{01}[k-1])$  or  $\text{rot}_{01}(\text{csa}_{01}[k+1])$ , then after at most

$$\max(\text{clcp}_{01}[k-1], \text{clcp}_{01}[k], \text{clcp}_{01}[k+1], \text{clcp}_{01}[k+2]) \quad (15)$$

iterations it will be in a size-2 non-monochrome block together with its identical rotation. In either case, the block containing  $\text{rot}_{01}(\text{csa}_{01}[k])$  will become irrelevant and it will be no longer processed in successive iterations. Hence, the overall cost of handling  $\text{rot}_{01}(\text{csa}_{01}[k])$  over all iterations is proportional to (15), and the overall cost of handling all rotations is bounded by  $\mathcal{O}(n \text{ avecLcp})$  as claimed. Note that the final bitarray  $Z^{(h)}$  describes also how  $\text{length}_0$  and  $\text{length}_1$  must be interleaved to get  $\text{length}_{01}$ .  $\square$

### 6.2. Compressed permuterm indices

Finally, we consider the case in which  $\text{cbwt}_{01}$  is to be used as the core of a compressed permuterm index [7]. In this case we do not have the  $\text{length}_0$  and  $\text{length}_1$  data structures, but each string in the collection is terminated by a unique  $\#$  symbol. In this case, to recognize whether a size-2 non-monochrome block contains two identical rotations, we make use of the following lemma.

**Lemma 14.** *Let  $\mathbf{t}$  and  $\mathbf{s}$  denote two strings each one containing a single occurrence of the symbol  $\#$ . If for some  $h > 0$  it is  $\mathbf{t}^\infty[1, h] = \mathbf{s}^\infty[1, h]$  and  $\mathbf{t}^\infty[1, h]$  contains two occurrences of  $\#$ , then  $\mathbf{t} = \mathbf{s}$ .*



**Proof:** Let  $\delta$  denote the distance between the two occurrences of  $\#$  in  $t^\infty[1, h]$ . Since  $t$  contains a single  $\#$ , we have  $t = t^\infty[1, \delta] = s^\infty[1, \delta] = s$ .  $\square$

The above lemma suggests to design a **#Gap** algorithm to merge compressed permuterm indices in which the arrays  $Z^{(\cdot)}$  are arrays of pairs so that they keep track also of the number of  $\#$ s in each prefix. In the following  $Z^{(h)}[k] = \langle b, m \rangle$  means that the  $k$ -th rotation belongs to  $\text{csa}_b$ , and among the first  $h$  symbols of the infinite form of that rotation there are exactly  $m$  occurrences of  $\#$ . Formally, for  $h = 0, 1, 2, \dots$  the array  $Z^{(h)}$  satisfies the following property.

**Property 4.** *At the end of iteration  $h$  of #Gap Property 3 holds and if  $Z^{(h)}[k] = \langle b, m \rangle$  is the  $i$ -th  $b$  in  $Z^{(h)}$  then  $\text{rot}_b(\text{csa}_b[i])^\infty[1, h]$  contains exactly  $m$  copies of symbol  $\#$ .*  $\square$

Initially we set  $Z^{(0)} = \langle \mathbf{0}, 0 \rangle^{n_0} \langle \mathbf{1}, 0 \rangle^{n_1}$  which clearly satisfies Property 4. At each iteration **#Gap** reads  $Z^{(h-1)}$  and updates  $Z^{(h)}$  using Lines 7–11 below instead of Lines 7–14 of Figure 3:

```

7:  $\langle b, m \rangle \leftarrow Z^{(h-1)}[k]$ 
8:  $c \leftarrow \text{cbwt}_b[k_b++]$  ▷ Get  $c$  according to  $b$ 
9:  $j \leftarrow F[c]++$  ▷ Get destination for  $b$  according to  $c$ 
10: if  $c = \#$  then  $m \leftarrow m + 1$  ▷ Update number of  $\#$ 
11:  $Z^{(h)}[j] \leftarrow \langle b, m \rangle$ 

```

Reasoning as in the previous sections, one can prove by induction that with this modification the array  $Z^{(h)}$  computed by **#Gap** satisfies Property 4. In the **#Gap** algorithm a block becomes irrelevant when it is monochrome or it is a size-2 non-monochrome block  $Z^{(h)}[k, k+1]$  such that  $Z^{(h)}[k] = \langle \mathbf{0}, 2 \rangle$  and  $Z^{(h)}[k+1] = \langle \mathbf{1}, 2 \rangle$ . By Lemma 14 such block corresponds to two identical rotations  $\text{rot}_0(\text{csa}_0[i]) = \text{rot}_1(\text{csa}_1[j])$  and after being processed a final time it can be ignored in successive iterations.

In the practical implementation of the **#Gap** algorithm, instead of maintaining the pairs  $\langle b, m \rangle$ , we maintain two bit arrays  $Z^{(h-1)}, Z^{(h)}$  as in **Gap**, and an additional 2-bit array  $C$  containing the second component of the pairs. For such array  $C$  two bits per entry are sufficient since the values stored in each entry  $C[k]$  never decrease and they are no longer updated when they reach the value 2.

**Theorem 15.** *The #Gap algorithm merges two compressed permuterm indices  $\text{cbwt}_0$  and  $\text{cbwt}_1$  in  $\mathcal{O}(n \text{ aveCLcp})$  time, where  $n = n_0 + n_1$ . The working space is  $6n + \mathcal{O}(\log n)$  bits plus the space required for handling irrelevant blocks.*

**Proof:** We reason as in the proof of Theorem 13 except that if  $\text{rot}_0(\text{csa}_0[k]) = \text{rot}_1(\text{csa}_1[k+1])$  we are guaranteed that the corresponding size-2 block will become irrelevant only after iteration  $h = 2 \lceil \text{rot}_0(\text{csa}_0[k]) \rceil = 2 \text{clcp}_0[k+1]$ . Hence, the cost of handling  $\text{rot}_0(\text{csa}_0[k])$  is still proportional to (15) and the overall cost of the algorithm is  $\mathcal{O}(n \text{ aveCLcp})$  time. The space usage is the same as in Theorem 13, except for the  $2n$  additional bits for the  $C$  array.  $\square$

## References

- [1] M. Burrows, D. Wheeler, A block-sorting lossless data compression algorithm, Tech. Rep. 124, Digital Equipment Corporation (1994).
- [2] P. Ferragina, G. Manzini, Indexing compressed text, *Journal of the ACM* 52 (4) (2005) 552–581.
- [3] A. J. Cox, F. Garofalo, G. Rosone, M. Sciortino, Lightweight LCP construction for very large collections of strings, *J. Discrete Algorithms* 37 (2016) 17–33.
- [4] S. Mantaci, A. Restivo, G. Rosone, M. Sciortino, An extension of the Burrows-Wheeler transform, *Theor. Comput. Sci.* 387 (3) (2007) 298–312.
- [5] M. J. Bauer, A. J. Cox, G. Rosone, Lightweight algorithms for constructing and inverting the BWT of string collections, *Theor. Comput. Sci.* 483 (2013) 134–148.
- [6] S. Bonomo, S. Mantaci, A. Restivo, G. Rosone, M. Sciortino, Sorting conjugates and suffixes of words in a multiset, *Int. J. Found. Comput. Sci.* 25 (8) (2014) 1161.
- [7] P. Ferragina, R. Venturini, The compressed permuterm index, *ACM Trans. Algorithms* 7 (1) (2010) 10:1–10:21.
- [8] P. Ferragina, F. Luccio, G. Manzini, S. Muthukrishnan, Structuring labeled trees for optimal succinctness, and beyond, in: *Proc. 46th IEEE Symposium on Foundations of Computer Science (FOCS)*, 2005, pp. 184–193.
- [9] P. Ferragina, F. Luccio, G. Manzini, S. Muthukrishnan, Compressing and indexing labeled trees, with applications, *J. ACM* 57 (1) (2009) 4:1–4:33.
- [10] A. Bowe, T. Onodera, K. Sadakane, T. Shibuya, Succinct de Bruijn graphs, in: *WABI*, Vol. 7534 of *Lecture Notes in Computer Science*, Springer, 2012, pp. 225–235.
- [11] M. D. Muggli, B. Alipanahi, C. Boucher, Building large updatable colored de Bruijn graphs via merging, *Bioinformatics* 35 (14) (2019) i51–i60. doi:10.1093/bioinformatics/btz350.
- [12] J. Sirén, Indexing variation graphs, in: *Proc. 19th Meeting on Algorithm Engineering and Experiments (ALENEX '17)*, SIAM, 2017, pp. 13–27.
- [13] J. C. Na, H. Kim, H. Park, T. Lecroq, M. Lonard, L. Mouchard, K. Park, FM-index of alignment: A compressed index for similar strings, *Theoretical Computer Science* 638 (2016) 159–170.
- [14] J. C. Na, H. Kim, S. Min, H. Park, T. Lecroq, M. Léonard, L. Mouchard, K. Park, FM-index of alignment with gaps, *Theoretical Computer Science* 710 (2018) 148–157. doi:10.1016/j.tcs.2017.02.020.

- [15] T. Gagie, G. Manzini, J. Sirén, Wheeler graphs: A framework for bwt-based data structures, *Theor. Comput. Sci.* 698 (2017) 67–78.
- [16] D. Belazzougui, Linear time construction of compressed text indices in compact space, in: *STOC*, ACM, 2014, pp. 148–193.
- [17] J. Fuentes-Sepúlveda, G. Navarro, Y. Nekrich, Space-efficient computation of the Burrows-Wheeler Transform, in: *Proc. 29th Data Compression Conference (DCC)*, 2019, pp. 132–141.
- [18] J. I. Munro, G. Navarro, Y. Nekrich, Space-efficient construction of compressed indexes in deterministic linear time, in: *SODA*, SIAM, 2017, pp. 408–424.
- [19] H. Li, Fast construction of FM-index for long sequence reads, *Bioinformatics* 30 (22) (2014) 3274–3275.
- [20] J. Sirén, Compressed suffix arrays for massive data, in: *Proc. 16th Int. Symp. on String Processing and Information Retrieval (SPIRE '09)*, Springer Verlag LNCS n. 5721, 2009, pp. 63–74.
- [21] J. Sirén, Burrows-Wheeler transform for Terabases, in: *IEEE Data Compression Conference (DCC)*, 2016, pp. 211–220.
- [22] J. Holt, L. McMillan, Merging of multi-string BWTs with applications, *Bioinformatics* 30 (24) (2014) 3524–3531.
- [23] J. Holt, L. McMillan, Constructing Burrows-Wheeler transforms of large string collections via merging, in: *BCB*, ACM, 2014, pp. 464–471.
- [24] M. Léonard, L. Mouchard, M. Salson, On the number of elements to reorder when updating a suffix array, *J. Discrete Algorithms* 11 (2012) 87–99. doi:10.1016/j.jda.2011.01.002.
- [25] G. Navarro, V. Mäkinen, Compressed full-text indexes, *ACM Computing Surveys* 39 (1).
- [26] S. Gog, E. Ohlebusch, Compressed suffix trees: Efficient computation and storage of LCP-values, *ACM Journal of Experimental Algorithmics* 18. doi:10.1145/2444016.2461327.
- [27] J. Kärkkäinen, D. Kempa, LCP array construction in external memory, *ACM Journal of Experimental Algorithmics* 21 (1) (2016) 1.7:1–1.7:22.
- [28] F. A. Louza, G. P. Telles, C. D. A. Ciferri, External memory generalized suffix and LCP arrays construction, in: *CPM*, Vol. 7922 of *Lecture Notes in Computer Science*, Springer, 2013, pp. 201–210.
- [29] P. Bonizzoni, G. D. Vedova, Y. Pirola, M. Previtali, R. Rizzi, Multithread multistring burrows-wheeler transform and longest common prefix array, *Journal of Computational Biology* 26 (9) (2019) 948–961. doi:10.1089/cmb.2018.0230.

- [30] W. Hon, C. Lu, R. Shah, S. Thankachan, Succinct indexes for circular patterns, in: Proc. 22nd International Symposium on Algorithms and Computation, (ISAAC '11), 2011, pp. 673–682. doi:10.1007/978-3-642-25591-5\_69.
- [31] L. Egidi, F. A. Louza, G. Manzini, G. P. Telles, External memory BWT and LCP computation for sequence collections with applications, Algorithms for Molecular Biology 14 (1) (2019) 6:1–6:15. doi:10.1186/s13015-019-0140-0.
- [32] L. Egidi, F. A. Louza, G. Manzini, Space-efficient merging of succinct de Bruijn graphs, in: SPIRE, Vol. 11811 of Lecture Notes in Computer Science, Springer, 2019, pp. 337–351.
- [33] C. Boucher, A. Bowe, T. Gagie, S. J. Puglisi, K. Sadakane, Variable-order de Bruijn graphs, in: DCC, IEEE, 2015, pp. 383–392.
- [34] L. Egidi, G. Manzini, Lightweight BWT and LCP merging via the Gap algorithm, in: SPIRE, Vol. 10508 of Lecture Notes in Computer Science, Springer, 2017, pp. 176–190.
- [35] F. A. Louza, S. Gog, G. P. Telles, Induced suffix sorting for string collections, in: DCC, IEEE, 2016, pp. 43–52.
- [36] J. Kärkkäinen, G. Manzini, S. Puglisi, Permuted longest-common-prefix array, in: Proc. 20th Symposium on Combinatorial Pattern Matching (CPM), Springer-Verlag, LNCS n. 5577, 2009, pp. 181–192.
- [37] D. E. Knuth, Sorting and Searching, 2nd Edition, Vol. 3 of The Art of Computer Programming, Addison-Wesley, Reading, MA, USA, 1998.
- [38] G. Manzini, XBWT tricks, in: SPIRE, Vol. 9954 of Lecture Notes in Computer Science, 2016, pp. 80–92.
- [39] E. Ohlebusch, S. Stauß, U. Baier, Trickier XBWT tricks, in: SPIRE, Vol. 11147 of Lecture Notes in Computer Science, Springer, 2018, pp. 325–333.
- [40] H. Wilf, N. Fine, Uniqueness theorem for periodic functions, Proc. Amer. Math. Soc. 16 (1965) 109–114.
- [41] W. Hon, T. Ku, C. Lu, R. Shah, S. V. Thankachan, Efficient algorithm for circular burrows-wheeler transform, in: CPM, Vol. 7354 of Lecture Notes in Computer Science, Springer, 2012, pp. 257–268.