

Technical Report TR-INF-2004-02-03-UNIPMN

Two space saving tricks for linear time LCP computation

Giovanni Manzini*

February 10, 2004

Abstract

In this paper we consider the linear time algorithm of Kasai *et al.* [10] for the computation of the LCP array given the text and the suffix array. We show that this algorithm can be implemented without any auxiliary array in addition to the ones required for the input (the text and the suffix array) and the output (the LCP array). Thus, for a text of length n , we reduce the space occupancy of this algorithm from $13n$ bytes to $9n$ bytes.

We also consider the problem of computing the LCP array “overwriting” the suffix array. For this problem we propose an algorithm whose space occupancy depends on the regularity of the text. Experiments show that for linguistic texts our algorithm uses roughly $7n$ bytes. Our algorithm makes use of the Burrows-Wheeler Transform even if it does not represent any data in compressed form. To our knowledge this is the first application of the Burrows-Wheeler Transform outside the domain of data compression.

1 Introduction

The *suffix array* [14] is a simple and elegant data structure used for several fundamental string matching problems involving both linguistic texts and biological data. Another important application of the suffix array is the computation of the Burrows-Wheeler Transform [3] which is a powerful tool used for data compression and for the construction of compressed indexes [4, 5, 7, 18]. The vitality of this data structure is proven by the large number of suffix array construction algorithms developed in the last two years [2, 9, 11, 12, 16].

The suffix array of a text $t[1, n]$ is the lexicographically sorted list of all its suffixes. The suffix array is often used together with the *LCP array* which contains the length of the longest common prefix between every pair of lexicographically consecutive suffixes. The LCP information can be used to speed up suffix array algorithms and to simulate the more powerful, but more resource consuming, suffix tree data structure [1, 10, 14].

*Dipartimento di Informatica, Università del Piemonte Orientale, Italy. Email: manzini@mf.n.unipmn.it. Partially supported by the Italian MIUR projects “Algorithmics for Internet and the Web (ALINWEB)” and “Technologies and Services for Enhanced Content Delivery (ECD)”.

In [10] Kasai et al. describe a simple (13 lines of C code) and elegant linear time algorithm for computing the LCP array given the text and the suffix array. This was an important result for several reasons. First, although many suffix array construction algorithms can be modified to return the LCP array as well, this is not true for every algorithm. Having decoupled the two problems allows one to choose the suffix array construction algorithm which better suits his/her needs without the constraint of considering only those algorithms which also provide the LCP array. Moreover, in some applications one may need the LCP array later than the suffix array: if one has to compute them simultaneously some temporary storage must be used for the LCP array. Another advantage of the separate construction of these arrays is that it leads to simpler algorithms with (possibly) an improved locality of reference.

The only drawback of the algorithm of Kasai et al. is its large space occupancy. Assuming a “real world” model in which each text symbol takes one byte and each suffix array or LCP array entry takes 4 bytes, the algorithm of Kasai et al. uses $13n$ bytes, where n is the length of the input text. Considering that the output of the computation (text, suffix array, and LCP array) takes $9n$ bytes, we have a $4n$ bytes overhead which is a serious issue since it is nowadays common to work with files hundreds of megabytes long. Indeed, space occupancy is currently the main bottleneck for this and other string matching data structures: a prototypical example is the suffix tree data structure [17] which provides an optimal solution for many string matching problems but it is seldom used in practice because of its huge space occupancy and because of the even larger amount of space required for its construction.

In this paper we present a modified version of the algorithm of Kasai et al. which only uses $9n$ bytes of storage. Our algorithm, called **Lcp9**, runs in linear time and has the same simplicity and elegance of the original algorithm. Experiments with several files of different size and structure show that **Lcp9** is only 5%–10% slower than the algorithm of Kasai et al..

In our “real world” model, a space occupancy of $9n$ bytes is optimal if we assume that at the end of the computation we need the text, the suffix array, and the LCP array. However, this is no longer true if one is interested *only* in the LCP array, that is, if at the end of the computation we no longer need the suffix array. In this case, the space initially used for storing the suffix array can be reused during the computation and for the storage of the LCP array. In this scenario we can aim to a space occupancy as low as $5n$ bytes. The problem of computing the LCP array discarding the suffix array has applications in the fields of string matching, data compression and text analysis. For example, using the algorithm described in [10, Sect. 5] with a single pass over the LCP array we can simulate a post order visit of the suffix tree of the text t . In some applications, for example for the construction of the compression booster described in [6], such visit does not need the information stored in the suffix array.

If we only need the LCP array, even the $9n$ bytes space occupancy of algorithm **Lcp9** becomes the space bottleneck of the whole computation since for the construction of the suffix array there are “lightweight” algorithms [2, 16] which only use $(5 + \epsilon)n$ bytes with $\epsilon \ll 1$. In this paper we address this issue proposing a simple linear time algorithm, called **Lcp6**, which computes the LCP array “overwriting” the suffix array. The space used by **Lcp6** depends on the regularity of the input text $t[1, n]$. If t is highly compressible the space occupancy of **Lcp6** can be as small as $6n$ bytes. Vice versa, if t is a “random” string the space required by our algorithm can be as large as $10n$ bytes. Note however that in the first step of **Lcp6** we can evaluate exactly how much space it will need: if such space turns out to be larger than $9n$ bytes we can quit **Lcp6** and compute the LCP array using **Lcp9**. Thus, combining **Lcp6** and **Lcp9** we get an algorithm with a space occupancy between $6n$ and $9n$ bytes. The experimental results show that for linguistic texts,

source code, and xml/html documents `Lcp6` always uses less than $8n$ bytes and for the largest files it often uses less than $7n$ bytes. For DNA sequences `Lcp6` uses between $8n$ and $9n$ bytes, and—not surprisingly—for compressed files it uses close to $10n$ bytes. The experimental results also show that `Lcp6` is roughly two times slower than the algorithm of Kasai et al..

One of the crucial ingredients of algorithm `Lcp6` is the Burrows-Wheeler Transform. However, `Lcp6` does not represent any data in compressed form. Instead, it makes use of a structural property of the transformed text to reduce the amount of auxiliary information needed for computing the LCP array. To our knowledge this is the first use of the Burrows-Wheeler Transform outside the domain of data compression.

Finally, we point out that both our algorithms only have a “practical” interest: from the theoretical point of view their working space of $\Theta(n \log n)$ bits is not optimal. Indeed, the optimal space/time tradeoff can be obtained combining the results in [8] and [18] which allow one to build the suffix array and LCP array in linear time using $O(n)$ bits of auxiliary storage. Unfortunately the algorithms in [8, 18] are quite complex and it is still unclear whether they will lead to competitive practical algorithms.

2 Background and notation

Let Σ denote a finite ordered alphabet. Without loss of generality, in the following we assume that Σ consists of the integers $1, 2, \dots, |\Sigma|$. Let $\mathbf{t}[1, n]$ denote a text over Σ . For $i = 1, \dots, n$ we write $\mathbf{t}[i, n]$ to denote the suffix of \mathbf{t} of length $n - i + 1$ that is $\mathbf{t}[i, n] = \mathbf{t}[i]\mathbf{t}[i + 1] \cdots \mathbf{t}[n]$.

The suffix array [14] for \mathbf{t} is the array $\text{SA}[1, n]$ such that $\mathbf{t}[\text{SA}[1], n], \mathbf{t}[\text{SA}[2], n], \dots, \mathbf{t}[\text{SA}[n], n]$ is the list of suffixes of \mathbf{t} sorted in lexicographic order. To define unambiguously the lexicographic order of the suffixes it is customary to logically append at the end of \mathbf{t} a special end-of-string symbol $\#$ which is smaller than any symbol in Σ . For example, for $\mathbf{t} = \text{baaba}$, $\text{SA} = [5, 2, 3, 4, 1]$ since $\mathbf{t}[5, 5] = \text{a}$ is the suffix with the lowest lexicographic rank, followed by $\mathbf{t}[2, 5] = \text{aaba}$, followed by $\mathbf{t}[3, 5] = \text{aba}$ and so on.

The rank array $\text{RANK}[1, n]$ of \mathbf{t} is the inverse of the suffix array. That is, $\text{RANK}[i] = j$ if and only if $\text{SA}[j] = i$. Note that $\text{RANK}[i]$ is the rank of the suffix $\mathbf{t}[i, n]$ in the lexicographic order of the suffixes. The LCP array $\text{LCP}[1, n]$ of \mathbf{t} is an array such that $\text{LCP}[i]$ contains the length of the longest common prefix between the suffix $\mathbf{t}[\text{SA}[i], n]$ and its predecessor in the lexicographic order (which is $\mathbf{t}[\text{SA}[i - 1], n]$). Note that $\text{LCP}[1]$ is undefined since $\mathbf{t}[\text{SA}[1], n]$ is the lexicographically smallest suffix and therefore it has no predecessor.

Finally, we define the `RANKNEXT` map such that:

$$\text{RANKNEXT}(i) = \text{RANK}[\text{SA}[i] + 1], \quad \text{for } i = 1, \dots, n, i \neq \text{RANK}[n]. \quad (1)$$

$\text{RANKNEXT}(i)$ is the rank of the suffix $\mathbf{t}[\text{SA}[i] + 1, n]$, that is, the rank of the suffix obtained removing the first character from the suffix of rank i . Note that $\text{RANKNEXT}(\cdot)$ is not defined for $i = \text{RANK}[n]$ because in this case $\mathbf{t}[\text{SA}[i] + 1, n]$ is the empty string.

2.1 The Burrows-Wheeler Transform

In 1994, Burrows and Wheeler [3] introduced a transform that turns out to be very elegant in itself and extremely useful for data compression. Given a string \mathbf{t} , the transform consists of three basic steps (see Fig. 1): (1) append to the end of \mathbf{t} a special symbol $\#$ smaller than any other symbol

mississippi#		# mississippi i
ississippi#m		i #mississipp
ssissippi#mi		i ppi#missis s
sissippi#mis		i ssiippi#mis s
issippi#miss		i ssissippi# m
ssippi#missi	⇒	m ississippi #
sippi#missis		p i#mississi p
ippi#mississ		p pi#mississ i
ppi#mississi		s ippi#missi s
pi#mississip		s issippi#mi s
i#mississipp		s sippi#miss i
#mississippi		s sissippi#m i

Figure 1: The Burrows-Wheeler Transform for the string $t = \text{mississippi}$. The matrix on the right has the rows sorted in lexicographic order. The output of the Burrows-Wheeler Transform is the last column of the matrix, i.e., $\text{bwt} = \text{ipssm\#pissii}$.

in Σ ; (2) form a *conceptual* matrix \mathcal{M} whose rows are the cyclic shifts of the string $t\#$, sorted in lexicographic order; (3) construct the transformed text bwt by taking the last column of \mathcal{M} . Notice that every column of \mathcal{M} , hence also the transformed text bwt , is a permutation of $t\#$.

If the input string t has length n , the transformed string bwt has length $n + 1$ because of the presence of the $\#$ symbol. In the following we assume that the transformed string is stored in an array indexed from 0 to n . For example, in Fig. 1 we have $\text{bwt}[0] = \text{i}$, $\text{bwt}[5] = \#$, $\text{bwt}[11] = \text{i}$. Using this notation and observing that the rows of the matrix \mathcal{M} are precisely the suffixes of t in lexicographic order, the computation of bwt given t and the suffix array can be easily accomplished with the code of Fig. 2 (procedure `Sa2Bwt`).

In [3] Burrows and Wheeler proved that from bwt we can always recover t . The inverse transform is based on the following remarkable property. Let $F[0, n]$ and $L[0, n]$ denote respectively the first and last column of the matrix \mathcal{M} (hence, $L \equiv \text{bwt}$). Then, for any $\sigma \in \Sigma$ we have that the k -th occurrence of σ in F corresponds to the k -th occurrence of σ in L . For example, in Fig. 1 we have that the second i in F (that is, $F[2]$) corresponds to the second i in L (that is, $L[7]$) since they both are the eighth symbol of `mississippi`. Similarly, the third s in F ($F[10]$) corresponds to the third s in L ($L[8]$) since they both are the sixth symbol of `mississippi`.

Assume now that the character $F[j]$ corresponds to $L[i]$. This means that row i of \mathcal{M} consists of a (rightward) cyclic shift of row j . Because of the relationship between rows of \mathcal{M} and suffixes of t this is equivalent to stating that the i -th suffix in the lexicographic order is equal to the j -th suffix with the first symbol removed. In terms of the map `RANKNEXT` defined by (1) we have `RANKNEXT(j) = i`. From this latter relationship it follows that from bwt we can obtain the `RANKNEXT` map. Indeed, we only need to scan the array bwt (which coincides with column L) finding, for $i = 1, \dots, n$ the character $F[j]$ corresponding to $\text{bwt}[i] \equiv L[i]$. The resulting code is given in Fig. 2 (procedure `Bwt2RankNext`). Note that column F is not represented explicitly (since it would take $O(n)$ space). Instead we use the array `count[1, |\Sigma|]`: at the beginning of the i -th iteration `count[k]` contains the number of occurrences in t of the characters $1, 2, \dots, k - 1$ plus the number of occurrences of character k in $\text{bwt}[0] \cdots \text{bwt}[i - 1]$.

Given the `RANKNEXT` map and the array bwt , we can recover t as follows. The position of the end-of-string symbol $\#$ in bwt gives us `RANK(1)`, that is, the position of $t[1, n]$ in the suffix array. By (1), setting $i = \text{RANK}(j)$ we get

$$\text{RANKNEXT}(\text{RANK}(j)) = \text{RANK}(\text{SA}[\text{RANK}(j)] + 1) = \text{RANK}(j + 1). \quad (2)$$

Procedure Sa2Bwt

```
1. bwt[0]=t[n];
2. for(i=1;i<=n;i++) {
3.   if(sa[i] == 1)
4.     bwt[i]='#';
5.   else
6.     bwt[i]=t[sa[i]-1];
7. }
```

Procedure Bwt2RankNext

```
1. for(i=0;i<=n;i++) {
2.   c = bwt[i];
3.   if(c == '#')
4.     eos_pos = i;
5.   else {
6.     j = count[c]++;
7.     rank_next[j]=i;
8.   }
9. }
10. return eos_pos;
```

Procedure RankNext2Text

```
1. k = eos_pos; i=1;
2. do {
3.   k = rank_next[k];
4.   t[i++] = bwt[k];
5. } while(k!=0);
```

Procedure RankNext2SuffixArray

```
1. k = eos_pos; i=1;
2. while(k!=0) {
3.   nextk = rank_next[k];
4.   sa[k] = i++;
5.   k = nextk;
6. }
```

Figure 2: Algorithms related to the Burrows-Wheeler Transform. Procedure `Sa2Bwt` computes the array `bwt` given the text `t` and the suffix array `sa`. Procedure `Bwt2RankNext` stores in `rank_next` the RANKNEXT map and returns the value `RANK(1)`. The procedure uses the auxiliary array `count[1, |Σ|]` which initially contains in `count[i]` the number of occurrences in `bwt` (and therefore in `t`) of the characters $1, \dots, i-1$. Procedure `RankNext2Text` recovers the text `t` given the arrays `bwt` and `rank_next` and the value `RANK(1)` stored in `eos_pos`. Procedure `RankNext2Sa` computes the suffix array given `rank_next` and the value `RANK(1)` stored in `eos_pos`.

Hence, given `RANKNEXT` and `RANK(1)` we can generate the sequence of values `RANK(2), RANK(3), \dots, RANK(n)` using the recurrence

$$\text{RANK}(2) = \text{RANKNEXT}(\text{RANK}(1)), \quad \text{RANK}(3) = \text{RANKNEXT}(\text{RANK}(2)), \quad \dots \quad (3)$$

From the sequence `RANK(1), RANK(2), \dots, RANK(n)` we recover `t` using the relationship $t[i] = \text{bwt}[\text{RANK}(i+1)]$. The corresponding code is shown in Fig. 2 (procedure `RankNext2Text`).

We conclude this section observing that from the sequence `RANK(1), RANK(2), \dots, RANK(n)` we can also recover the suffix array since $k = \text{RANK}(i)$ implies $\text{SA}[k] = i$. The corresponding code is shown in Fig. 2 (procedure `RankNext2SuffixArray`). Note that in `RankNext2SuffixArray` as soon as we have read `rank_next[k]` in Line 3 that entry is no longer needed. Therefore, if we replace Line 4 of `RankNext2SuffixArray` with the instruction `rank_next[k] = i++`; we get a procedure which stores the suffix array entries in the array `rank_next` overwriting the old content of the array (the RANKNEXT map). This property will be used in Section 4.

2.2 The algorithm of Kasai et al.

The algorithm of Kasai et al. (algorithm `Lcp13` from now on) takes as input the text `t[1, n]` and the corresponding suffix array `SA[1, n]` and returns the LCP array. For $i = 1, \dots, n$ let ℓ_i denote the LCP between `t[i, n]` and the suffix immediately preceding it in the lexicographic order (ℓ_i is

Procedure Lcp13

```
1. for(i=1;i<=n;i++) rank[sa[i]] = i;
2. h=0;
3. for(i=1;i<=n;i++) {
4.   k = rank[i];
5.   if(k==1) lcp[k]=-1;
6.   else {
7.     j = sa[k-1];
8.     while(i+h<=n && j+h<=n && t[i+h]==t[j+h]):
9.       h++;
10.    lcp[k] = h;
11.  }
12.  if(h>0) h--;
13. }
```

Figure 3: Algorithm of Kasai et al. for the linear time computation of the LCP array. The algorithm takes as input the text t and the suffix array sa and stores in lcp the LCP array. The algorithm uses an auxiliary array $rank$ to store the rank array (which is the inverse of the suffix array).

undefined when $t[i, n]$ is the lexicographically smallest suffix). The algorithm **Lcp13** computes the LCP values in the order $\ell_1, \ell_2, \dots, \ell_n$.

The code of **Lcp13** is shown in Fig. 3. As a first step (Line 1) the algorithm computes the rank array $RANK[1, n]$. Then, at the i -th iteration of the main loop (Lines 3–13) **Lcp13** computes ℓ_i as follows. At Line 4 the value $RANK[i]$ is stored in the variable k . If $RANK[i] = 1$ then $t[i, n]$ is the smallest suffix in the lexicographic order and ℓ_i is undefined (we set it to -1 at Line 5). If $RANK[i] > 1$, we compute $j = SA[RANK[i] - 1]$ (Line 7). $t[j, n]$ is the suffix preceding $t[i, n]$ in the lexicographic order, hence ℓ_i is the longest common prefix between $t[i, n]$ and $t[j, n]$.

The crucial observation, which ensures that **Lcp13** runs in $O(n)$ time, is that whenever ℓ_i and ℓ_{i-1} are both defined we have $\ell_i \geq \ell_{i-1} - 1$ (Theorem 1 in [10]). To use this property **Lcp13** maintains the invariant that at the beginning of the i -th iteration the variable h contains the value $\ell_{i-1} - 1$. Hence, ℓ_i is computed comparing $t[i, n]$ and $t[j, n]$ starting from position h (Lines 8–9). Note that at Line 10 **Lcp13** stores ℓ_i in $LCP[RANK[i]]$ since the definition of LCP array states that $LCP[t]$ contains the LCP between $t[SA[t], n]$ and $t[SA[t - 1], n]$.

In our “real world” model, algorithm **Lcp13** requires n bytes for the array t and $4n$ bytes for each one of the arrays SA , $RANK$, and LCP . Therefore its peak space occupancy is $13n$ bytes.

3 LCP computation in $9n$ bytes of storage

In this section we show how to modify the algorithm of Kasai et al. for computing the LCP array in linear time without using any auxiliary array. As a result get an algorithm which only uses $9n$ bytes of storage. This amount is the minimum possible if we assume that at the end of the computation we want an explicit representation of the text, the suffix array, and the LCP array. Our approach consists in using the array lcp for storing both “rank information” and “LCP information”. Initially the array contains only “rank information”. Then, at each iteration of the main loop one item of rank information is used and replaced by one item of LCP information. At the end of the computation the array lcp only contains LCP information.

<pre> Procedure Lcp9 1. k = Sa2RankNext(lcp); 2. h=0; 3. for(i=1;i<=n;i++) { 4. nextk = lcp[k]; 5. if(k==1) lcp[k]=-1; 6. else { 7. j = sa[k-1]; 8. while(i+h<=n && j+h<=n && t[i+h]==t[j+h]) 9. h++; 10. lcp[k] = h; 11. } 12. if(h>0) h--; 13. k=nextk; 14. }</pre>	<pre> Procedure Sa2RankNext(rank_next) 1. j = count[t[n]]++; 2. rank_next[j]=0; 3. for(i=1;i<=n;i++) { 4. if(sa[i]==1) 5. eos_pos = i; 6. else { 7. c = t[sa[i]-1]; 8. j = count[c]++; 9. rank_next[j]=i; 10. } 11. } 12. return eos_pos;</pre>
---	---

Figure 4: Algorithm `Lcp9` for linear time computation of the LCP array using $9n$ bytes of storage. The algorithm takes as input the text t and the suffix array sa and stores in lcp the LCP array. The procedure `Sa2RankNext` computes the `RANKNEXT` map given t and sa . After the procedure call at Line 1 of `Lcp9` the `RANKNEXT` map is stored in the array lcp and the value `RANK(1)` is stored in the variable k .

Our starting point is the observation that algorithm `Lcp13` (Fig. 3) uses the rank information only in Line 4 where, during the i -th iteration of the main loop, the algorithm retrieves the value `RANK(i)`. Therefore, `Lcp13` uses the sequence of rank values `RANK(1)`, `RANK(2)`, \dots , `RANK(n)` exactly in this order. Moreover, after the i -th iteration of the main loop the value `RANK(i)` is no longer needed.

In Section 2.1 we have shown that using the recurrence (2) we can generate the sequence `RANK(1)`, `RANK(2)`, \dots , `RANK(n)` given the `RANKNEXT` map and the value `RANK(1)`. The above observations suggest the algorithm `Lcp9` whose code is shown in Fig. 4. In first step of `Lcp9` (Line 1) we call the procedure `Sa2RankNext` which, for $j = 1, \dots, n$, stores the value `RANKNEXT(j)` in `LCP[j]`, and returns the value `RANK(1)`. Then, in the i -th iteration of the main loop (Lines 3–14) given `RANK(i)` we retrieve `RANK($i + 1$)` from entry `LCP[RANK(i)]`. Note that as soon as we have retrieved `RANK($i + 1$)` we can use the entry `LCP[RANK(i)]` for storing the LCP relative to $t[i, n]$.

Summing up, the main loop of algorithm `Lcp9` (Lines 3–14) works as follows. At the beginning of the i -th iteration the variable k contains the value `RANK(i)`. In the body of the loop we store in `nextk` the value `lcp[k]` which is `RANK($i + 1$)`; then we compute ℓ_i (the LCP between $t[i, n]$ and the suffix preceding it) and we store it in `lcp[k]`, which is the right place since $k = \text{RANK}(i)$. Finally, we update k (Line 13) and we start the next iteration. Note that the actual computation of ℓ_i is done as in the `Lcp13` algorithm; indeed, lines 5–12 are identical in both algorithms. The only difference between our algorithm and the one of Kasai et al. is the computation of the rank information using the `RANKNEXT` map rather than the `rank` array.

We conclude observing that the correctness of the procedure `Sa2RankNext` follows from the correctness of `Bwt2RankNext` in Fig. 2 and by the relationship between the suffix array and the Burrows-Wheeler Transform (see the procedure `Sa2Bwt` in Fig. 2). Indeed, `Sa2RankNext` is a transposition of `Bwt2RankNext` in which `bwt[i]` has been replaced by $t[sa[i]-1]$ when $i > 1$ (Line 7) and by $t[n]$ when $i = 0$ (Line 1), and the test `bwt[i]=='#'` has been replaced by the test `sa[i]==1` (Line 4).

4 LCP computation in $(6 + \delta)n$ bytes of storage

In this section we describe the algorithm `Lcp6` which computes the LCP array “overwriting” the suffix array in the sense that the LCP array is stored in the same array which initially contains the suffix array entries.

For $i = 1, \dots, n$ let ℓ_i denote the LCP between $\mathbf{t}[i, n]$ and the suffix preceding it in the lexicographic order (ℓ_i is undefined when $\mathbf{t}[i, n]$ is the lexicographically smallest suffix). The correctness of the algorithm of Kasai et al. follows from the observation that whenever ℓ_i and ℓ_{i-1} are both defined we have $\ell_i \geq \ell_{i-1} - 1$ (see Section 2.2). We now show that using the Burrows-Wheeler Transform of \mathbf{t} we can say something more on the relationship between ℓ_i and ℓ_{i-1} .

Lemma 1 *Let \mathbf{bwt} denote the Burrows-Wheeler Transform of \mathbf{t} , and let $k = \text{RANK}(i)$. If $k > 1$ and $\mathbf{bwt}[k] = \mathbf{bwt}[k - 1]$ then $\ell_i = \ell_{i-1} - 1$.*

Proof: Let $\mathbf{t}[j, n]$ (resp. $\mathbf{t}[j', n]$) denote the suffix immediately preceding $\mathbf{t}[i, n]$ (resp. $\mathbf{t}[i - 1, n]$) in the lexicographic order. Since $k = \text{RANK}(i)$ we know that the suffixes $\mathbf{t}[j, n]$ and $\mathbf{t}[i, n]$ are in the rows $k - 1$ and k of the Burrows-Wheeler matrix \mathcal{M} . By hypothesis these rows end with the same character $\alpha = \mathbf{bwt}[k - 1] = \mathbf{bwt}[k]$. Since there is only one occurrence of $\#$ in \mathbf{bwt} we have that $\alpha \in \Sigma$. Hence, the strings $\alpha\mathbf{t}[j, n]$ and $\alpha\mathbf{t}[i, n]$ are both suffixes of \mathbf{t} .

By construction $\alpha\mathbf{t}[i, n]$ is equal to $\mathbf{t}[i - 1, n]$. Since $\mathbf{t}[j, n]$ immediately precedes $\mathbf{t}[i, n]$ in the lexicographic order, $\alpha\mathbf{t}[j, n]$ must be the suffix immediately preceding $\alpha\mathbf{t}[i, n] \equiv \mathbf{t}[i - 1, n]$. Hence, $j' = j - 1$. This means that the longest common prefix between $\mathbf{t}[i, n]$ and $\mathbf{t}[j, n]$ is equal to the longest common prefix between $\mathbf{t}[i - 1, n]$ and $\mathbf{t}[j - 1, n]$ with the first character removed. Thus, $\ell_i = \ell_{i-1} - 1$ as claimed. \blacksquare

Assume now that the array \mathbf{bwt} is available, and consider the main loop of `Lcp9` (Lines 3–14 in Fig. 4). At the beginning of the i -th iteration the variable \mathbf{k} contains the value $k = \text{RANK}(i)$. By Lemma 1, if $\mathbf{bwt}[k] = \mathbf{bwt}[k - 1]$ we know that $\ell_i = \ell_{i-1} - 1$. Since ℓ_{i-1} is stored in the variable \mathbf{h} , we conclude that, if $\mathbf{bwt}[k] = \mathbf{bwt}[k - 1]$, we can skip Lines 8–9 and proceed with the next iteration. This means that for computing the LCP array we only need the values $\text{SA}[k - 1]$'s for all k 's such that $\mathbf{bwt}[k] \neq \mathbf{bwt}[k - 1]$. This observation is the starting point of our algorithm.

Let z' denote the number of consecutive equal characters in \mathbf{bwt} and let $z = n - z'$. In the algorithm `Lcp6` (see Figure 5) we evaluate z with a scan of \mathbf{bwt} and we allocate an array `sa_aux` of size z for storing those suffix array entries that are needed for computing the LCP array (Lines 2–4). Although we already know which suffix array entries must be stored in `sa_aux`, to retrieve these entries efficiently we must store them in the proper order (recall that we are trying to use as small space as possible). Let k_1, k_2, \dots, k_z , with $k_1 < k_2 < \dots < k_z$ denote the indexes such that $\mathbf{bwt}[k_i] \neq \mathbf{bwt}[k_i - 1]$. By the above discussion we know that we must store in `sa_aux` the values $\text{SA}[k_1 - 1], \text{SA}[k_2 - 1], \dots, \text{SA}[k_z - 1]$. Note, however, that the value $\text{SA}[k_i - 1]$ is needed when we process the suffix $\mathbf{t}[\text{SA}[k_i], n]$. Since the main loop of the LCP algorithm considers the suffixes in the order $\mathbf{t}[1, n], \mathbf{t}[2, n], \dots, \mathbf{t}[n, n]$ in `Lcp6` we store in `sa_aux[i]` the value $\text{SA}[k_{\pi(i)} - 1]$ where π is a permutation of $1, \dots, z$ such that

$$\text{SA}[k_{\pi(1)}] < \text{SA}[k_{\pi(2)}] < \dots < \text{SA}[k_{\pi(z)}]. \quad (4)$$

In other words, we store in `sa_aux` the suffix array entries in the order in which they will be used by the LCP algorithm. This will make the retrieval a very simple task.

Algorithm Lcp6

```
1. // ----- count how many suffix array entries we need -----
2. for(z=0,i=2;i<=n;i++)
3.   if(bwt[i-1]!=bwt[i]) z++;
4. sa_aux=malloc(z*sizeof(int));           // allocate sa_aux[0,z-1]
5. // ----- determine order in which suffix array entries are needed -----
6. k = Bwt2RankNext(sa);                   // store RankNext in sa[]
7. for(v=0,i=2;i<=n;i++) {
8.   if(bwt[k-1]!=bwt[k]) sa_aux[v++]=k-1;
9.   k=lcp[k];
10. }
11. // ----- store needed suffix array entries in sa_aux -----
12. RankNext2Sa(sa);                       // store Suffix Array in sa[]
13. for(v=0;v<z;v++)
14.   sa_aux[v] = sa[sa_aux[v]];
15. // ----- compute the lcp array as usual -----
16. k = Bwt2RankNext(sa);                   // store RankNext in sa[]
17. v=h=0;
18. for(i=1;i<=n;i++) {
19.   nextk = sa[k];
20.   if(k==1) sa[k]=-1;
21.   else if(bwt[k-1]==bwt[k]) sa[k]=h;
22.   else {
23.     j = sa_aux[v++];                     // retrieve sa[k-1]
24.     while(i+h<=n && j+h<=n && t[i+h]==t[j+h])
25.       h++;
26.     sa[k] = h;
27.   }
28.   if(h>0) h--;
29.   k=nextk;
30. }
```

Figure 5: Algorithm Lcp6 for linear time computation of the LCP array using $(6 + \delta)n$ bytes of storage. The algorithm takes as input the text τ , the Burrows-Wheeler Transform bwt , and the suffix array sa and stores the LCP values in sa (thus overwriting the suffix array entries). The algorithm uses an auxiliary array sa_aux whose size depends on the structure of the input text. Note that after the procedure call at Lines 6 and 16 the RANKNEXT map is stored in the array sa and the value $\text{RANK}(1)$ is stored in the variable k . The procedure call at Line 12 has the effect of storing the suffix array information in the array sa overwriting the RANKNEXT map (see comment at the end of Sect. 2.1).

To obtain such a convenient arrangement of the suffix array entries within `sa_aux`, the algorithm `Lcp6` uses the following two-step procedure. In the first step (Lines 6–10) the algorithm computes the `RANKNEXT` map storing it in the array `sa`. Then, it generates the sequence `RANK(1), RANK(2), . . . , RANK(n)` thus traversing the suffix array entries in the order in which they will be considered by the LCP computation. When `Lcp6` finds an index k such that `bwt[k - 1] ≠ bwt[k]` it stores $k - 1$ in the next empty position of `sa_aux` (Line 8). Hence, at the end of this first step, for $i = 1, \dots, z$, the entry `sa_aux[i]` contains the value $k_{\pi(i)} - 1$, where π is the permutation defined by (4). In the second step (Lines 12–14) the algorithm recomputes the suffix array and, with a simple scan over `sa_aux`, stores in `sa_aux[i]` the value `SA[kπ(i) - 1]`. Note that we use this elaborate two step procedure simply because we do not want to store at the same time both the suffix array and the `RANKNEXT` map.

Once the array `sa_aux` is properly initialized, the computation of the LCP array proceeds as in algorithm `Lcp9`. First, we store the `RANKNEXT` map in the array `sa` (Line 16). Then, at each iteration of the main loop (Lines 18–30) a `RANKNEXT` value in `sa` is replaced by a LCP value so that at the end of the loop `sa` contains the LCP array. The computation of the value ℓ_i makes use of Lemma 1. At the beginning of the i -th iteration the variable `k` contains `RANK(i)`; if `bwt[k-1]==bwt[k]` then ℓ_i is equal to $\ell_{i-1} - 1$ (which is readily available since it is stored in the variable `h`); otherwise we retrieve from `sa_aux` the value `SA[k - 1]` (Line 23) and we compute ℓ_i with the `while` loop of Lines 24–25.

In our “real world” model the total space occupancy of the above algorithm is $6n + 4z$ bytes: we use $2n$ bytes for the arrays `t` and `bwt`, $4n$ bytes for the array `sa` (which is used for storing the suffix array, the `RANKNEXT` map, and the LCP array), and $4z$ bytes for `sa_aux`. This latter amount depends on the structure of the input. It is well known that for linguistic texts and other “structured” texts the Burrows-Wheeler Transform usually contains many repetitions and consequently z is relatively small. For example, if $z \approx n/2$ (which is not an unusual value) the total space occupancy of `Lcp6` is $\approx 8n$ bytes. On the other hand, in the worst case we have $z = n$ and our algorithm uses $10n$ bytes. However, if at Line 4 we find that $z > 3n/4$ —which would yield a space occupancy larger than $9n$ bytes—we can quit `Lcp6` and use `Lcp9` instead.

We conclude observing that the entries of `sa_aux` are always accessed in sequential order. Hence, is one is really tight on space, such array could be stored in secondary memory with only a “reasonable” slowdown. Note that such option is not available for `Lcp9` which accesses the three arrays `t`, `sa`, and `lcp` in random order.

5 Experimental results

In this section we report the results of an experimental comparison between the algorithms `Lcp13`, `Lcp9`, and `Lcp6`. We ran these algorithms on a collection of files with different lengths and structures (see Table 1). For all tests we used a 1700 Mhz Pentium 4 running GNU/Linux with 1.25Gb main memory and 256Kb L2 cache. The compiler was `gcc` ver. 3.2 with options `-O3 -fomit-frame-pointer -march=pentium4`.

For each file we built the suffix array using the `ds` algorithm [15, 16] which is currently one of the fastest suffix array construction algorithm (at least for non-pathological inputs). Then, the text and the suffix array were given as input to the algorithms `Lcp13`, `Lcp9`, and `Lcp6` and their running times were measured considering (user+system) time averaged over five runs.

In Table 2 we report, for each file and for each algorithm, running time over file length. In other words the table shows the average time (in microseconds) per input byte spent by each

File	Size (Kb)	Ave. LCP	Description
<i>calgary.zip</i>	1,043	2.00	Zip archive containing the files of the Calgary corpus
<i>war&peace</i>	3,142	9.45	“War and Peace” novel from Project Gutenberg
<i>texbook</i>	1,351	10.87	T _E X source of Knuth’s T _E Xbook
<i>bible</i>	3,952	13.97	The King James version of the Bible
<i>ecoli</i>	4,529	17.38	Complete genome of the E.Coli bacterium
<i>world192</i>	2,415	23.01	The 1992 CIA world fact book
<i>pic</i>	501	2,353.32	Black and white bitmap of an image in the CCITT test set.

File	Size (Kb)	Ave. LCP	Description
<i>etext99.gz</i>	38,747	2.65	The file <i>etext99</i> (see below) in gzipped format
<i>sprot</i>	107,048	89.08	Swiss prot database (original file name <i>sprot34.dat</i>)
<i>rfc</i>	113,693	93.02	Concatenation of RFC text files
<i>howto</i>	38,498	267.56	Concatenation of Linux Howto text files
<i>reuters</i>	112,022	282.07	Reuters news in XML format
<i>linux</i>	113,530	479.00	Tar archive containing the Linux kernel 2.4.5 source files
<i>jdk13</i>	68,094	678.94	Concatenation of <i>html</i> and <i>java</i> files from the JDK 1.3 doc.
<i>etext99</i>	102,809	1,108.63	Concatenation of Project Gutenberg <i>etext99/*</i> .txt files
<i>chr22</i>	33,743	1,979.25	Genome assembly of human chromosome 22
<i>gcc</i>	84,600	8,603.21	Tar archive containing the gcc 3.0 source files
<i>w3c</i>	101,759	42,299.75	Concatenation of <i>html</i> files from www.w3c.org

Table 1: Collection of small files (top) and large files (bottom) used in our experiments. Files are sorted in order of increasing average LCP.

algorithm. For **Lcp6** we also report the space occupancy expressed as total space occupancy over file length. The files in Table 2 are ordered by increasing average LCP: a large average LCP indicates that the input file contains many long repeated substrings. Note that the files *calgary.zip* and *etext99.gz* have a very small average LCP: the reason is that these are compressed files and therefore their content has very little regularity and essentially consists of a “random” sequence over the alphabet $\{0, 1, \dots, 255\}$. The files *ecoli* and *chr22* are DNA sequences and consist of apparently random strings over the alphabet $\{a, c, g, t\}$: their relatively high average LCP is due to the small cardinality of the underlying alphabet. Another file with an unusual structure is *pic* which is the bitmap of a black and white image and contains long runs of zeroes.

Our first observation is that **Lcp9** is roughly 10% slower than **Lcp13**. We also notice that for most files both LCP algorithms are faster than the suffix array construction algorithm. Thus, if we consider the combined time required to compute the suffix array and the LCP array, the overhead for using **Lcp9** is usually less than 5% of the total running time. We stress once more that the $13n$ space occupancy of algorithm **Lcp13** is the space bottleneck of the combined—suffix array + LCP array—computation since there are many efficient suffix array construction algorithms which use $9n$ bytes or less (see for example [2, 13, 16]).

For the algorithm **Lcp6** we observe that it is roughly two times slower than **Lcp13**. However, we also notice that for most files **Lcp6** uses less than $8n$ bytes. The exceptions are, as expected, the compressed files (*calgary.zip* and *etext99.gz*) and the DNA sequences (*ecoli* and *chr22*). We conclude that, although **Lcp6** is slower than **Lcp13** and **Lcp9**, for most files it yields a significant saving in the peak space occupancy. When working with very large files the combination of **Lcp6** with a “lightweight” suffix sorter [2, 16] can be the only way to avoid the (deleterious) use of secondary memory.

File	Size (Kb)	Ave. LCP	SA time	Lcp13 time	Lcp9 time	Lcp6 time	Lcp6 space
<i>calgary.zip</i>	1,043	2.00	0.69	0.75	0.84	1.50	9.98
<i>war&peace</i>	3,142	9.45	0.72	0.62	0.71	1.32	7.71
<i>texbook</i>	1,351	10.87	0.51	0.55	0.58	1.07	7.51
<i>bible</i>	3,952	13.97	0.71	0.59	0.67	1.22	7.31
<i>ecoli</i>	4,529	17.38	0.73	0.70	0.81	1.47	8.83
<i>world192</i>	2,415	23.01	0.60	0.57	0.64	1.19	6.99
<i>pic</i>	501	2,353.32	0.18	0.28	0.27	0.45	6.50

File	Size (Kb)	Ave. LCP	SA time	Lcp13 time	Lcp9 time	Lcp6 time	Lcp6 space
<i>etext99.gz</i>	38,747	2.65	0.97	1.07	1.18	2.08	9.97
<i>sprot</i>	107,048	89.08	1.49	1.00	1.03	1.90	7.01
<i>rfe</i>	113,693	93.02	1.18	0.89	0.92	1.66	6.86
<i>howto</i>	38,498	267.56	0.99	0.77	0.84	1.48	7.29
<i>reuters</i>	112,022	282.07	2.65	0.91	0.96	1.77	6.58
<i>linux</i>	113,530	479.00	1.04	0.76	0.76	1.35	6.88
<i>jdk13</i>	68,094	678.94	2.67	0.69	0.75	1.33	6.26
<i>etext99</i>	102,809	1,108.63	1.55	1.07	1.10	2.02	7.57
<i>chr22</i>	33,743	1,979.25	0.96	0.92	1.01	1.76	8.34
<i>gcc</i>	84,600	8,603.21	1.87	0.69	0.73	1.30	6.75
<i>w3c</i>	101,759	42,299.75	2.11	0.72	0.79	1.40	6.31

Table 2: Experimental results for small files (top) and large files (bottom). The second and third column show the file size and average LCP. The fourth column reports the time (microseconds per input byte) for the construction of the suffix array. The next three columns report the time (microseconds per input byte) for the computation of the LCP array using the algorithms Lcp13, Lcp9, and Lcp6 respectively. The last column shows the space used by Lcp6 expressed as total space occupancy over file length. The running times were measured considering (user + system) time averaged over five runs. The running times do not include the time spent for reading the input files. The test files are ordered by increasing average LCP.

6 Conclusions

In this paper we have addressed the problem of devising a “lightweight” algorithm for computing the LCP array given the text and the suffix array. We have considered the algorithm of Kasai et al. [10] and we have shown how to significantly reduce the space occupancy of this algorithm maintaining its simplicity and robustness.

Although we do not represent any data in compressed form, one of our space saving tricks makes use of the Burrows-Wheeler Transform. To our knowledge this is the first use of the Burrows-Wheeler Transform outside the domain of data compression. Our results show that we are still far from understanding all the potential applications of this fascinating mathematical tool.

References

- [1] M. I. Abouelhoda, E. Ohlebusch, and S. Kurtz. Optimal exact string matching based on suffix arrays. In *Proc. 9th International Symposium on String Processing and Information Retrieval (SPIRE '02)*, pages 31–43. Springer-Verlag LNCS n. 2476, 2002.
- [2] S. Burkhardt and J. Kärkkäinen. Fast lightweight suffix array construction and checking. In *Proc. 14th Symposium on Combinatorial Pattern Matching (CPM '03)*, pages 55–69. Springer-Verlag LNCS n. 2676, 2003.
- [3] M. Burrows and D. Wheeler. A block sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994.
- [4] P. Ferragina and G. Manzini. Opportunistic data structures with applications. In *Proc. of the 41st IEEE Symposium on Foundations of Computer Science*, pages 390–398, 2000.
- [5] P. Ferragina and G. Manzini. An experimental study of an opportunistic index. In *Proc. 12th ACM-SIAM Symposium on Discrete Algorithms*, pages 269–278, 2001.
- [6] P. Ferragina and G. Manzini. Compression boosting in optimal linear time using the Burrows-Wheeler transform. In *Proc. 15th ACM-SIAM Symposium on Discrete Algorithms (SODA '04)*, 2004.
- [7] R. Grossi, A. Gupta, and J. Vitter. High-order entropy-compressed text indexes. In *Proc. 14th Annual ACM-SIAM Symp. on Discrete Algorithms (SODA '03)*, pages 841–850, 2003.
- [8] W. Hon, K. Sadakane, and W. Sung. Breaking a time-and-space barrier in constructing full-text indices. In *Proc. of the 44th IEEE Symposium on Foundations of Computer Science*, pages 251–260, 2003.
- [9] J. Kärkkäinen and P. Sanders. Simple linear work suffix array construction. In *Proc. 30th International Colloquium on Automata, Languages and Programming (ICALP '03)*, pages 943–955. Springer-Verlag LNCS n. 2719, 2003.
- [10] T. Kasai, G. Lee, H. Arimura, S. Arikawa, and K. Park. Linear-time longest-common-prefix computation in suffix arrays and its applications. In *Proc. 12th Symposium on Combinatorial Pattern Matching (CPM '01)*, pages 181–192. Springer-Verlag LNCS n. 2089, 2001.

- [11] D. K. Kim, J. S. Sim, H. Park, and K. Park. Linear-time construction of suffix arrays. In *Proc. 14th Symposium on Combinatorial Pattern Matching (CPM '03)*, pages 186–199. Springer-Verlag LNCS n. 2676, 2003.
- [12] P. Ko and S. Aluru. Space efficient linear time construction of suffix arrays. In *Proc. 14th Symposium on Combinatorial Pattern Matching (CPM '03)*, pages 200–210. Springer-Verlag LNCS n. 2676, 2003.
- [13] N. J. Larsson and K. Sadakane. Faster suffix sorting. Technical Report LU-CS-TR:99-214, LUNDFD6/(NFCS-3140)/1-43/(1999), Department of Computer Science, Lund University, Sweden, 1999.
- [14] U. Manber and G. Myers. Suffix arrays: a new method for on-line string searches. *SIAM Journal on Computing*, 22(5):935–948, 1993.
- [15] G. Manzini and P. Ferragina. Lightweight suffix sorting home page. <http://www.mfn.unipmn.it/~manzini/lightweight>.
- [16] G. Manzini and P. Ferragina. Engineering a lightweight suffix array construction algorithm. In *Proc. of the 10th European Symposium on Algorithms (ESA '02)*, pages 698–710. Springer Verlag LNCS n. 2461, 2002.
- [17] E. M. McCreight. A space-economical suffix tree construction algorithm. *Journal of the ACM*, 23(2):262–272, 1976.
- [18] K. Sadakane. Succinct representations of LCP information and improvements in the compressed suffix arrays. In *Proc. 13th Annual ACM-SIAM Symp. on Discrete Algorithms (SODA '02)*, pages 225–232, 2002.