

# Opportunistic Data Structures with Applications

Paolo Ferragina\*  
Università di Pisa

Giovanni Manzini†  
Università del Piemonte Orientale

## Abstract

*In this paper we address the issue of compressing and indexing data. We devise a data structure whose space occupancy is a function of the entropy of the underlying data set. We call the data structure opportunistic since its space occupancy is decreased when the input is compressible and this space reduction is achieved at no significant slowdown in the query performance. More precisely, its space occupancy is optimal in an information-content sense because a text  $T[1, u]$  is stored using  $O(H_k(T)) + o(1)$  bits per input symbol in the worst case, where  $H_k(T)$  is the  $k$ th order empirical entropy of  $T$  (the bound holds for any fixed  $k$ ). Given an arbitrary string  $P[1, p]$ , the opportunistic data structure allows to search for the occurrences of  $P$  in  $T$  in  $O(p + occ \log^\epsilon u)$  time (for any fixed  $\epsilon > 0$ ). If data are uncompressible we achieve the best space bound currently known [12]; on compressible data our solution improves the succinct suffix array of [12] and the classical suffix tree and suffix array data structures either in space or in query time or both.*

*We also study our opportunistic data structure in a dynamic setting and devise a variant achieving effective search and update time bounds. Finally, we show how to plug our opportunistic data structure into the Glimpse tool [19]. The result is an indexing tool which achieves sublinear space and sublinear query time complexity.*

## 1 Introduction

Data structure is a central concept in algorithmics and computer science in general. In the last decades it has been investigated from different points of view and its basic ideas enriched by new functionalities with the aim to cope with the features of the peculiar setting of use: dynamic, persistent, self-adjusting, implicit, fault-tolerant, just to cite a few.

\* Dipartimento di Informatica, Università di Pisa, 56100 Pisa, Italy. E-mail: [ferragin@di.unipi.it](mailto:ferragin@di.unipi.it). Supported in part by Italian MURST project “Algorithms for Large Data Sets: Science and Engineering” and by UNESCO grant UVO-ROSTE 875.631.9.

† Dipartimento di Scienze e Tecnologie Avanzate, Università del Piemonte Orientale, 15100 Alessandria, Italy and IMC-CNR, 56100 Pisa, Italy. E-mail: [manzini@mfn.unipmn.it](mailto:manzini@mfn.unipmn.it). Supported in part by MURST 60% funds.

Space reduction in data structural design is an attractive issue, now more than ever before, because of the exponential increase of electronic data nowadays available, and because of its intimate relation with algorithmic performance improvements (see e.g. Knuth [16] and Bentley [5]). This has recently motivated an upsurging interest in the design of *implicit* data structures for basic searching problems (see [23] and references therein). The goal is to reduce as much as possible the *auxiliary information* kept together with the input data without introducing any significant slowdown in the query performance. However, input data are represented in their entirety thus taking no advantage of possible repetitiveness into them. The importance of those issues is well known to programmers who typically use various tricks to squeeze data as much as possible and still achieve good query performance. Their approaches, though, boil down to heuristics whose effectiveness is witnessed only by experimentation.

In this paper we address the issue of compressing and indexing data by studying it in a theoretical framework. From the best of our knowledge no other result is known in the literature about the study of the interplay between compression and indexing of data collections. The exploitation of data compressibility have been already investigated only with respect to its impact on algorithmic performance in the context of on-line algorithms (e.g. caching and prefetching [15, 17]), string-matching algorithms (see e.g. [1, 2, 9]), sorting and computational geometry algorithms [8].

**The scenario.** Most of the research in the design of indexing data structures has been directed to devise solutions which offer a good trade-off between query and update time versus space usage. The two main approaches are *word-based* indices and *full-text* indices. The former achieve succinct space occupancy at the cost of being mainly limited to index linguistic texts [27], the latter achieve versatility and guaranteed performance at the cost of requiring large space occupancy (see e.g. [10, 18, 21]). Some progress on full-text indices has been recently achieved [12, 23], but an asymptotical linear space *seems* unavoidable and this makes word-based indices much more appealing when space occupancy is a primary concern. In this context compression appears always as an attractive choice, if not mandatory. Processing speed is currently improving at a faster rate than disk speed. Since compression decreases the demand of

storage at the expenses of processing, it is becoming more economical to store data in a compressed form rather than uncompressed.

Starting from these promising considerations, many researchers have recently concentrated on the *compressed matching problem*, introduced in [1], as the task of performing string matching in a compressed text without decompressing it. A collection of algorithms is currently known to solve efficiently (possibly optimally) this problem on text compressed by means of various schemes: e.g. run-length [1], LZ77 [9], LZ78 [2], Huffman [24]. All of these results, although asymptotically faster than the classical scan-based methods, they rely on the scan of the whole *compressed* text and thus result still unacceptable for large text collections.

Approaches to combine compression and indexing techniques are nowadays receiving more and more attention, especially in the context of word-based indices, achieving *experimental* trade-offs between space occupancy and query performance (see e.g. [4, 19, 27]). An interesting idea towards the direct compression of the index data structure has been proposed in [13, 14] where the properties of the Lempel-Ziv's compression scheme have been exploited to reduce the number of *index points*, still supporting pattern searches. As a result, the overall index requires provably sublinear space but at the cost of either limiting the search to  $q$ -grams [13] or worsening significantly the query performance [14].

A natural question arises at this point: Do full-text indices need a space occupancy *linear* in the (uncompressed) text size in order to support effective search operations on arbitrary patterns? It is a common belief [27] that some space overhead must be paid to use the full-text indices, but is this actually a provable need?

**Our Results.** In this paper we answer the two questions above by providing a novel data structure for indexing and searching whose space occupancy is a *function of the entropy* of the underlying data set. The data structure is called *opportunistic* in that, although no assumption on a particular distribution is made, it takes advantage of the compressibility of the input data by decreasing the space occupancy at *no significant slowdown* in the query performance.<sup>1</sup> The data structure is provably space optimal in an *information-content* sense because it stores a text  $T[1, u]$  using  $O(H_k(T)) + o(1)$  bits per input symbol in the worst case (for any fixed  $k \geq 0$ ), where  $H_k(T)$  is the  $k$ th order empirical entropy.  $H_k$  expresses the maximum compression we can achieve using for each character a code which depends only on the  $k$  characters preceding it. We point out that in the case of an uncompressible string  $T$ , the space occupancy is  $\Theta(u)$  bits which is actually optimal [12]; for

<sup>1</sup>The concept of *opportunistic algorithm* has been introduced in [9] to characterize an algorithm which takes advantage of the compressibility of the text to speed up its (scan based) search operations. In our paper, we turn this concept into the one of *opportunistic data structure*.

a compressible string, our opportunistic data structure is the first to achieve sublinear space occupancy. Given an arbitrary pattern  $P[1, p]$ , such an opportunistic data structure allows to search for the *occ* occurrences of  $P$  in  $T$  in  $O(p + occ \log^\epsilon u)$  time, for any fixed  $\epsilon > 0$ .

The novelty of our approach resides in the careful combination of the Burrows-Wheeler compression algorithm [7] with the the suffix array data structure [18] to obtain a sort of *compressed* suffix array. We indeed show how to *augment* the information kept by the Burrows-Wheeler algorithm, in order to support effective *random accesses* to the compressed data without the need of uncompressing all of them at query time. We design two algorithms for operating on our opportunistic data structure. The first algorithm is an effective approach to search for an arbitrary pattern  $P[1, p]$  in a *compressed* suffix array, taking  $O(p)$  time in the worst case (Section 3.1). The second algorithm exploits compression to speed up the retrieval of the actual positions of the pattern occurrences, thus incurring only in a *sublogarithmic*  $O(\log^\epsilon u)$  time slowdown for any fixed  $\epsilon > 0$  (Section 3.2).

In some sense, our result can be interpreted as a method to *compress* the suffix array, and still support effective searches for arbitrary patterns. In their seminal paper, Manber and Myers [18] introduced the suffix array data structure showing how to search for a pattern  $P[1, p]$  in  $O(p + \log u + occ)$  time in the worst case. The suffix array uses  $\Theta(u \log u)$  bits of storage. Recently, Grossi and Vitter [12] reduced the space usage of suffix arrays to  $\Theta(u)$  bits at the cost of requiring  $O(\log^\epsilon u)$  time to retrieve the  $i$ -th suffix. Hence, searching in this succinct suffix array via the classical Manber-Myers' procedure takes  $O(p + \log^{1+\epsilon} u + occ \log^\epsilon u)$  time. Our solution therefore improves the succinct suffix array of [12] both in space and query time complexity. The authors of [12] introduce also other hybrid indices which achieve better query-time complexity but yet require  $\Omega(u)$  bits of storage. As far as the problem of counting the pattern occurrences is concerned, our solution improves the classical suffix tree and suffix array data structures, because they achieve  $\Omega(p)$  time complexity and occupy  $\Omega(u \log u)$  bits of storage.

In Section 4, we investigate the modifiability of our opportunistic data structure by studying how to choreograph its basic ideas with a dynamic setting. We show that a dynamic text collection  $\Delta$  of size  $u$  can be stored in  $O(H_k(\Delta)) + o(1)$  bit per input symbol (for any fixed  $k \geq 0$  and not very short texts), support insert operations on individual texts  $T[1, t]$  in  $O(t \log u)$  amortized time, delete operations on  $T[1, t]$  in  $O(t \log^2 u)$  amortized time, and search for a pattern  $P[1, p]$  in  $O(p \log^3 u + occ \log u)$  time in the worst case. We point out that even in the case of an uncompressible text  $T$ , our space bounds are the best known ones since the data structures in [12] do not support updates (the dynamic case is left as open in their Section 4).

Finally, we investigate applications of our ideas to the development of novel text retrieval systems based on the

concept of block addressing (first introduced in the Glimpse tool [19]). The notable feature of block addressing is that it can achieve both sublinear space overhead and sublinear query time, whereas inverted indices achieve only the second goal [4]. Unfortunately, up to now all the known block addressing indices [4, 19] achieve time and space sublinearity only under some restrictive conditions on the block size. We show how to use our opportunistic data structure to devise a novel block addressing scheme, called *CGlimpse* (standing for *Compressed Glimpse*), which always achieves time and space sublinearity.

## 2 Background

Let  $T[1, u]$  be a text drawn from a constant-size alphabet  $\Sigma$ . A central concept in our discussion is the *suffix array* data structure [18]. The suffix array  $\mathcal{A}$  built on  $T[1, u]$  is an array containing the lexicographically ordered sequence of the suffixes of  $T$ , represented via pointers to their starting positions (i.e., *integers*). For instance, if  $T = ababc$  then  $\mathcal{A} = [1, 3, 2, 4, 5]$ . Clearly  $\mathcal{A}$  requires  $u \log_2 u$  bits, actually a lot when indexing large text collections. It is a long standing belief that suffix arrays are uncompressible because of the “apparently random” permutation of the suffix pointers. Recent results in the data compression field have opened the door to revolutionary ways to compress suffix arrays and are basic tools of our data structure.

In [7], Burrows and Wheeler propose a transformation (BWT from now on) consisting of a reversible permutation of the text characters which gives a new string that is “easier to compress”. The BWT tends to group together characters which occur adjacent to similar text substrings. This nice property is exploited by locally-adaptive compression algorithms, such as move-to-front coding [6], in combination with statistical (i.e. Huffman or Arithmetic coders) or structured coding models. The BWT-based compressors are among the best compressors currently available since they achieve a very good compression ratio using relatively small time and space.

**The reversible BWT.** We distinguish between a *forward* transformation, which produces the string to be compressed, and a *backward* transformation which gives back the original text from the transformed one. The forward BWT consists of three basic steps: (1) Append to the end of  $T$  a special character  $\#$  smaller than any other text character; (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  $L$  by taking the last column of  $\mathcal{M}$ . Notice that every column of  $\mathcal{M}$  is a permutation of the last column  $L$ , and in particular the first column of  $\mathcal{M}$ , call it  $F$ , is obtained by lexicographically sorting the characters in  $L$ .

There is a strong relation between the matrix  $\mathcal{M}$  and the suffix array  $\mathcal{A}$  of the string  $T$ . When sorting the rows of the matrix  $\mathcal{M}$  we are essentially sorting the suffixes of  $T$ .

Consequently, entry  $\mathcal{A}[i]$  points to the suffix of  $T$  occupying (a prefix of) the  $i$ th row of  $\mathcal{M}$ . The cost of performing the forward BWT is given by the cost of constructing the suffix array  $\mathcal{A}$ , and this requires  $O(u)$  time [21].

The cyclic shift of the rows of  $\mathcal{M}$  is crucial to define the backward BWT, which is based on two easy to prove observations [7]:

- a. Given the  $i$ th row of  $\mathcal{M}$ , its last character  $L[i]$  precedes its first character  $F[i]$  in the original text  $T$ , namely  $T = \dots L[i]F[i] \dots$ .
- b. Let  $L[i] = c$  and let  $r_i$  be the number of occurrences of  $c$  in the prefix  $L[1, i]$ . Let  $\mathcal{M}[j]$  be the  $r_i$ -th row of  $\mathcal{M}$  starting with  $c$ . The character in the first column  $F$  corresponding to  $L[i]$  is located at  $F[j]$ . We call this *LF-mapping* (Last-to-First mapping) and set  $LF[i] = j$ .

We are now ready to describe the backward BWT:

1. Compute the array  $C[1 \dots |\Sigma|]$  storing in  $C[c]$  the number of occurrences in  $T$  of the characters  $\{\#, 1, \dots, c-1\}$ . Notice that  $C[c] + 1$  is the position of the first occurrence of  $c$  in  $F$  (if any).
2. Define the LF-mapping  $LF[1 \dots u + 1]$  as follows:  $LF[i] = C[L[i]] + r_i$ , where  $r_i$  equals the number of occurrences of the character  $L[i]$  in the prefix  $L[1, i]$  (see observation (b) above).
3. Reconstruct  $T$  backward as follows: set  $s = 1$  and  $T[u] = L[1]$  (because  $\mathcal{M}[1] = \#T$ ); then, for each  $i = u - 1, \dots, 1$  do  $s = LF[s]$  and  $T[i] = L[s]$ .

In [26] it is shown how to derive the suffix array  $\mathcal{A}$  from  $L$  in linear time; however, in the context of pattern searching, the algorithm in [26] is no better than the known scan-based opportunistic algorithms (such as [9]). Nonetheless, the *implicit* presence of the suffix array  $\mathcal{A}$  into  $L$  suggests to take full advantage of the structure of  $\mathcal{A}$  for fast searching, and of the high compressibility of  $L$  for space reduction. This is actually the ultimate hope of any indexer: succinct and fast! In the next section, we show that this result is achievable provided that a sublogarithmic slowdown (wrt the suffix array) is introduced in the cost of listing the pattern occurrences.

Let  $T^{bw} = \text{bwt}(T)$  denote the last column  $L$ , output of the BWT. Our indexing data structure consists of a compressed version of  $T^{bw}$  together with some other auxiliary array-based data structures that support random access to  $T^{bw}$ . We compress  $T^{bw}$  in three steps (see also [20]):

1. Use a move-to-front coder, briefly *mtf* [6], to encode a character  $c$  via the count of distinct characters seen since its previous occurrence. The structural properties of  $T^{bw}$ , mentioned above, imply that the string  $T^{mtf} = \text{mtf}(T^{bw})$  will be dominated by *low numbers*.

2. Encode each run of zeroes in  $T^{mtf}$  using run length encoding (rle). More precisely, replace the sequence  $0^m$  with the number  $(m + 1)$  written in binary, least significant bit first, discarding the most significant bit. For this encoding we use two new symbols **0** and **1** so that the resulting string  $T^{rl} = rle(T^{mtf})$  is over the alphabet  $\{0, 1, 1, 2, \dots, |\Sigma| - 1\}$ .
3. Compress  $T^{rl}$  by means of a variable-length prefix code, called PC, which encodes the symbols **0** and **1** using two bits (10 for **0**, 11 for **1**), and the symbol  $i$  using a variable-length prefix code of  $1 + 2 \lfloor \log(i + 1) \rfloor$  bits, the first one being a zero.

The resulting algorithm  $BW\_RLX = bwt + mtf + rle + PC$  is sufficiently simple so that in the rest of the paper we can concentrate on the searching algorithm without being distracted by the details of the compression. Despite of the simplicity of  $BW\_RLX$ , using the results in [20]<sup>2</sup> it is possible to show that (proof in the full paper), for any  $k \geq 0$  and for any  $T$  there exists a constant  $g_k$  such that

$$|BW\_RLX(T)| \leq 5 |T| H_k(T) + g_k \log |T| \quad (1)$$

where  $H_k$  is the  $k$ th order empirical entropy.  $H_k$  expresses the maximum compression we can achieve using for each character a code which depends only on the  $k$  characters preceding it.

### 3 Searching in BWT-compressed text

Let  $T[1, u]$  denote an arbitrary text over the alphabet  $\Sigma$ , and let  $Z = BW\_RLX(T)$ . In this section we describe an algorithm which, given a pattern  $P[1, p]$ , reports all occurrences of  $P$  in the uncompressed text  $T$  by looking only at  $Z$  and without uncompressing all of it. Our algorithm makes use of the relation between the suffix array  $\mathcal{A}$  and the matrix  $\mathcal{M}$ . Recall that the suffix array  $\mathcal{A}$  posses two nice structural properties which are usually exploited to support fast pattern searches: (i) all the suffixes of the text  $T$  prefixed by a pattern  $P$  occupy a contiguous portion (subarray) of  $\mathcal{A}$ ; (ii) that subarray has starting position  $sp$  and ending position  $ep$ , where  $sp$  is actually the *lexicographic position* of the string  $P$  among the ordered sequence of text suffixes.

#### 3.1 Step I: Counting the occurrences

We now describe an algorithm, called  $BW\_Count$ , which identifies the positions  $sp$  and  $ep$  by accessing only the compressed string  $Z$  and some auxiliary *array-based* data structures.

$BW\_Count$  consists of  $p$  phases each one preserving the following invariant: *At the  $i$ -th phase, the parameter  $sp$*

<sup>2</sup>The algorithm  $BW\_RLX$  corresponds to the procedure  $A^*$  described in [20]

---

#### Algorithm $BW\_Count(P[1, p])$

1.  $c = P[p], i = p;$
  2.  $sp = C[c] + 1, ep = C[c + 1];$
  3. **while**  $((sp \leq ep)$  **and**  $(i \geq 2))$  **do**
  4.      $c = P[i - 1];$
  5.      $sp = C[c] + Occ(c, 1, sp - 1) + 1;$
  6.      $ep = C[c] + Occ(c, 1, ep);$
  7.      $i = i - 1;$
  8. **if**  $(ep < sp)$  **then return** “pattern not found”  
**else return** “found  $(ep - sp + 1)$  occurrences”
- 

**Figure 1. Algorithm for counting the number of occurrences of  $P[1, p]$  in  $T[1, u]$ .**

points to the first row of  $\mathcal{M}$  prefixed by  $P[i, p]$  and the parameter  $ep$  points to the last row of  $\mathcal{M}$  prefixed by  $P[i, p]$ . The pseudo-code is given in Fig. 1. In the first phase (i.e.  $i = p$ ),  $sp$  and  $ep$  are determined via the array  $C$  defined in Section 2 (Step 2). The values  $sp$  and  $ep$  are updated at Steps 5 and 6 using the subroutine  $Occ(c, 1, k)$  which reports the number of occurrences of  $c$  in  $T^{bw}[1, k]$ . Note that at Steps 5 and 6 we are computing the LF-mapping for, respectively, the first and the last occurrence (if any) of  $P[i - 1]$  in  $T^{bw}[sp, ep]$ . If at the generic  $i$ th phase we have  $ep < sp$  we can conclude that  $P[i, p]$  does not occur in  $T$  and hence  $P$  does not too. After the final phase,  $sp$  and  $ep$  will delimit the portion of  $\mathcal{M}$  (and thus of the suffix array  $\mathcal{A}$ ) containing all the text suffixes prefixed by  $P$ . The integer  $(ep - sp + 1)$  will therefore account for the number of occurrences of  $P$  in  $T$ . The following lemma proves the correctness of  $BW\_Count$  assuming  $Occ$  works as claimed (proof in the full paper).

**Lemma 1** *For  $i = p, p - 1, \dots, 2$ , if  $P[i - 1, p]$  occurs in  $T$  then Step 5 (resp. Step 6) of  $BW\_Count$  correctly updates the value of  $sp$  (resp.  $ep$ ) thus pointing to the first (resp. last) row prefixed by  $P[i - 1, p]$ . ■*

The running time of  $BW\_Count$  depends on the cost of the procedure  $Occ$ . We now describe an algorithm for computing  $Occ(c, 1, k)$  in  $O(1)$  time, on a RAM with word size  $\Theta(\log u)$  bits.

We logically partition the transformed string  $T^{bw}$  into substrings of  $\ell$  characters each (called *buckets*), and denote them by  $BT_i = T^{bw}[(i - 1)\ell + 1, i\ell]$ , for  $i = 1, \dots, u/\ell$ . This partition naturally induces a partition of  $T^{mtf}$  into  $u/\ell$  buckets  $BT_1^{mtf}, \dots, BT_{u/\ell}^{mtf}$  of size  $\ell$  too. We assume that each run of zeroes in  $T^{mtf}$  is entirely contained in a single bucket and we describe our algorithm

for computing  $\text{Occ}(c, 1, k)$  under this simplifying assumption. The general case in which a sequence of zeroes may span several buckets is similar and thus its discussion is deferred to the full paper. Under our assumption, the buckets  $BT_i^{mtf}$ 's induce a partition of the compressed file  $Z$  into  $u/\ell$  compressed buckets  $BZ_1, \dots, BZ_{u/\ell}$ , defined as  $BZ_i = \text{PC}(\text{r1e}(BT_i^{mtf}))$ .

Let  $BT_i$  denote the bucket containing the character  $T^{bw}[k]$  (namely  $i = \lceil k/\ell \rceil$ ). The computation of  $\text{Occ}(c, 1, k)$  is based on a hierarchical decomposition of  $T^{bw}[1, k]$  in three substrings as follows: (i) the longest prefix of  $T^{bw}[1, k]$  having length a multiple of  $\ell^2$  (i.e.  $BT_1 \cdots BT_{i^*}$ , where  $i^* = \lfloor \frac{k-1}{\ell^2} \rfloor$ ), (ii) the longest prefix of the remaining suffix having length a multiple of  $\ell$  (i.e.  $BT_{i^*+1} \cdots BT_{i-1}$ ), and finally (iii) the remaining suffix of  $T^{bw}[1, k]$  which is indeed a prefix of the bucket  $BT_i$ . We compute  $\text{Occ}(c, 1, k)$  by summing the number of occurrences of  $c$  in each of these substrings. This can be done in  $O(1)$  time and sublinear space using the following auxiliary data structures.

For the calculations on the substring of point (i):

- For  $i = 1, \dots, u/\ell^2$ , the array  $NO_i[1, |\Sigma|]$  stores in the entry  $NO_i[c]$  the number of occurrences of the character  $c$  in  $BT_1 \cdots BT_{i\ell}$ .
- The array  $W[1, u/\ell^2]$  stores in the entry  $W[i]$  the value  $\sum_{j=1}^{i\ell} |BZ_j|$  equals to the sum of the sizes of the compressed buckets  $BZ_1, \dots, BZ_{i\ell}$ .

For the calculations on the substring of point (ii):

- For  $i = 1, \dots, u/\ell$ , the array  $NO'_i[1, |\Sigma|]$  stores in the entry  $NO'_i[c]$  the number of occurrences of the character  $c$  in the string  $BT_{i^*+1} \cdots BT_{i-1}$  (this concatenated string has length less than  $\ell^2$ ).
- The array  $W'[1, u/\ell]$  stores in the entry  $W'[i]$  the value  $\sum_{j=i^*+1}^{i-1} |BZ_j|$  equals to the overall size of the compressed buckets  $BZ_{i^*+1}, \dots, BZ_{i-1}$  (the value is bounded above by  $O(\ell^2)$ ).

For the calculations on the (compressed) buckets:

- The array  $MTF[1, u/\ell]$  stores in the entry  $MTF[i]$  a picture of the state of the MTF list at the beginning of the encoding of  $BT_i$ . Each entry takes  $|\Sigma| \log |\Sigma|$  bits (i.e.  $O(1)$  bits).
- The table  $S$  stores in the entry  $S[c, j, b, m]$  the number of occurrences of  $c$  among the first  $j$  characters of the compressed string  $b$ , assuming that  $m$  is the picture of the MTF list used to produce  $b$ . Thus, entry  $S[c, j, BZ_i, MTF[i]]$  stores the number of occurrences of  $c$  in  $BT_i[1, j]$ . Table  $S$  has  $O(\ell 2^{\ell'})$  entries each one occupying  $O(\log \ell)$  bits, where  $\ell'$  is the maximum length of a compressed bucket.

The computation of  $\text{Occ}(c, 1, k)$  therefore proceeds as follows. First, the bucket  $BT_i$  containing the character  $c =$

$T^{bw}[k]$  is determined via  $i = \lceil k/\ell \rceil$ , together with the position  $j = k - (i-1)\ell$  of this character in  $BT_i$  and the parameter  $i^* = \lfloor (k-1)/\ell^2 \rfloor$ . Then the number of occurrences of  $c$  in the prefix  $BT_1 \cdots BT_{i^*}$  (point (i)) is determined via  $NO_{i^*}[c]$ , and the number of occurrences of  $c$  in the substring  $BT_{i^*+1}, \dots, BT_{i-1}$  (point (ii)) is determined via  $NO'_i[c]$ . Finally, the compressed bucket  $BZ_i$  is retrieved from  $Z$  (notice that  $W[i^*] + W'[i] + 1$  is its starting position), and the number of occurrences of  $c$  within  $BT_i[1, j]$  are accounted accessing  $S[c, j, BZ_i, MTF[i]]$  in  $O(1)$  time. The sum of these three quantities gives  $\text{Occ}(c, 1, k)$ .

By construction any compressed bucket  $BZ_i$  has size at most  $\ell' = (1 + 2 \lfloor \log \Sigma \rfloor) \ell$  bits. We choose  $\ell = \Theta(\log u)$  so that  $\ell' = c \log u$  with  $c < 1$ . Under this assumption, every step of  $\text{Occ}$  consists of arithmetic operations or table lookup operations involving  $O(\log u)$ -bit operands. Consequently every call to  $\text{Occ}$  takes  $O(1)$  time on a RAM. As far as the space occupancy is concerned, the arrays  $NO$  and  $W$  take  $O((u/\ell^2) \log u) = O(u/\log u)$  bits. The arrays  $NO'$  and  $W'$  take  $O((u/\ell) \log \ell) = O((u/\log u) \log \log u)$  bits. The array  $MTF$  takes  $O(u/\ell) = O(u/\log u)$  bits. Table  $S$  consists of  $O(\ell 2^{\ell'}) \log \ell$ -bit entries and thus it occupies  $O(2^{\ell'} \ell \log \ell) = O(u^c \log u \log \log u)$  bits, where  $c < 1$ . We conclude that the auxiliary data structures used by  $\text{Occ}$  occupy  $O((u/\log u) \log \log u)$  bits (in addition to the compressed file  $Z$ ).

**Theorem 1** *Let  $Z$  denote the output of the algorithm BW\_RLX on input  $T[1, u]$ . The number of occurrences of a pattern  $P[1, p]$  in  $T[1, u]$  can be computed in  $O(p)$  time on a RAM. The space occupancy is  $|Z| + O\left(\frac{u}{\log u} \log \log u\right)$  bits in the worst case. ■*

### 3.2 Step II: Locating the occurrences

We now consider the problem of determining the exact position in the text  $T$  of all the occurrences of the pattern  $P[1, p]$ . This means that for  $s = sp, sp+1, \dots, ep$ , we want to find the text position  $pos(s)$  of the suffix which prefixes the  $s$ th row  $\mathcal{M}[s]$ . We propose two approaches: the first one is simple and slow, the second one is faster and relies on the very special properties of the the string  $T^{bw}$ .

In the first algorithm we *logically mark* the rows of  $\mathcal{M}$  which correspond to text positions having the form  $1 + i\eta$ , for  $\eta = \Theta(\log^2 u)$  and  $i = 0, 1, \dots, u/\eta$ . We store with these marked rows the starting positions of the corresponding text suffixes explicitly. This preprocessing is done at compression time. At query time we find  $pos(s)$  as follows. If  $s$  is a marked row, then there is nothing to be done and its position is directly available. Otherwise, we use the LF-mapping to find the row  $s'$  corresponding to the suffix  $T[pos(s) - 1, u]$ . We iterate this procedure  $v$  times until  $s'$  points to a marked row; at that point  $pos(s')$  is available and we set  $pos(s) = pos(s') + v$ . The crucial point of

the algorithm is the logical marking of the rows of  $\mathcal{M}$  corresponding to the text suffixes starting at positions  $1 + i\eta$ ,  $i = 0, \dots, u/\eta$ . Our solution consists in storing the row numbers in a two-level bucketing scheme. We partition the rows of  $\mathcal{M}$  into buckets of size  $\Theta(\log^2 u)$  each. For each bucket, we take all the marked rows lying into it, and store them into a Packet B-tree [3] using as a key their distance from the beginning of the bucket. Since a bucket contains at most  $O(\log^2 u)$  keys, each  $O(\log \log u)$  bits long, membership queries take  $O(1)$  time on a RAM. The overall space required for the logical marking is  $O((u/\eta) \log \log u)$  bits. In addition, for each marked row we also keep the starting position of the corresponding text suffix (i.e.  $pos()$ ), which requires additional  $O(\log u)$  bits per marked row. Consequently, the overall space occupancy is  $O((u/\eta) \log u) = O(u/\log u)$  bits. For what concerns the time complexity, our algorithm computes  $pos(s)$  in at most  $\eta = \Theta(\log^2 u)$  steps, each taking constant time. Hence the  $occ$  occurrences of a pattern  $P$  in  $T$  can be retrieved in  $O(occ \log^2 u)$  time, with a space overhead of  $O(u/\log u)$  bits. Combining the results of this section with (1) we have:

**Theorem 2** *A text  $T[1, u]$  can be preprocessed in  $O(u)$  time so that all the  $occ$  occurrences of a pattern  $P[1, p]$  in  $T$  can be listed in  $O(p + occ \log^2 u)$  time on a RAM. The space occupancy is bounded by  $5H_k(T) + O(\frac{\log \log u}{\log u})$  bits per input symbol in the worst case, for any fixed  $k \geq 0$ . ■*

We now refine the above algorithm in order to compute  $pos(s)$  in  $O(\log^\epsilon u)$  time for any fixed  $\epsilon > 0$ . We still use the idea of marking some of the rows in  $\mathcal{M}$ , however we introduce some *shortcuts* which allow to move in  $T$  by more than one character at a time, thus reducing the number of steps required to reach a marked position. The key ingredient of our new approach is a procedure for computing the LF-mapping over a string  $\tilde{T}$  drawn from an alphabet  $\Lambda$  of *non-constant* size (proof and details in the full paper):

**Lemma 2** *Given a string  $\tilde{T}[1, v]$  over an arbitrary alphabet  $\Lambda$ , we can compute the LF-mapping over  $\tilde{T}^{bw}$  in  $O(\log^\epsilon v)$  time using  $O(v(1 + H_k(\tilde{T})) + |\Lambda|^{k+1}(\log |\Lambda| + \log v))$  bits of storage, for any given  $\epsilon > 0$ . ■*

We use Lemma 2 to compute  $pos(s)$  in  $O(\log^{(1/2)+2\epsilon} u)$  time; this is an intermediate result that will be then refined to achieve the final  $O(\log^\epsilon u)$  time bound.

At compression time we logically mark the rows of  $\mathcal{M}$  which correspond to text positions of the form  $1 + i\gamma$  for  $i = 0, \dots, u/\gamma$  and  $\gamma = \Theta(\log^{(1/2)+\epsilon} u)$ . Then, we consider the string  $T_0$  obtained by grouping the characters of  $T$  into blocks of size  $\gamma$ . Clearly  $T_0$  has length  $u/\gamma$  and its characters belong to the alphabet  $\Sigma^\gamma$ . Let  $\mathcal{M}_0$  denote the cyclic-shift matrix associated to  $T_0$ ; notice that  $\mathcal{M}_0$  consists of the marked rows of  $\mathcal{M}$ . Now we mark the rows of  $\mathcal{M}_0$  corresponding to the suffixes of  $T_0$  starting at the positions  $1 + i\eta$ , for  $i = 0, \dots, |T_0|/\eta$  and  $\eta = \Theta(\log^{(1/2)+\epsilon} u)$ . For these

rows we explicitly keep the starting position of the corresponding text suffixes. To compute  $pos(s)$  we first compute the LF-mapping in  $\mathcal{M}$  until we reach a marked row  $s'$ . Then we compute  $pos(s')$  by finding its corresponding row in  $\mathcal{M}_0$  and computing the LF-mapping in  $\mathcal{M}_0$  (via Lemma 2) until we reach a marked row  $s''$  in  $\mathcal{M}_0$  (for which  $pos(s'')$  is explicitly available by construction). The marking of  $T$  and the counting of the number of marked rows in  $\mathcal{M}$  that precede a given marked row  $s'$  (this is required in order to determine the position in  $\mathcal{M}_0$  of  $\mathcal{M}[s']$ ) can be done in constant time and  $O(\frac{u}{\gamma} \log \log u)$  bits of storage using again a Packed B-tree and a two level bucketing scheme as before. In addition, for  $\Theta(|T_0|/\eta)$  rows of  $\mathcal{M}_0$  we keep explicitly their positions in  $T_0$  which take  $\Theta((|T_0|/\eta) \log u) = \Theta(u/\log^{2\epsilon} u)$  bits of storage. The space occupancy of the procedure for computing the LF-mapping in  $T_0^{bw}$  is given by Lemma 2. Since  $H_k(T_0) \leq \gamma H_{k\gamma}(T)$ , a simple algebraic calculation yields that the overall space occupancy is  $O(H_k(T) + \frac{1}{\log^{2\epsilon} u})$  bits per input symbol, for any fixed  $k$ . The time complexity of the algorithm is  $O(\gamma)$  (for finding a marked row in  $\mathcal{M}$ ) plus  $O(\eta \log^\epsilon u)$  (for finding a marked row in  $\mathcal{M}_0$ ), thus  $O(\log^{(1/2)+2\epsilon} u)$  time overall.

The final time bound of  $O(\log^\epsilon u)$  for the computation of  $pos(s)$  can be achieved by iterating the approach above as follows. The main idea is to take  $\gamma_0 = \Theta(\log^\epsilon u)$ , and apply the procedure for computing the LF-mapping in  $T_0$  for  $O(\log^\epsilon u)$  steps, thus identifying a row  $s_1$  of the matrix  $\mathcal{M}_0$  such that  $pos(s_1)$  has the form  $1 + i\gamma_1$  with  $\gamma_1 = \Theta(\log^{2\epsilon} u)$ . Next, we define the string  $T_1$  obtained by grouping the characters of  $T$  into blocks of size  $\gamma_1$  and we consider the corresponding matrix  $\mathcal{M}_1$ . By construction  $s_1$  corresponds to a row in  $\mathcal{M}_1$  and we can iterate the above scheme. At the  $j$ th step we operate on the matrix  $\mathcal{M}_{j-1}$  until we find a row  $s_j$  such that  $pos(s_j)$  has the form  $1 + i\gamma_j$  where  $\gamma_j = \Theta(\log^{(j+1)\epsilon} u)$ . This continues until  $j$  reaches the value  $\lceil 1/\epsilon \rceil$ . At that point the matrix  $\mathcal{M}_j$  consists of  $\Theta(u/\log^{1+\delta} u)$  rows, where  $\delta = \lceil 1/\epsilon \rceil \epsilon - 1$ . Since we can always choose  $\epsilon$  so that  $\delta > 0$ , we can store explicitly the starting positions  $pos()$  of the marked text suffixes in  $\mathcal{M}_j$  using sublinear space, i.e.  $o(u)$  bits. Summing up, the algorithm computes  $pos(s)$  in  $\lceil 1/\epsilon \rceil = \Theta(1)$  iterations, each taking  $\Theta(\log^{2\epsilon} u)$  time. Since  $\epsilon$  is an arbitrary positive constant, it is clear that we can rewrite the previous time bound as  $\Theta(\frac{\log^\epsilon u}{\epsilon}) = \Theta(\log^\epsilon u)$ . The space occupancy is dominated by the one required for the marking of  $\mathcal{M}$ .

**Theorem 3** *A text  $T[1, u]$  can be indexed so that all the  $occ$  occurrences of a pattern  $P[1, p]$  in  $T$  can be listed in  $O(p + occ \log^\epsilon u)$  time on a RAM. The space occupancy is  $O(H_k(T) + \frac{\log \log u}{\log^\epsilon u})$  bits per input symbol in the worst case, for any fixed  $k \geq 0$ . ■*

## 4 Dynamizing our approach

Let  $\Delta$  be a dynamic collection of texts  $\{T_1, \dots, T_m\}$  having arbitrary lengths and total size  $u$ . Collection  $\Delta$  may shrink or grow over the time due to **insert** and **delete** operations which allow to add or remove from  $\Delta$  an individual text string. Our aim is to store  $\Delta$  in succinct space, perform the update operations efficiently, and support fast searches for the occurrences of an arbitrary pattern  $P[1, p]$  into  $\Delta$ 's texts. This problem can be solved in optimal time complexity and  $\Theta(u \log u)$  bits of storage [10, 21]. In the present section we aim at *dynamizing* our compressed index in order to keep  $\Delta$  in a reduced space and be able to efficiently support update and search operations. Our result exploits an elegant technique proposed in [22, 25], here adapted to manage items of variable lengths (i.e. texts).

In the following we bound the space occupancy of our data structure in terms of the entropy of the concatenation of  $\Delta$ 's texts. A better overall space reduction might be possibly achieved by compressing separately the texts  $T_i$ 's. However, if the texts  $T_i$ 's have similar statistics, the entropy of the concatenated string is a reasonable lower bound to the compressibility of the collection. Furthermore, in the probabilistic setting where we assume that every text is generated by the same probabilistic source, the entropy of the concatenated string coincides with the entropy of the single texts and therefore provides a tight lower bound to the compressibility of the collection.

In the following we focus on the situation in which the length  $p$  of the searched pattern is  $O(\frac{u}{\log^2 u})$  because for the other range of  $p$ 's values, the search operation can be implemented in a brute-force way by first decompressing the text collection and by then searching for  $P$  into it using a scan-based string matching algorithm in  $O(p \log^3 u + occ)$  time. We partition the texts  $T_i$ 's into  $\eta = \Theta(\log^2 u)$  collections  $\mathcal{C}^1, \dots, \mathcal{C}^\eta$ , each containing texts of overall length  $O(\frac{u}{\log^2 u})$ . This is always possible, independently of the lengths of the text strings in  $\Delta$ , since the upper bound on the length of the searchable patterns allows us to split very long texts (i.e. texts of lengths  $\Omega(\frac{u}{\log^2 u})$ ) into  $2 \log^2 u$  pieces overlapping for  $\Theta(\frac{u}{\log^2 u})$  characters. This covering of a single long text with many shorter ones still allows us to find the occurrences of the searched patterns.

Every collection  $\mathcal{C}^h$  is then partitioned into a series of subsets  $\mathcal{S}_i^h$  defined as follows:  $\mathcal{S}_i^h$  contains some texts of  $\mathcal{C}^h$  having overall length in the range  $[2^i, 2^{i+1})$ , where  $i = O(\log u)$ . Each set  $\mathcal{S}_i^h$  is simultaneously indexed and compressed using our opportunistic data structure. Searching for an arbitrary pattern  $P[1, p]$  in  $\Delta$ , with  $p = O(\frac{u}{\log^2 u})$  can be performed by searching for  $P$  in all the  $O(\log^3 u)$  subsets  $\mathcal{S}_i^h$  via the compressed index built on each of them. This takes  $O(p \log^3 u + occ \log^\epsilon u)$  time overall.

Inserting a new text  $T[1, t]$  into  $\Delta$  consists of inserting  $T$  into one of the sets  $\mathcal{C}^h$ , the most empty one. Then,

the subset  $\mathcal{S}_i^h$  is selected, where  $i = \lfloor \log t \rfloor$ , and  $T$  is inserted into it using the following approach. If  $\mathcal{S}_i^h$  is empty then the compressed index is built for  $T$  and associated to this subset, thus taking  $O(t)$  time. Otherwise the new set  $\mathcal{S}_i^h \cup \{T\}$  is formed and inserted in  $\mathcal{S}_{i+1}^h$ . If the latter subset is not empty then the insertion process is propagated until an empty subset  $\mathcal{S}_{i+j}^h$  is found. At this point, the compressed index is built over the set  $\mathcal{S}_i^h \cup \dots \cup \mathcal{S}_{i+j-1}^h \cup \{T\}$ , by concatenating all the texts contained in this set to form a unique string, texts are separated by a special symbol (as usual). By noticing that these texts have overall length  $\Theta(2^{i+j})$ , we conclude that this propagation process has a complexity proportional to the overall length of the *moved* texts. Although each single insertion may be very costly, we can amortize this cost by charging  $O(\log u)$  credits per text character (since  $i, j = O(\log u)$ ), thus obtaining an overall amortized cost of  $O(t \log u)$  to insert  $T[1, t]$  in  $\Delta$ . Some care must be taken to evaluate the space occupied during the reconstruction of the set  $\mathcal{S}_i^h$ . In fact, the construction of our compressed index over the set  $\mathcal{S}_i^h$  requires the use of the suffix tree data structure (to compute the BWT) and thus  $O(2^i \log 2^i)$  bits of auxiliary storage. This could be too much, but we ensured that every collection  $\mathcal{C}^h$  contains texts having overall length  $O(\frac{u}{\log^2 u})$ . So that at most  $O(\frac{u}{\log u})$  bits suffices to support any reconstruction process.

We now show how to support text deletions from  $\Delta$ . The main problem here is that from one side we would like to physically cancel the texts in order to avoid the listing of *ghost* occurrences belonging to texts no longer in  $\Delta$ ; but from the other side a physical deletion would be too much time-consuming to be performed on-the-fly. Amortization can still be used but much care must be taken when answering a query to properly deal with texts which have been *logically* deleted from the  $\mathcal{S}_i^h$ 's. For the sake of presentation let  $T^{bw}$  be the BWT of the texts stored in some set  $\mathcal{S}_i^h$ . We store in a balanced search tree the set  $\mathcal{I}_i^h$  of interval positions in  $T^{bw}$  occupied by deleted text suffixes. If a pattern occurrence is found in  $T^{bw}$  using our compressed index, we can check in  $O(\log u)$  time if it is a real or a ghost occurrence. Every time a text  $T[1, t]$  must be deleted from  $\mathcal{S}_i^h$ , we search for all of its suffixes in  $\mathcal{S}_i^h$  and then update accordingly  $\mathcal{I}_i^h$  in  $O(t \log u)$  time. The additional space required to store the balanced search tree is  $O(|\mathcal{I}_i^h| \log u) = O(\frac{u}{\log^3 u})$  bits, where we are assuming to physically delete the texts from  $\mathcal{S}_i^h$  as soon as a fraction of  $\Theta(\frac{1}{\log^2 u})$  suffixes is logically marked. Hence, each set  $\mathcal{S}_i^h$  may undergo  $O(\log^2 u)$  reconstructions before it shrinks enough to move back to the previous set  $\mathcal{S}_{i-1}^h$ . Consequently the amortized cost of **delete** is  $O(t \log u + t \log^2 u) = O(t \log^2 u)$ , where the first term denotes the cost of  $\mathcal{I}_i^h$ 's update and the second term accounts for the credits to be left in order to pay for the physical deletions.

Finally, to identify a text to be deleted we append to every text in  $\Delta$  an identifier of  $O(\log u)$  bits, and we keep

track of the subset  $S_i^h$  containing a given text via a table. This introduces an overhead of  $O(m \log u)$  bits which is  $o(u)$  if we reasonably assume that the texts are not too short, i.e.  $\omega(\log u)$  bits each.

**Theorem 4** *Let  $\Delta$  be a dynamic collection of texts  $\{T_1, T_2, \dots, T_m\}$  having total length  $u$ . All the occurrences of a pattern  $P[1, p]$  in the texts of  $\Delta$  can be listed in  $O(p \log^3 u + occ \log u)$  time in the worst case. Operation insert adds a new text  $T[1, t]$  to  $\Delta$  in  $O(t \log u)$  amortized time. Operation delete removes a text  $T[1, t]$  from  $\Delta$  in  $O(t \log^2 u)$  amortized time. The space occupancy is  $O\left(H_k(\Delta) + m \frac{\log u}{u}\right) + o(1)$  bits per input symbol in the worst case for any fixed  $k \geq 0$ . ■*

## 5 A simple application

Glimpse [19] is an effective tool to index linguistic texts. From a high level point of view, it is a hybrid between inverted files and scan-based approaches with *no* index. It relies on the observation that there is no need to index every word with an exact location (as it occurs in inverted files); but only pointers to an area where the word occurs (called a *block*) should be maintained. Glimpse assumes that the text  $T[1, u]$  is *logically* partitioned into  $r$  blocks of size  $b$  each, and thus its index consists of two parts: a *vocabulary*  $V$  containing all the distinct words of the text; and for each word  $w \in V$ , a *list*  $L(w)$  of blocks where the word  $w$  occurs. This blocking scheme induces two space savings: pointers to word occurrences are shorter, and the occurrences of the same word in a single block are represented only once. Typically the index is very compact: 2-4% of the original text size [19].

Given this index structure, the search scheme proceeds in two steps: first the queried word  $w$  is searched in the vocabulary  $V$ , then all candidate blocks of  $L(w)$  are *sequentially* examined to find all the  $w$ 's occurrences. Complex queries (e.g. approximate or regular expression searches) can be supported by using `Agrep` [28] both in the vocabulary and in the block searches. Clearly, the search is efficient if the vocabulary is small, if the query is enough selective, and if the block size is not too large. The first two requirements are usually met in practice, so that the main constraint to the effective use of Glimpse remains the strict relation between block-pointer sizes and text sizes. Theoretical and experimental analysis of such block-addressing scheme [4, 19] have shown that the Glimpse approach is effective only for medium sized texts (roughly up to 200Mb). Recent papers tried to overcome this limitation by compressing each text block individually and then searching it via proper opportunistic string-matching algorithms [19, 24]. The experimental results showed an improvement of about 30-50% in the final performance, thus implicitly proving that the second searching step dominates Glimpse's query performance.

Our opportunistic index naturally fits in this block-addressing framework and allows us to extend its applicability to larger text databases. The new approach, named *Compressed Glimpse* (shortly `CGlimpse`), consists in using our opportunistic data structure to index each text block individually; this way, each candidate block is not fully scanned at query time but its index is employed to fasten the detection of the pattern occurrences. In some sense `CGlimpse` is a compromise between a full-text index (like a suffix array) and a word-based index (like an inverted list) over a compressed text.

A theoretical investigation of the performance of `CGlimpse` is feasible using a model generally accepted in Information Retrieval [4]. It assumes the Heaps law to model the vocabulary size (i.e.  $V = O(u^\beta)$  with  $0 < \beta < 1$ ), the generalized Zipf's law to model the frequency of words in the text collection (i.e. the largest  $i$ th frequency of a word is  $u/(i^\theta H_V^{(\theta)})$ , where  $H_V^{(\theta)}$  is a normalization term and  $\theta$  is a parameter larger than 1), and assumes that  $O(u^\rho)$  is the number of matches for a given word with  $k \geq 1$  errors (where  $\rho < 1$ ). Under these hypothesis we can show that `CGlimpse` achieves *both sublinear space overhead and sublinear query time independent of the block size* (proof in the full paper). Conversely, inverted indices achieve only the second goal [27], and classical Glimpse achieves both goals but under some restrictive conditions on the block size [4].

## 6 Conclusions

Some issues remain still to be investigated in various models of computation. In external memory, it would be interesting to devise a compressed index which takes advantage of the blocked access to the disk and thus achieves  $O(occ/B)$  I/Os for locating the pattern occurrences, where  $B$  is the disk-page size. In the RAM, it would be interesting to avoid the  $o(\log u)$  overhead incurred in the listing of the pattern occurrences. In the full paper we will show how to use known techniques (see e.g. [11]) for designing hybrid indices which achieve  $O(occ)$  retrieval time cost under restrictive conditions either on the pattern length or on the number of pattern occurrences. Guaranteeing the  $\Theta(occ)$  retrieval cost in the general case is an open problem also in the uncompressed setting [12].

## References

- [1] A. Amir and G. Benson. Efficient two-dimensional compressed matching. *Proceedings of IEEE Data Compression Conference*, pages 279–288, 1992.
- [2] A. Amir, G. Benson, and M. Farach. Let sleeping files lie: Pattern matching in Z-compressed files. *Journal of Computer and System Sciences*, 52(2):299–307, 1996.
- [3] A. Andersson. Sorting and searching revisited. In R. G. Karlsson and A. Lingas, editors, *Proceedings of the 5th*



- Scandinavian Workshop on Algorithm Theory*, pages 185–197. Springer-Verlag LNCS n. 1097, 1996.
- [4] R. Baeza-Yates and G. Navarro. Block addressing indices for approximate text retrieval. *Journal of the American Society for Information Science*, 51(1):69–82, 2000.
- [5] J. Bentley. *Programming Pearls*. Addison-Wesley, USA, 1989.
- [6] J. Bentley, D. Sleator, R. Tarjan, and V. Wei. A locally adaptive compression scheme. *Communication of the ACM*, 29(4):320–330, 1986.
- [7] M. Burrows and D. Wheeler. A block sorting lossless data compression algorithm. Technical Report 124, Digital Equipment Corporation, 1994.
- [8] S. Chen and J. Reif. Using difficulty of prediction to decrease computation: Fast sort, priority queue and convex hull on entropy bounded inputs. In *Proceedings of the 34th IEEE Symposium on Foundations of Computer Science*, pages 104–112, 1993.
- [9] M. Farach and M. Thorup. String matching in Lempel-Ziv compressed strings. *Algorithmica*, 20(4):388–404, 1998.
- [10] P. Ferragina and R. Grossi. The string B-tree: A new data structure for string search in external memory and its applications. *Journal of the ACM*, 46:236–280, 1999.
- [11] P. Ferragina, S. Muthukrishnan, and M. deBerg. Multi-method dispatching: A geometric approach with applications to string matching problems. In *Proceedings of the 31st ACM Symposium on the Theory of Computing*, pages 483–491, 1999.
- [12] R. Grossi and J. Vitter. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. In *Proceedings of the 32nd ACM Symposium on Theory of Computing*, 2000.
- [13] J. Kärkkäinen and E. Sutinen. Lempel-Zip index for  $q$ -grams. In J. Díaz and M. Serna, editors, *Proceedings of the 4th European Symposium on Algorithms*, pages 378–391. Springer-Verlag LNCS n. 1136, 1996.
- [14] J. Kärkkäinen and E. Ukkonen. Lempel-Ziv parsing and sublinear-size index structures for string matching. In N. Ziviani, R. Baeza-Yates, and K. Guimarães, editors, *Proceedings of the 3rd South American Workshop on String Processing*, pages 141–155. Carleton University Press, 1996.
- [15] A. Karlin, S. Phillips, and P. Raghavan. Markov paging (extended abstract). In *Proceedings of the 33rd IEEE Symposium on Foundations of Computer Science*, pages 208–217, 24–27 Oct. 1992.
- [16] D. E. Knuth. *Sorting and Searching*, volume 3 of *The Art of Computer Programming*. Addison-Wesley, Reading, MA, USA, second edition, 1998.
- [17] P. Krishnan and J. Vitter. Optimal prediction for prefetching in the worst case. *SIAM Journal on Computing*, 27(6):1617–1636, Dec. 1998.
- [18] U. Manber and G. Myers. Suffix arrays: a new method for on-line string searches. *SIAM Journal on Computing*, 22(5):935–948, 1993.
- [19] U. Manber and S. Wu. GLIMPSE: A tool to search through entire file systems. In *Proceedings of the USENIX Winter 1994 Technical Conference*, pages 23–32, 1994.
- [20] G. Manzini. An analysis of the Burrows-Wheeler transform. In *Proceedings of the 10th ACM-SIAM Symposium on Discrete Algorithms*, pages 669–677, 1999. Full version in [www.imc.pi.cnr.it/~manzini/tr-99-13/](http://www.imc.pi.cnr.it/~manzini/tr-99-13/).
- [21] E. M. McCreight. A space-economical suffix tree construction algorithm. *Journal of the ACM*, 23(2):262–272, 1976.
- [22] K. Mehlhorn and M. H. Overmars. Optimal dynamization of decomposable searching problems. *Information Processing Letters*, 12(2):93–98, Apr. 1981.
- [23] J. I. Munro. Succinct data structures. In *Proceeding of the 19th Conference on Foundations of Software Technology and Theoretical Computer Science*. Springer-Verlag LNCS n. 1738, 1999.
- [24] G. Navarro, E. de Moura, M. Neubert, N. Ziviani, and R. Baeza-Yates. Adding compression to block addressing inverted indexes. *Information Retrieval Journal*, 2000, (to appear).
- [25] M. H. Overmars and J. van Leeuwen. Worst-case optimal insertion and deletion methods for decomposable searching problems. *Information Processing Letters*, 12(4):168–173, Aug. 1981.
- [26] K. Sadakane. A modified Burrows-Wheeler transformation for case-insensitive search with application to suffix array compression. In *Proceedings of IEEE Data Compression Conference*, 1999.
- [27] I. H. Witten, A. Moffat, and T. C. Bell. *Managing Gigabytes: Compressing and Indexing Documents and Images*. Morgan Kaufmann Publishers, Los Altos, CA 94022, USA, second edition, 1999.
- [28] S. Wu and U. Manber. AGREP - A fast approximate pattern-matching tool. In *Proceedings of the Usenix Winter 1992 Technical Conference*, pages 153–162. Usenix Association, 1992.