

# Dynamic Entropy-Compressed Sequences and Full-Text Indexes

VELI MÄKINEN  
University of Helsinki

and  
GONZALO NAVARRO  
University of Chile

---

---

First author funded by the Academy of Finland under grant 108219. Second author partially funded by Fondecyt Grant 1-050493, Chile. A preliminary partial version of this paper appeared in *Proc. CPM 2006, LNCS 4009*.

Author's address: Gonzalo Navarro, Department of Computer Science, University of Chile, Blanco Encalada 2120, Santiago, Chile. [gnavarro@dcc.uchile.cl](mailto:gnavarro@dcc.uchile.cl). Veli Mäkinen, Department of Computer Science, P. O. Box 68 (Gustaf Hällströmin katu 2b), FIN-00014 University of Helsinki, Finland. [vmakinen@cs.helsinki.fi](mailto:vmakinen@cs.helsinki.fi).

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 20YY ACM 0000-0000/20YY/0000-0001 \$5.00

We give new solutions to the SEARCHABLE PARTIAL SUMS WITH INDELS problem. Given a sequence of  $n$   $k$ -bit numbers, we present a structure taking  $kn + o(kn)$  bits of space, able of performing operations *sum*, *search*, *insert*, and *delete*, all in  $O(\log n)$  worst-case time, for any  $k = O(\log n)$ . This extends previous results by Hon et al. [ISAAC 2003] achieving the same space and  $O(\log n / \log \log n)$  time complexities for the queries, yet offering complexities for *insert* and *delete* that are amortized and worse than ours, and supported only for  $k = O(1)$ . Our result matches an existing lower bound for large values of  $k$ .

We also give new solutions to the DYNAMIC SEQUENCE problem. Given a sequence of  $n$  symbols in the range  $[1, \sigma]$  with binary zero-order entropy  $H_0$ , we present a dynamic data structure that requires  $nH_0 + o(n \log \sigma)$  bits of space, which is able of performing *rank* and *select*, as well as inserting and deleting symbols at arbitrary positions, in  $O(\log n \log \sigma)$  time. Our result is the *first* entropy-bound dynamic data structure for *rank* and *select* over general sequences.

In the case  $\sigma = 2$ , where both previous problems coincide, we improve the dynamic solution of Hon et al. in that we compress the sequence. The only previous result with entropy-bound space for dynamic binary sequences is by Blandford and Blelloch [SODA 2004], which has the same complexities as our structure, but does not achieve constant 1 multiplying the entropy term in the space complexity.

Finally, we present a new dynamic compressed full-text self-index, for a collection of texts over an alphabet of size  $\sigma$ , of overall length  $n$  and  $h$ -th order empirical entropy  $H_h$ . The index requires  $nH_h + o(n \log \sigma)$  bits of space, for any  $h \leq \alpha \log_\sigma n$  and constant  $0 < \alpha < 1$ . It can count the number of occurrences of a pattern of length  $m$  in time  $O(m \log n \log \sigma)$ . Each such occurrence can be reported in  $O(\log^2 n \log \log n)$  time, and displaying a context of length  $\ell$  from a text takes time  $O(\log n (\ell \log \sigma + \log n \log \log n))$ . Insertion/deletion of a text to/from the collection takes  $O(\log n \log \sigma)$  time per symbol. This significantly improves the space of a previous result by Chan et al. [CPM 2004] in exchange for a slight time complexity penalty. We achieve at the same time the *first* dynamic index requiring essentially  $nH_h$  bits of space, and the *first* construction of a compressed full-text self-index within that working space. Previous results achieve at best  $O(nH_h)$  space with constants larger than 1 (Ferragina and Manzini [FOCS 2000], Arroyuelo and Navarro [ISAAC 2005]) and higher time complexities.

An important result we prove in this paper is that the wavelet tree of the Burrows-Wheeler transform of a text, if compressed with a technique that achieves zero-order compression locally (e.g., Raman et al. [SODA 2002]), automatically achieves  $h$ -th order entropy space for any  $h$ . This unforeseen relation is essential for the results of the previous paragraph, but it also derives into significant simplifications on many existing static compressed full-text self-indexes that build on wavelet trees.

Categories and Subject Descriptors: E.1 [Data structures]: ; E.2 [Data storage representations]: ; E.4 [Coding and information theory]: Data compaction and compression; F.2.2 [Analysis of algorithms and problem complexity]: Nonnumerical algorithms and problems—*Pattern matching, Computations on discrete structures, Sorting and searching*; H.2.1 [Database management]: Physical design—*Access methods*; H.3.2 [Information storage and retrieval]: Information storage—*File organization*; H.3.3 [Information storage and retrieval]: Information search and retrieval—*Search process*

General Terms: Algorithms

Additional Key Words and Phrases: Compressed dynamic data structures, compressed text databases, entropy, partial sums, sequences.

## 1. INTRODUCTION AND RELATED WORK

The study of compressed data structures aims to represent classical structures like trees, graphs, text indexes, etc., in the smallest possible space without challenging the functionality of the structure; the original operations should be supported efficiently without decompressing the whole structure.

The well-known SEARCHABLE PARTIAL SUMS problem consists in maintaining a sequence  $A$  of nonnegative integers  $a_1 \dots a_n$ , each of  $k$  bits, supporting queries on the prefix sums and limited updates on the values. An extension [Hon et al. 2003b] called SEARCHABLE PARTIAL SUMS WITH INDELS allows insertions and deletions of values as well.

The restricted version where  $k = 1$ , and thus the numbers are actually bits, is called the DYNAMIC BIT VECTOR problem. In this case the update operation means flipping the bit. Also, queries on the prefix sums correspond to the well-known *rank* and *select* queries on bit sequences. The extension to the DYNAMIC BIT VECTOR WITH INDELS problem is immediate.

Some recent articles deal with the problem of designing succinct dynamic data structures, requiring  $kn + o(kn)$  bits of space, for these problems [Raman et al. 2001; Hon et al. 2003b]. The best current result for  $k = O(1)$  [Hon et al. 2003b]  $kn + o(kn)$  bits of space,  $O(\log_b n)$  time for queries, and  $O(b)$  time for updates, for any  $b = \Omega(\log n / \log \log n)$ . Insertions and deletions can be supported in  $O(b)$  amortized time, for  $b = \Omega(\log n / \log \log n)^2$ . These results scale well to larger  $k$ , obtaining the same worst-case complexities for prefix sum queries as well as for updates, but insertions and deletions are not supported anymore<sup>1</sup>. The best time complexities achieved for this problem are worst-case  $O(\log n / \log(w/\delta))$ , where updates that add/subtract a number of  $\delta$  bits are permitted [Patrascu and Demaine 2004]. They show that this complexity is optimal, but insertions and deletions are not considered.

In this paper we extend this result by achieving  $kn + o(kn)$  bits of space and  $O(\log n)$  worst case time complexity for all the operations, for any  $k = O(\log n)$ . For larger  $k = O(w)$ , where  $w$  is the machine word size under the RAM model of computation, our time complexity raises to  $O(w / \log^{1-\varepsilon} n + \log n)$ , for any constant  $\varepsilon > 0$ . Moreover, we refine this result for the case of strictly positive numbers: the total space is not anymore  $kn$  bits, but the sum of the exact number of bits necessary to represent each number. For the case  $k = \Theta(\log n) = \Theta(w)$  the obtained time complexity is optimal, given the lower bound  $\Omega(\log n / \log(w/k)) = \Omega(\log n)$  for the more restricted problem with (arbitrary) updates [Patrascu and Demaine 2004]<sup>2</sup>.

Let us return to the DYNAMIC BIT VECTOR WITH INDELS problem. We go further than representing them using  $n + o(n)$  bits, and achieve a representation that uses  $nH_0 + o(n)$  bits of space, where  $0 < H_0 \leq 1$  is the empirical zero-order entropy of the sequence. This is an improvement over an  $O(nH_0)$  result by Blandford and Blelloch [Blandford and Blelloch 2004] since, although their results are more general, they do not achieve constant 1 multiplying the entropy term in the space complexity. The comparison of time complexities against partial sums with  $k = 1$  [Hon et al. 2003b] stays as above, but now we compress the sequence. Our space complexity has been only achieved in the static scenario [Raman et al. 2001], where no updates are possible.

<sup>1</sup>They can scale up to  $k = O(\log \log n)$  under some extra assumptions.

<sup>2</sup>Note that, even if we wish to restrict our updates to  $\delta < k$  bits, insertions and deletions would permit simulating general updates, thus the stated lower bound must hold anyway if insertions and deletions are permitted, unless we restrict these in a rather unnatural way.

We further generalize this result to sequences of symbols over an alphabet  $[1, \sigma]$ , where operations *rank* and *select* generalize to  $rank_a(A, i)$ , counting the occurrences of symbol  $a$  in  $A[1, i]$ , and  $select_a(A, j)$ , giving the position of the  $j$ -th  $a$  in  $A$ . By using wavelet trees [Grossi et al. 2003] we achieve  $nH_0 + o(n \log \sigma)$  bits of space and  $O(\log n \log \sigma)$  time complexities. This space has been previously achieved only for static data structures, with query time  $O(\lceil \log \sigma / \log \log n \rceil)$  time for *rank* and *select* but no support for updates [Raman et al. 2002; Ferragina et al. 2007]. By using multiary wavelet trees [Ferragina et al. 2007], we can reduce the query times to  $O(\frac{1}{\epsilon} \log n \lceil \log \sigma / \log \log n \rceil)$  in exchange for increasing the update times to  $O(\frac{1}{\epsilon} \log^{1+\epsilon} n / \log \log n)$ , for any  $\epsilon > 0$ .

Moreover, all our results work under weaker assumptions on the RAM model than many previous results on dynamic settings. Instead of assuming  $\log n = \Theta(w)$  as it is frequent in the literature, we opt for the weaker condition  $\log n = O(w)$ . We show that the results stay the same in terms of time and space complexities, except for  $O(w)$  extra bits of space that are required for a constant number of system pointers.

Let us now regard sequences of symbols with another semantics. The indexed string matching problem is that of, given a long text  $T[1, n]$  over an alphabet  $\Sigma$  of size  $\sigma$ , building a data structure called *full-text index* on it, to solve two types of queries: (a) Given a short pattern  $P[1, m]$  over  $\Sigma$ , *count* the occurrences of  $P$  in  $T$ ; (b) *locate* those *occ* positions in  $T$ . There are several classical full-text indexes requiring  $O(n \log n)$  bits of space which can answer counting queries in  $O(m \log \sigma)$  time (like suffix trees [Apostolico 1985]) or  $O(m + \log n)$  time (like suffix arrays [Manber and Myers 1993]). Both locate each occurrence in constant time once the counting is done. Similar complexities are obtained with modern compressed data structures [Ferragina and Manzini 2000; Grossi et al. 2003; Ferragina et al. 2007], requiring space  $nH_h + o(n \log \sigma)$  bits (for some small  $h$ ), where  $H_h \leq \log \sigma$  is the  $h$ -th order empirical entropy of  $T$ .<sup>3</sup> These indexes are often called *compressed self-indexes* referring to their space requirement and to their ability to work without the text and even to fully replace it, by delivering any text substring without accessing  $T$ . A *dynamic* self-index permits managing a collection of texts and inserting/deleting texts to/from the collection.

There exist dynamic self-indexes by Chan et al. [Chan et al. 2004; Chan et al. 2007]. One version requires  $O(\sigma n)$  bits of space, and it can count the number of occurrences of a pattern of length  $m$  in time  $O(m \log n)$ . Insertions and deletions require  $O(\sigma \log n)$  and  $O((\sigma + \log n) \log n)$  time per character, respectively. A second version requires  $O(n \log \sigma)$  bits of space and counts in time  $O(m \log^2 n)$ . Insertions and deletions take  $O(\log^2 n)$  time per character. In either case, each occurrence position can be retrieved in  $O(\log^2 n)$  time. Both structures can be combined so as to get  $O(\sigma n)$  bits of space,  $O(m \log n)$  counting time, and  $O(\sigma \log n)$  insertion and deletion time per character.

The main building block in compressed self-indexes is function  $rank_a$ . Actually, our dynamic compressed  $rank_a$  structure can be used to implement a dynamic compressed self-index that takes  $nH_h + o(n \log \sigma)$  bits of space, for any  $h \leq \alpha \log_\sigma n$

<sup>3</sup>In this paper  $\log$  stands for  $\log_2$ .

and constant  $0 < \alpha < 1$ . This is obtained by plugging our structure for symbol sequences in the dynamic self-index of [Chan et al. 2004]. Our counting time is  $O(m \log n \log \sigma)$ . We can locate each occurrence in time  $O(\log^2 n \log \log n)$ , and display a text context of length  $\ell$  in time  $O(\log n (\ell \log \sigma + \log n \log \log n))$ . Insertion and deletion of a text to the collection takes  $O(\log n \log \sigma)$  time per symbol. Compared with the original indexes of Chan et al., we obtain a significant space saving and, depending on the case, better update or better counting times.

The fact that plugging our  $nH_0$ -bits sequence representation into the self index of [Chan et al. 2004] yields  $h$ -th order compression stems from the fact that the sequence we are representing is the Burrows-Wheeler transform of the text collection [Burrows and Wheeler 1994; Manzini 2001]. This is a striking result we prove in this paper. For several years, much effort has been spent in designing sophisticated (static) data structures on top of the plain wavelet tree so as to reduce its  $nH_0$ -bit size to  $nH_h$  bits [Grossi et al. 2003; Ferragina et al. 2007; Mäkinen and Navarro 2005]. In this paper we show that this is automatically achieved by the *original* wavelet tree without any further effort! Thus, as a byproduct, we obtain a significant simplification in the design of static data structures for this problem.

Ours is the *first* dynamic compressed self-index with space essentially equal to the  $h$ -th order empirical entropy of the text collection, which in addition can be built within this working space. We know only of two previous dynamic full-text self-indexes. The older [Ferragina and Manzini 2000] requires  $O(nH_h)$  bits of space (with constant 5 at least),  $O(m \log^3 n)$  counting time,  $O(\log n)$  amortized insertion time per character, and  $O(\log^2 n)$  amortized deletion time per character. A newer one [Chan et al. 2004; Chan et al. 2007] requires  $O(\sigma n)$  bits of space,  $O(m \log n)$  counting time,  $O(\log^2 n)$  locating time per occurrence, and  $O(\sigma \log n)$  insertion/deletion time per character.

As a plus, we obtain an  $O(n \log n \log \sigma)$  time construction algorithm for static self-index requiring  $nH_h + o(n \log \sigma)$  bits *working space* during construction (the same as the final structure). Previous construction algorithms within entropy space achieve  $O(nH_0)$  bits of space and  $O(n \log n)$  time [Hon et al. 2003], or  $O(nH_h)$  bits of space (with constant larger than 4) and  $O(\sigma n)$  time [Arroyuelo and Navarro 2005].

Several other compressed indexes can be obtained using our algorithm. Moreover, it is very easy to obtain the Burrows-Wheeler transform of  $T$  from the index we build, within the same  $O(n \log n \log \sigma)$  time. A recent result [Kärkkäinen 2004] achieves  $n \log \sigma + O(n)$  bits and  $O(n \log^2 n)$  time.

## 2. DEFINITIONS

To simplify notation, we ignore roundings. When referring to number of bits, we use simply  $\log n$  to refer to  $\lfloor (\log n) + 1 \rfloor$ . That is,  $\log \log n$  bits means actually  $\lfloor (\log \lfloor (\log n) + 1 \rfloor) + 1 \rfloor$  bits. Similarly  $(\log n)/2$  is the integer nearest to  $\lfloor (\log n) + 1 \rfloor / 2$ , and so on.

We assume our sequence  $A = a_1 \dots a_n$  to be drawn from an alphabet  $\{0, 1, \dots, \sigma - 1\}$ . Let  $n_c$  denote the number of occurrences of symbol  $c$  in  $A$ , i.e.,  $n_c = |\{i \mid a_i = c\}|$ . Then the zero-order *empirical entropy* is defined as  $H_0(A) = \sum_{0 \leq c < \sigma} \frac{n_c}{n} \log \frac{n}{n_c}$ . This is the lower bound for the average codeword length of any compressor that

fixes the codewords to the symbols independently of the context they appear in. A tighter lower bound for the compressibility of sequences is the *h-th order empirical entropy*  $H_h(A)$ , where the compressor can fix the codeword based on the *h*-symbol context following the symbol to be coded.<sup>4</sup> Formally, it can be defined as  $H_h(A) = \sum_{x \in \Sigma^h} \frac{n_x}{n} H_0(A|x)$ , where  $n_x$  denotes the number of occurrences of substring  $x$  in  $A$  and  $A|x$  denotes the concatenation of the symbols appearing immediately before those  $n_x$  occurrences [Manzini 2001]. Substring  $x = A[i+1, i+h]$  is called a *h-context* of symbol  $a_i$ . We take  $A$  here as a *cyclic string*, such that  $a_n$  precedes  $a_1$ , and thus the amount of *h*-contexts is exactly  $n$ .

We assume a random access machine with word size  $w$ ; typical arithmetic operations on  $w$ -bit integers are assumed to take constant time. We make the minimal assumption that  $\log n = O(w)$ , instead of the common stronger assumption  $\log n = \Theta(w)$ .

We study the following problems in this paper:

The DYNAMIC SEQUENCE WITH INDELS problem is to maintain a (virtual) sequence  $A = a_1 \dots a_n$ ,  $a_i \in \{0, 1, \dots, \sigma - 1\}$ , supporting the operations:

- *read*( $A, i$ ) obtains symbol  $a_i$ ;
- *rank<sub>c</sub>*( $A, i$ ) returns the number of occurrences of symbol  $c$  in  $a_1 \dots a_i$ ;
- *select<sub>c</sub>*( $A, j$ ) returns the index  $i$  containing  $j$ -th occurrence of  $c$ ;
- *insert*( $A, c, i$ ) inserts  $c \in \{0, 1, \dots, \sigma - 1\}$  between  $a_{i-1}$  and  $a_i$ ; and
- *delete*( $A, i$ ) deletes  $a_i$  from the sequence.

The DYNAMIC BIT VECTOR WITH INDELS problem is a restriction of the above to alphabet  $\{0, 1\}$  (i.e.,  $\sigma = 2$ ). Then we use short-hand notation  $\text{rank}(A, i) = \text{rank}_1(A, i)$  and  $\text{select}(A, i) = \text{select}_1(A, i)$ . Notice that  $\text{rank}_0(A, i) = i - \text{rank}_1(A, i)$ , but the same does not apply for  $\text{select}_0(A, j)$ ; so both *select* queries must be handled.

The SEARCHABLE PARTIAL SUMS problem consists in maintaining a sequence  $A$  of nonnegative integers  $a_1 \dots a_n$ , each of  $k$  bits, so that we can perform the following queries and operations on them:

- *sum*( $A, i$ ) returns  $\sum_{t=1}^i a_t$ ;
- *search*( $A, j$ ) returns the smallest  $i$  such that  $\text{sum}(A, i) \geq j$ ; and
- *update*( $A, i, \Delta$ ) increases  $a_i$  by  $\Delta$ , assuming  $a_i + \Delta$  is within bounds and  $\Delta = O(\text{polylog}(n))$ .

A more general problem called SEARCHABLE PARTIAL SUMS WITH INDELS includes also the following operations:

- *insert*( $A, i, x$ ) inserts  $x$  between  $a_{i-1}$  and  $a_i$ .
- *delete*( $A, i$ ) deletes  $a_i$  from the sequence.

<sup>4</sup>It is more logical (and hence customary) to define the context as the  $h$  symbols preceding a symbol, but we use the reverse definition for technical convenience. If this is an issue, the sequences can be handled in reverse order to obtain results on the more standard definition. It is anyway known that both definitions differ by lower order terms only [Ferragina and Manzini 2005].

Notice that the SEARCHABLE PARTIAL SUMS WITH INDELS problem with  $k = 1$  is equivalent to the DYNAMIC BIT VECTOR WITH INDELS problem (*sum* being *rank*, *search* being *select*).

A problem related to the DYNAMIC SEQUENCE WITH INDELS problem is the DYNAMIC TEXT COLLECTION problem, defined as follows: Maintain a dynamic collection  $\mathcal{C}$  of texts  $\{T_1, T_2, \dots, T_m\}$ , where each  $T_i \in \{1, 2, \dots, \sigma\}^*$ , supporting the following operations:

- count*( $\mathcal{C}, P$ ) returns the number of times pattern  $P$  occurs as a substring in the collection;
- locate*( $\mathcal{C}, P$ ) returns the occurrence positions of  $P$  in the collection;
- substring*( $\mathcal{C}, j, l, r$ ) returns  $T_j[l, r]$ ;
- $j = \textit{insert}(\mathcal{C}, T)$  inserts text  $T$  into the collection, returning a handle  $j$  to it (that is, from now on  $T = T_j$ ); and
- delete*( $\mathcal{C}, j$ ) deletes text  $T_j$  from the collection.

### 3. PREVIOUS RESULTS

Our new solutions build on top of various previous results. We explain part of these previous results in detail, in order to present our contribution in a self-contained manner.

#### 3.1 Static Entropy-Bound Structures for Bit Vectors

Raman et al. [Raman et al. 2002] proposed a data structure to solve *rank* and *select* queries in constant time over a static bit vector  $A = a_1 \dots a_n$  with binary zero-order entropy  $H_0$ . The structure requires  $nH_0 + o(n)$  bits.

The idea is to split  $A$  into *superblocks*  $S_1 \dots S_{n/s}$  of  $s = \log^2 n$  bits. Each superblock  $S_i$  is in turn divided into  $2 \log n$  blocks  $B_i(j)$ , of  $t = (\log n)/2$  bits each, thus  $1 \leq j \leq s/t$ . Each such block  $B_i$  is said to belong to *class*  $c$  if it has exactly  $c$  bits set, for  $0 \leq c \leq t$ . For each class  $c$ , a universal table  $G_c$  of  $\binom{t}{c}$  entries is precomputed. Each entry corresponds to a possible block belonging to class  $c$ , and it stores all the local *rank* answers for that block. Overall all the  $G_c$  tables add up  $2^t = \sqrt{n}$  entries, and  $O(\sqrt{n} \text{ polylog}(n))$  bits.

Each block  $B_i(j)$  of the sequence is represented by a pair  $D_i(j) = (c, o)$ , where  $c$  is its class and  $o$  is the index of its corresponding entry in table  $G_c$ . A block of class  $c$  thus requires  $\log(t+1) + \log \binom{t}{c}$  bits. The first term is  $O(\log \log n)$ , whereas all the second terms add up  $nH_0 + O(n/\log n)$  bits. To see this, note that  $\log \binom{t}{c_1} + \log \binom{t}{c_2} \leq \log \binom{2t}{c_1+c_2}$ , and that  $nH_0 \geq \log \binom{t(n/t)}{c_1+\dots+c_{n/t}}$ . The pairs  $D_i(j)$  are of variable length and are all concatenated into a single sequence.

Each superblock  $S_i$  stores a pointer  $P_i$  to its first block description in the sequence (that is, the first bit of  $D_i(1)$ ) and the *rank* value at the beginning of the superblock,  $R_i = \textit{rank}(A, (i-1)s)$ .  $P$  and  $R$  add up  $O(n/\log n)$  bits. In addition,  $S_i$  contains  $s/t$  numbers  $L_i(j)$ , giving the initial position of each of its blocks in the sequence, relative to the beginning of the superblock. That is,  $L_i(j)$  is the position of  $D_i(j)$  minus  $P_i$ . Similarly,  $S_i$  stores  $s/t$  numbers  $Q_i(j)$  giving the *rank* value at the beginning of each of its blocks, relative to the beginning of the superblock. That

is,  $Q_i(j) = \text{rank}(A, (i-1)s + (j-1)t) - R_i$ . As those relative values are  $O(\log n)$ , sequences  $L$  and  $Q$  require  $O(n \log \log n / \log n)$  bits.

To solve  $\text{rank}(A, p)$ , we compute the corresponding superblock  $i = 1 + \lfloor (p-1)/s \rfloor$  and block  $j = 1 + \lfloor (p - (i-1)s - 1)/t \rfloor$ . Then we add the  $\text{rank}$  value of the corresponding superblock,  $R_i$ , the relative  $\text{rank}$  value of the corresponding block,  $Q_i(j)$ , and complete the computation by fetching the description  $(c, o)$  of the block where  $p$  belongs (from bit position  $P_i + L_i(j)$ ) and performing a (precomputed) local  $\text{rank}$  query in the universal table,  $\text{rank}(G_c(o), p - (i-1)s - (j-1)t)$ .

The overall space requirement is  $nH_0 + O(n \log \log n / \log n)$  bits, and  $\text{rank}$  is solved in constant time. We do not cover  $\text{select}$  because it is not necessary to follow this paper.

The scheme extends to sequences over small alphabets as well [Ferragina et al. 2007]. Let  $B = a_1 \dots a_t$  be the symbols in a block, and call  $n_a$  the number of occurrences of symbol  $a \in [1, q]$  in  $B$ . We call  $(n_1, \dots, n_q)$  the *class* of  $B$ . Thus, in our  $(c, o)$  pairs,  $c$  will be a number identifying the class of  $B$  and  $o$  an index within the class. A simple upper bound to the number of classes is  $(t+1)^q$  (as a class is a tuple of  $q$  numbers in  $[0, t]$ , although they have to add up  $t$ ). Thus  $O(q \log \log n)$  bits suffice for  $c$  (a second bound on the number of classes is  $q^t$  as there cannot be more classes than different sequences). Just as in the binary case, the sum of the sizes of all  $o$  fields adds up  $nH_0(A) + O(n/\log_q n)$  [Ferragina et al. 2007].

### 3.2 Static Wavelet Trees and Entropy-Bound Structures for Sequences

We now extend the result of the previous section to larger alphabets. The idea is to build a wavelet tree [Grossi et al. 2003] over sequences represented using  $\text{rank}$  and  $\text{select}$  structures for small alphabets.

A binary wavelet tree is a balanced binary tree whose leaves represent the symbols in the alphabet. The root is associated with the whole sequence  $A = a_1 \dots a_n$ , its left child with the subsequence of  $A$  obtained by concatenating all positions  $i$  having  $a_i < \sigma/2$ , and its right child with the complementary subsequence (symbols  $a_i \geq \sigma/2$ ). This subdivision is continued recursively, until each leaf contains a repeat of one symbol. The sequence at each node is represented by a bit vector that tells which positions (those marked with 0) go to the left child, and which (marked with 1) go to the right child. It is easy to see that the bit vectors alone are enough to determine the original sequence: To recover  $a_i$ , start at the root and go left or right depending on the bit vector value  $B_i$  at the root. When going to the left child, replace  $i \leftarrow \text{rank}_0(B, i)$ , and similarly  $i \leftarrow \text{rank}_1(B, i)$  when going right. When arriving at the leaf of character  $c$  it must hold that the original  $a_i$  is  $c$ . This requires  $O(\log \sigma)$   $\text{rank}$  operations over bit vectors.

It also turns out that operations  $\text{rank}$  and  $\text{select}$  on the original sequence can be carried out via  $O(\log \sigma)$  operations of the same type on the bit vectors of the wavelet tree [Grossi et al. 2003]. For example, to solve  $\text{rank}_c(A, i)$ , start at the root and go to the left child if  $c < \sigma/2$  and to the right child otherwise. When going down, update  $i$  as in the previous paragraph. When arriving at the leaf of  $c$ , the current  $i$  value is the answer. For  $\text{select}_c(A, j)$ , the algorithm starts at the leaf of  $c$  and goes upwards until the root, updating  $j \leftarrow \text{select}_b(B, j)$  with  $b = 0$  or  $1$  depending on whether we descend from the parent (owning vector  $B$ ) from the left or right child.



A multiary wavelet tree, of arity  $q$ , is used in [Ferragina et al. 2007]. In this case the sequence of each wavelet tree node ranges over alphabet  $[1, q]$ , and symbol rank/select queries are needed over those sequences. One needs  $O(\log_q \sigma)$  operations on those sequences to perform the corresponding operation on the original sequence.

Either for binary or general wavelet trees, it can be shown that the  $H_0$  entropies in the representations of the sequences at each level add up to  $nH_0(A)$  bits [Grossi et al. 2003; Ferragina et al. 2007]. The space occupancy of the sublinear structures adds up to  $o(n \log \sigma)$  if  $\sigma = o(n)$ . Overall, the structure requires  $nH_0(A) + o(n \log \sigma)$  bits<sup>5</sup>, solving *rank* in  $O(\log \sigma)$  time. Within the same bounds one can solve *select* as well [Raman et al. 2002; Grossi et al. 2003].

It is more space-efficient for the  $o(n \log \sigma)$  part to concatenate all the bitmaps of each level, so that we handle  $\log \sigma$  bitmaps of length  $n$ . It is then possible to do *rank* without any tree pointer [Ferragina et al. 2007], yet one needs equivalent upward pointers for *select*.

One can also use multiary wavelet trees and represent the sequences with alphabet size  $q$  using the techniques for small alphabets (see the end of previous section). With a suitable value for  $q$ , one obtains a structure requiring the same  $nH_0(A) + o(n \log \sigma)$  bits of space, but answering *rank* and *select* in constant time when  $\sigma = O(\text{polylog}(n))$ , and  $O(\lceil \log \sigma / \log \log n \rceil)$  time in general [Ferragina et al. 2007].

### 3.3 Dynamic Structures for Bit Vectors

Hon et al. [Hon et al. 2003b] show how to handle a bit vector  $A = a_1 \dots a_n$  in  $n + o(n)$  bits of space, so that *rank* and *select* can be solved in  $O(\log_b n)$  time, while insertions and deletions to the sequence can be handled in  $O(b)$  time, for any parameter  $b = \Omega(\log n / \log \log n)^2$ . Hence, they provide a solution to the DYNAMIC BIT VECTOR WITH INDELS problem. Their main structure is a weight-balanced B-tree (WBB) [Dietz 1989; Raman et al. 2001].

Our goal is to obtain  $nH_0 + o(n)$  bits of space and  $O(\log n)$  time for all the operations above. We build over a simplified version of their structure, which uses standard balanced trees and achieves  $O(\log n)$  time and  $O(n)$  bits of space [Chan et al. 2004]. This is described below.

Consider a balanced binary tree on  $A$  whose leftmost leaf contains bits  $a_1 a_2 \dots a_{\log n}$ , second left-most leaf contains bits  $a_{\log n+1} a_{\log n+2} \dots a_{2 \log n}$ , and so on. Each node  $v$  contains counters  $p(v)$  and  $r(v)$  telling the number of positions stored and the number of bits set in the subtree rooted at  $v$ , respectively. Note that this tree, with all its  $\log n$ -size pointers and counters, requires  $O(n)$  bits.

To perform  $\text{rank}(A, i)$ , we enter the tree to find the leaf containing position  $i$ . We start with  $\text{rank} \leftarrow 0$ . If  $p(\text{left}(v)) \geq i$  we enter the left subtree, otherwise we enter the right subtree with  $i \leftarrow i - p(\text{left}(v))$  and  $\text{rank} \leftarrow \text{rank} + r(\text{left}(v))$ . In  $O(\log n)$  time we reach the desired leaf and complete the rank query in  $O(\log n)$  time by scanning the bit sequence corresponding to that node. For *select* we proceed similarly, except that the roles of  $p()$  and  $r()$  are reversed. For  $\text{select}_0$  the computation is analogous.

Insertions and deletions are handled by entering to the correct leaf as in *rank*, and

<sup>5</sup>Note that  $o(n \log \sigma)$  is sublinear in the size of  $A$  measured in bits.

replacing its bit sequence with the new content. Then the  $p(v)$  and  $r(v)$  counters in the path from the leaf to the root are changed accordingly. When a leaf is updated to contain  $2 \log n$  bits, it is split into two leaves, each containing  $\log n$  bits. When a leaf is updated to contain  $(\log n)/2$  bits, it is merged with its sibling. If this merging produces a leaf with more than  $2 \log n - 1$  bits, this leaf is again split into two equal-size halves. After splitting and merging, the tree needs to be rebalanced and the counters updated in the nodes on the way to the root.

### 3.4 Dynamic Entropy-Bound Structures for Bit Vectors

Blandford and Blelloch [Blandford and Blelloch 2004] design a general scheme to convert a space-demanding data structure into one that requires  $O(nH_0 + \log n)$  bits of space. The data structures considered can solve a subset of a wide range of problems related to ordered sets, and the main idea is to represent such sets using *gap encoding* (see next). In particular, if one applies the idea to the structure described above for solving bit vector *rank*, *select*, *insert*, and *delete* in  $O(\log n)$  time, the only difference is that bit vectors are represented in compressed form in the leaves of the binary tree.

We note that gap encoding has also been used to achieve zero-order entropy in static schemes. In [Grossi and Vitter 2006] they explore the idea of inserting some information into the encoding so as to permit solving *rank* and *select* queries in logarithmic time via binary searches. In [Gupta et al. 2006a] they improve this result and reach time  $o((\log \log n)^2)$ , close to the lower bound on the predecessor problem when the space depends on the number of bits set and only logarithmically on the total number of bits. In [Mäkinen and Navarro 2007] they achieve constant time on gap encoding, yet they have a higher  $o(n)$ -type dependence on the total number of bits.

**3.4.1 Gap Encoding.** Let  $A = 0^{g_0}10^{g_1}1 \dots 0^{g_{\ell-1}}10^{g_\ell}$ , where  $0^{g_i}$  represents a sequence of  $g_i$  0-bits (called a gap). Gap encoding represents  $A$  as  $\delta(g_0)\delta(g_1) \dots \delta(g_{\ell-1})\delta(g_\ell)$ , where  $\delta(x)$  is an encoding for the nonnegative integer  $x$ . This encoding must satisfy two properties: (i)  $|\delta(x)| = \log x + o(\log x)$ ; (ii) we can univocally distinguish  $x$  and  $D$  from  $\delta(x)D$ , being  $D$  any bit sequence.

A well-known encoding satisfying the above properties is Elias'  $\delta$  [Elias 1975; Bell et al. 1990]. To represent  $x$ , let  $l = \lceil \log(x+1) \rceil$  be the number of bits necessary to encode  $x$ , and let  $ll = \lceil \log(l+1) \rceil$  be the number of bits necessary to code  $l$ . Then  $\delta(x)$  is formed by three parts: (a)  $ll$  0-bits followed by a 1-bit, (b) the  $ll - 1$  least significant bits of the binary representation of  $l$  (this part is empty if  $l < 2$ ), and (c) the  $l - 1$  least significant bits of the binary representation of  $x$  (this part is empty if  $x < 2$ ). For example, for  $x = 0$ , we have  $l = 0$  and  $ll = 0$ , thus  $\delta(0) = 1$ ; for  $x = 1$ ,  $l = 1$  and  $ll = 1$ , thus  $\delta(1) = 01$ ; for  $x = 2$ ,  $l = 2$  and  $ll = 2$ , thus  $\delta(2) = 001 0 0$ ; whereas  $\delta(3) = 001 0 1$ ;  $\delta(4) = 001 1 00$  since  $l = 3$  and  $ll = 2$ ; and so on.

It is clear that  $|\delta(x)| = \log x + 2 \log \log x + O(1) = \log x + o(\log x)$ . The total

ACM Journal Name, Vol. V, No. N, Month 20YY.

length of this representation of  $A$  is therefore

$$\begin{aligned} \sum_{i=0}^{\ell} \log g_i + o(\log g_i) &\leq \sum_{i=0}^{\ell} \log \frac{n-\ell}{\ell+1} + o(\log \frac{n-\ell}{\ell+1}) \\ &\leq (\ell+1) \log \frac{n}{\ell} + (\ell+1) o(\log \frac{n}{\ell}), \end{aligned}$$

where we have used the fact that  $\sum_{i=0}^{\ell} g_i = n - \ell$ , and thus the summation of those convex functions achieves its maximum value when all  $g_i = \frac{n-\ell}{\ell+1}$ . As the binary entropy of  $A$  is  $H_0 = \frac{\ell}{n} \log \frac{n}{\ell} + \frac{n-\ell}{n} \log \frac{n}{n-\ell}$ , the first term of the result is  $\ell \log \frac{n}{\ell} + O(\log n) \leq nH_0 + O(\log n)$ . The second term is  $O(\ell)$  if  $\ell = \Theta(n)$ , and  $o(nH_0 + \log n)$  otherwise. Therefore, the total size of the gap representation is  $n' = nH_0(1 + o(1)) + O(\ell + \log n)$  bits. For Elias' representation, this is more precisely  $n' = \ell \log \frac{n}{\ell} + O(\ell \log \log \frac{n}{\ell} + \log n)$ .

**3.4.2 A Dynamic Structure based on Gap Encoding.** Consider the balanced binary search tree of Section 3.3 built on the gap encoded bit vector  $A$ : The encoded bit vector  $\delta(g_0)\delta(g_1) \dots \delta(g_{\ell-1})\delta(g_{\ell})$  of length  $n'$  is partitioned into blocks of approximately  $\log n$  bits, each block containing as many full  $\delta(g_i)$  codes as can be accommodated into  $\log n$  bits. These blocks form the leaves of the binary search tree. To answer *rank* and *select* queries one can proceed just like in the uncompressed case, except that the final scanning in the leaves requires decoding the gap encoding. The time complexity remains  $O(\log n)$ .

To support *insert* and *delete* on this gap-encoded sequence one can proceed as in the uncompressed case, splitting and merging leaves when necessary. To insert a bit  $a$  preceding position  $i$  inside a block, we sequentially look for the gap where  $a$  should be inserted. Say that  $a$  must be inserted at relative position  $i'$  within  $0^{g_k}1$ ,  $1 \leq i' \leq g_k + 1$ . If  $a = 1$  we must replace  $\delta(g_k)$  by  $\delta(i' - 1)\delta(g_k - i' + 1)$ . Otherwise, if  $a = 0$  we must replace  $\delta(g_k)$  by  $\delta(g_k + 1)$ . All the  $\delta$ -codes that follow must be shifted to make room for the new code. The replacement and shifting can be easily done in  $O(\log n)$  time if there is enough empty space left in the block. If not, the block needs to be split into two. Notice that on a single insert the space needed can at most double. Deletions are handled analogously.

Space is improved from  $O(n)$  to  $O(nH_0 + \log n)$  bits. The gap encoding itself takes essentially  $nH_0$  bits, the unused empty space in the leaves occupies in the worst case other  $nH_0 + \log n$  bits, as blocks can be half-full<sup>6</sup>, and the tree pointers occupy  $\lceil nH_0 / \log n \rceil \times O(\log n) = O(nH_0 + \log n)$  bits.

### 3.5 Static Full-Text Self-Indexes

Many static full-text self-indexes are based on representing the Burrows-Wheeler transform [Burrows and Wheeler 1994] of a text using wavelet trees to support efficient substring searches. We will later consider dynamic wavelet trees to solve the DYNAMIC TEXT COLLECTION problem, hence we introduce the basic concepts here. We follow closely the description given in [Mäkinen and Navarro 2005].

<sup>6</sup>As described blocks could be 25% full, but it is easy to force them to be half-full.

3.5.1 *The Burrows-Wheeler Transform.* The *Burrows-Wheeler transform (BWT)* [Burrows and Wheeler 1994] of a text  $T$  produces a permutation of  $T$ , denoted by  $bwt(T) = T^{bwt}$ . We assume that  $T$  is terminated by an endmarker “\$”  $\in \Sigma$ , smaller than other symbols. String  $T^{bwt}$  is the result of the following transformation: (1) Form a *conceptual* matrix  $\mathcal{M}$  whose rows are the cyclic shifts of the string  $T$ , call  $F$  its first column and  $L$  its last column; (2) sort the rows of  $\mathcal{M}$  in lexicographic order; (3) the transformed text is  $T^{bwt} = L$ .

The BWT is reversible, that is, given  $T^{bwt}$  we can obtain  $T$ . Note the following properties [Burrows and Wheeler 1994]:

- Given the  $i$ -th row of  $\mathcal{M}$ , its last character  $L[i]$  precedes its first character  $F[i]$  in the original text  $T$ , that is,  $T = \dots L[i]F[i] \dots$
- Let  $L[i] = c$  and let  $r_i$  be the number of occurrences of  $c$  in  $L[1, i]$ . Let  $\mathcal{M}[j]$  be the  $r_i$ -th row of  $\mathcal{M}$  starting with  $c$ . Then the character corresponding to  $L[i]$  in the first column  $F$  is located at  $F[j]$  (this is called the *LF mapping*:  $LF(i) = j$ ). This is because the occurrences of character  $c$  are sorted both in  $F$  and  $L$  using the same criterion: by the text following the occurrences.

The BWT can then be reversed as follows:

- Compute the array  $C[1, \sigma]$  storing in  $C[c]$  the number of occurrences of characters  $\{\$, 1, \dots, c-1\}$  in the text  $T$ . Notice that  $C[c] + 1$  is the position of the first occurrence of  $c$  in  $F$  (if any).
- Define the *LF mapping* as follows:  $LF(i) = C[L[i]] + rank_{L[i]}(L, i)$ .
- Reconstruct  $T$  backwards as follows: set  $s = 1$  (since  $\mathcal{M}[1] = \$t_1t_2 \dots t_{n-1}$ ) and, for each  $i \in n-1, \dots, 1$  do  $T[i] \leftarrow L[s]$  and  $s \leftarrow LF[s]$ . Finally put the endmarker  $T[n] = \$$ .

The BWT transform by itself does not compress  $T$ , it just permutes its characters. However, this permutation is more compressible than the original  $T$ . Actually, it is not hard to compress  $T^{bwt}$  to  $O(nH_h + \sigma^{h+1} \log n)$  bits, for any  $h \geq 0$  [Manzini 2001]. The idea is as follows (we will reuse it in Section 7.3): Partition  $L$  into minimum number of pieces  $L^1L^2 \dots L^\ell$  such that the symbols inside each piece  $L^k = L[i_k, j_k]$  have the same  $h$ -context. Note that the  $h$ -context of a symbol  $L[i]$  is  $\mathcal{M}[i][1, h]$ . By the definition of  $h$ -th order entropy, it follows that  $|L^1|H_0(L^1) + |L^2|H_0(L^2) + \dots + |L^\ell|H_0(L^\ell) = nH_h$ . That is, if one is able to compress each piece up to its zero-order entropy, then the end result is  $h$ -th order entropy. Using, say, arithmetic coding on each piece, one achieves  $nH_h + \sigma^{h+1} \log n$  bits encoding of  $T$ . The latter term comes from the encoding of the symbol frequencies in each piece separately.

3.5.2 *Suffix Arrays.* The *suffix array*  $\mathcal{A}[1, n]$  of text  $T$  is an array of pointers to all the suffixes of  $T$  in lexicographic order. Since  $T$  is terminated by the endmarker “\$”, all lexicographic comparisons are well defined. The  $i$ -th entry of  $\mathcal{A}$  points to text suffix  $T[\mathcal{A}[i], n] = t_{\mathcal{A}[i]}t_{\mathcal{A}[i]+1} \dots t_n$ , and it holds  $T[\mathcal{A}[i], n] < T[\mathcal{A}[i+1], n]$  in lexicographic order.

Given the suffix array, the occurrences of the pattern  $P = p_1p_2 \dots p_m$  can be counted in  $O(m \log n)$  time. The occurrences form an interval  $\mathcal{A}[sp, ep]$  such that suffixes  $t_{\mathcal{A}[i]}t_{\mathcal{A}[i]+1} \dots t_n$ , for all  $sp \leq i \leq ep$ , contain the pattern  $P$  as a prefix. This

interval can be searched for using two binary searches in time  $O(m \log n)$ . Once the interval is obtained, a locating query is solved simply by listing all its pointers in constant time each.

We note that the suffix array  $\mathcal{A}$  is essentially the matrix  $\mathcal{M}$  of the BWT (Section 3.5.1), as sorting the cyclic shifts of  $T$  is the same as sorting its suffixes given the endmarker “\$”:  $\mathcal{A}[i] = j$  if and only if the  $i$ -th row of  $\mathcal{M}$  contains the string  $t_j t_{j+1} \dots t_{n-1} \$ t_1 \dots t_{j-1}$ .

**3.5.3 Backward Search.** The FM-index [Ferragina and Manzini 2000] is a self-index based on the Burrows-Wheeler transform. It solves counting queries by finding the interval of  $\mathcal{A}$  that contains the occurrences of pattern  $P$ . The FM-index uses the array  $C$  and function  $rank_c(L, i)$  of the  $LF$  mapping to perform backward search for the pattern. Fig. 1 shows the counting algorithm. Using the properties of the BWT, it is easy to see that the algorithm maintains the following invariant [Ferragina and Manzini 2000]: At the  $i$ -th phase, variables  $sp$  and  $ep$  point, respectively, to the first and last row of  $\mathcal{M}$  prefixed by  $P[i, m]$ . The correctness of the algorithm follows from this observation. Note that  $P$  is processed backwards, from  $p_m$  to  $p_1$ .

---

**Algorithm** FMCount( $P[1, m], L[1, n]$ )

- (1)  $i \leftarrow m$ ;
  - (2)  $sp \leftarrow 1$ ;  $ep \leftarrow n$ ;
  - (3) **while** ( $sp \leq ep$ ) **and** ( $i \geq 1$ ) **do**
  - (4)      $c \leftarrow P[i]$ ;
  - (5)      $sp \leftarrow C[c] + rank_c(L, sp - 1) + 1$ ;
  - (6)      $ep \leftarrow C[c] + rank_c(L, ep)$ ;
  - (7)      $i \leftarrow i - 1$ ;
  - (8) **if** ( $ep < sp$ ) **then return** 0 **else return**  $ep - sp + 1$ ;
- 

Fig. 1. FM-index algorithm for counting the number of occurrences of  $P[1, m]$  in  $T[1, n]$ .

Note that array  $C$  can be explicitly stored in little space, and for  $rank_c(L, i)$  we can directly use the wavelet tree as explained in Section 3.2. The space usage is  $nH_0 + o(n \log \sigma)$  bits and the  $m$  steps of backward search take overall  $O(m \log \sigma)$  time [Mäkinen and Navarro 2005].

Let us now consider how to locate the positions in  $\mathcal{A}[sp, ep]$ . The idea is that  $T$  is sampled at regular intervals, so that we explicitly store the positions in  $\mathcal{A}$  pointing to the sampled positions in  $T$  (note that the sampling is not regular in  $\mathcal{A}$ ). Hence, using the  $LF$  mapping, we move backward in  $T$  until finding a position that is known in  $\mathcal{A}$ . Then it is easy to infer our original text position. Fig. 2 shows the pseudocode.

We note that, in addition to  $C$  and  $rank$ , we need access to characters  $L[i']$  as well. These can be found using the same wavelet tree built for  $rank$ . Finally, if we sample one out of  $\log^{1+\varepsilon} n / \log \sigma$  positions in  $T$ , for any constant  $\varepsilon > 0$ , and use  $\log n$  bits to represent each corresponding  $\mathcal{A}$  value, we require  $O(n \log \sigma / \log^\varepsilon n) = o(n \log \sigma)$  additional bits of space and can locate each occurrence of  $P$  in  $O(\log^{1+\varepsilon} n)$  time.

---

**Algorithm** FMlocate( $i, L[1, n]$ )

- (1)  $i' \leftarrow i, t \leftarrow 0;$
- (2) **while**  $\mathcal{A}[i']$  is not known **do**
- (3)      $i' \leftarrow LF(i') = C[L[i']] + \text{rank}_{L[i']}(L, i');$
- (4)      $t \leftarrow t + 1;$
- (5) **return**  $\mathcal{A}[i'] + t;$

---

Fig. 2. FM-index algorithm for locating the occurrence  $\mathcal{A}[i]$  in  $T$ .

Finally, let us consider displaying text contexts. To retrieve  $T[l, r]$ , we start at the position in  $\mathcal{A}$  that points to the lowest marked text position following  $r$ . This position in  $\mathcal{A}$  is known from the sampling. From there, we perform  $O(\log^{1+\varepsilon} n / \log \sigma)$  steps, using the  $LF$  mapping, until reaching  $r$ . Then we perform  $\ell = r - l$  additional  $LF$  steps to collect the desired text characters. The resulting complexity is  $O(\ell \log \sigma + \log^{1+\varepsilon} n)$ .

All the  $O(\log \sigma)$  terms in the time complexities can be made  $O(\lceil \log \sigma / \log \log n \rceil)$  (which is constant if  $\sigma = O(\text{polylog}(n))$ ), by using multiary wavelet trees.

### 3.6 Dynamic Full-Text Self-Indexes

Chan, Hon, and Lam [Chan et al. 2004] show how to use a solution to DYNAMIC SEQUENCE WITH INDELS problem to obtain a solution to the DYNAMIC TEXT COLLECTION problem. One of the ideas is to simulate the above backward search algorithm: They use  $A = \text{bwt}(\mathcal{C})$  on a text collection  $\mathcal{C}$  (seen as a concatenation of texts over  $[0, \sigma]$ , where alphabet symbol 0 is reserved for separating contiguous texts and somehow plays the role of “\$”).

They show that one can dynamically maintain a collection of texts, by keeping a data structure supporting *rank*, *insert* and *delete* on  $A$ , in addition to a dynamic version of table  $C$  of Section 3.5. Adding new text  $T$  (preceded by 0) triggers  $|T| + 1$  insertions to  $A$ : The insertion points can be found in  $O(|T|g(|\mathcal{C}|))$  time, where  $g(n)$  is the time to access  $C$  and to answer *rank* on a collection of length  $n$ . The process consists of inserting the suffixes of  $T = t_1 \dots t_{|T|}$  one by one in backward fashion. The initial insertion point, for  $t_{|T|}$ , can be at  $C[t_{|T|}] + 1$  (note that, as there might be repeated suffixes, the position is not unique; this does not affect the correctness of the BWT-based scheme), thus we have to *insert* symbol  $t_{|T|-1}$  (preceding  $t_{|T|}$ ) at  $A[C[t_{|T|}] + 1]$ . In general, once we have inserted  $t_i$  at position  $j$  in  $A$  corresponding to  $T_{i+1, |T|}$ , we use the  $LF$  mapping to find the next insertion point (that is, the  $A$  position corresponding to  $T_{i, |T|}$ ). This is done with an access to  $C$  and a *rank* operation. At the end, the symbol 0 preceding  $t_1$  is inserted and we finally insert  $t_{|T|}$  at the position in  $A$  corresponding to suffix  $0 \cdot T$  ( $T$  is seen as a circular string).

Deleting a text  $T$  triggers  $|T| + 1$  deletions from  $A$ : The deletion points can be found in  $O(|T|g(|\mathcal{C}|) \log |\mathcal{C}|)$  time. In principle one would mimic the very same insertion process, this time doing deletions. The problem is that the operation receives the text  $T$  to delete but cannot know which of the 0’s of  $A$  correspond to it. Thus they search backwards for  $T$  in  $A$ , and assuming it is unique, they find the position in  $A$  corresponding to  $t_1$ . Then they have to find  $t_2$  and so on using the inverse of  $LF$ , which can be computed in  $O(\log |\mathcal{C}|)$  time.

If we call  $n = |\mathcal{C}|$ , their original structure (called **COUNT**) takes  $O(\sigma n)$  bits of space (it consists of one bitmap per alphabet symbol, marking its positions in  $A$ ), supports *rank* in  $g(n) = O(\log n)$  time, *insert* in  $O(\sigma \log n)$  time per symbol, and *delete* in  $O((\sigma + \log n) \log n)$  time per symbol. They give another structure that uses  $O(n \log \sigma)$  bits of space in exchange for supporting *rank* in  $O(\log^2 n)$  time. Insertion and deletion take  $O(\log^2 n)$  time per symbol. Using both structures in conjunction they still have  $O(\sigma n)$  bits of space and can handle both insertion and deletion in  $O(|T| \sigma \log n)$  time. Searches for  $P$  cost  $O(|P| \log n)$  time.

Their index is extended to support locating of occurrences using a structure called **MARK**, which samples one out of  $\log n$  collection positions and stores the position in  $A$  that corresponds to the sampled collection position. This requires  $O(\log n)$  bits per sample, for a total of  $O(n)$  further bits. To locate an occurrence at  $A[i]$ , they look for it in **MARK**. If it is not present, they use the *LF* mapping repeatedly (which traverses  $\mathcal{C}$  backwards) until a marked position is found. Then the original value  $A[i]$  is the sampled one plus the number of *LF* steps performed. This takes overall  $O(\log^2 n)$  time using either variant of their structure (the variant that moves backward in the text in  $O(\log^2 n)$  time per step can move forward in  $O(\log n)$  time, and this is equally good).

Finally, to recover a substring of some text  $T$  in the collection, the user must somehow know the lexicographic position of  $T$  within the other texts (this is a strong assumption!). The lexicographically  $j$ -th text has its character  $t_{|T|}$  at  $a_j$ .<sup>7</sup> By locating  $t_{|T|}$  in  $\mathcal{C}$  they can translate a relative position within a text into an absolute position in  $\mathcal{C}$ . Then they take the next sampled position in  $\mathcal{C}$  and use the *LF* mapping to discover the desired characters backward from there. This takes  $O(\log n(\ell + \log n))$  time, where  $\ell$  is the length of the text piece to display.

#### 4. DYNAMIC SUCCINCT STRUCTURES FOR BIT VECTORS AND PARTIAL SUMS

In this section we design a data structure to represent a bit sequence  $A = a_1 \dots a_n$  using  $n + o(n)$  bits of space and performing operations *rank*, *select*, *insert* and *delete* all in  $O(\log n)$  time. This already improves previous results [Hon et al. 2003b], and serves as a basis for the entropy-bound structures developed in the next sections.

##### 4.1 High-Level Hierarchy

Section 3.3 shows how to obtain the desired time complexities using  $O(n)$  bits of space. To achieve  $n + o(n)$  bits, we use the same tree organization, except that it is built on  $\omega(\log n)$ -size leaves. Thus the tree has  $o(n/\log n)$  internal nodes which, with all their pointers, require only  $o(n)$  bits of space. The two problems to solve are (i) we cannot process the leaves bitwise in  $O(\log n)$  time; (ii) we cannot store  $(1 + \epsilon)s$  bits for leaves and use only  $s$  bits from those. This is essentially the technique used in Section 3.3 to ensure that bit insertions/deletions are handled

<sup>7</sup>To see this, note that  $t_{|T|}$  in  $a_j = L[j]$  of the BWT corresponds to  $F[j] = 0$ . The 0's of all the texts are at the beginning of  $F$ , sorted lexicographically by the texts. This observation, plus the assumption that one knows the lexicographic position of the text of interest, is sufficient to remove the  $O(\log n)$  factor for deletions: To delete the  $j$ -th text, start the deletion from  $a_j$  and go on with backward steps until returning to the original position.

with  $O(1)$  tree node updates.

We divide  $A$  into blocks and superblocks, as in Section 3.1. Each superblock  $S$  will maintain  $s = f(n) \log n$  bits (for some  $f(n) = O(\log n)$  to be determined later), and will be stored in a tree leaf, without any extra bits of space. Each superblock will hold exactly  $2f(n)$  whole blocks of  $t = (\log n)/2$  bits each. All the  $s$  bits of the superblock will thus be stored contiguously in plain form, without any extra structure in principle (later we add a few pointers per leaf). From now on we will use the term “leaf” and “superblock” interchangeably.

#### 4.2 Queries Inside a Superblock

A  $\text{rank}(S, i)$  query inside a superblock is handled in  $O(\log n)$  time by using a universal table  $R$ , which receives a  $t$ -bit sequence and gives the total number of 1-bits in it. Thus we traverse the superblock in a blockwise manner, adding 1-bits until we reach the block that contains position  $i$ . Within that block we count the 1-bits one by one until reaching the  $i$ -th position. The whole process takes  $O(f(n) + t) = O(\log n)$  time. More formally, if  $a_1 \dots a_s$  are the bits of the leaf, then  $b = 1 + \lfloor (i - 1)/t \rfloor$  is the block  $i$  belongs to, and we compute

$$\sum_{1 \leq q < b} R[a_{(q-1)t+1} \dots a_{qt}] + \sum_{(b-1)t+1 \leq q \leq i} a_q.$$

A  $\text{select}(S, j)$  query is solved similarly: We add up successive  $R$  values until we exceed  $j$  at some block  $b$ . Then we rescan block  $b$  bitwise until reading the  $j'$ -th 1-bit, where  $j' = j - \sum_{1 \leq q < b} R[a_{(q-1)t+1} \dots a_{qt}]$ . This also takes  $O(f(n) + \log n)$  time. We remark that table  $R$  is universal and does not depend on the sequence  $A$  nor on the particular leaf; just on  $t$ . Moreover, table  $R$  is very small, requiring just  $O(\sqrt{n} \log \log n)$  bits.

#### 4.3 Updates Inside a Superblock

To insert a bit  $a$  at position  $i$  of  $S$ , we simply shift to the right the bit vector  $a_i \dots a_s$ , to make room for  $a_i = a$ . This shift can be done by chunks of  $\Theta(\log n)$  bits under a RAM model that permits bit shifts, or multiplication and division (as, say, dividing by two is equivalent to shifting the bits to the right). Otherwise, it is easy to build a small universal table to perform the shifts by chunks of  $t$  bits. Thus the new bit is accommodated within the leaf in  $O(f(n))$  time.

We note that, after the shift, former bit  $a_s$  overflows to the next leaf. We now insert  $a_s$  at the beginning of the next leaf, which causes a new overflow, and so on. This brings two problems: (i) doing the propagation efficiently within those next leaves, and (ii) limiting the propagation to a reasonable number of leaves.

To achieve efficiency within each leaf, we redefine them as circular arrays of bits. A pointer telling the position of the first bit in the leaf is stored within the leaf, and it requires  $O(\log \log n)$  bits. This amounts to  $O(n \log \log n / (f(n) \log n)) = o(n/f(n))$  wasted space. The advantage is that, if leaves are circular arrays, then inserting the overflow bit at their beginning and taking out their last bit that overflows to the next leaf is easily done in constant time.

To limit the propagation of overflows across leaves, every  $f(n)$  leaves we permit the formation of a *partial* leaf, which reserves  $f(n) \log n$  bits but might be partially full. Those partial leaves amount to  $n/f(n)$  wasted bits overall. Partial leaves are



not circular, so the pointer of the previous paragraph can be reused to store their current number of bits. An additional bit, stored for each leaf, tells whether it is full or partial.

Partial leaves ensure that we never traverse more than  $f(n)$  leaves in the overflow propagation process. Thus the overall insertion time (considering just the work within leaves) is  $O(f(n))$  in the leaf that receives the insertion, plus  $O(1)$  per leaf to propagate the overflow across  $O(f(n))$  full leaves, plus  $O(f(n))$  to insert the overflowed bit in the partial leaf (as it is not circular, it is necessary to shift its values). This adds up  $O(f(n))$ .

To ensure the desired density of partial leaves, we first check whether there is a partial leaf among the next  $2f(n)$  leaves. If there is one, we carry out the propagation up to it. Otherwise, we propagate  $f(n)$  leaves and create a new empty partial leaf. In both cases we work over  $O(f(n))$  leaves, and guarantee that every partial leaf is  $f(n)$  leaves away from any other. We note that partial leaves may end up overflowing, at which point they are not considered partial anymore.

For deletions we proceed similarly, bringing back the first bit of the next leaf and propagating the underflow. If there is a partial leaf within the next  $2f(n)$  leaves, we propagate the underflow until there. Otherwise, we propagate the underflow for  $f(n)$  leaves, and declare the  $f(n)$ -th leaf partial. A partial leaf that gets empty should be removed. Thus deletions are also handled in  $O(f(n))$  time.

#### 4.4 Operations in the Tree

Section 3.3 shows how to perform *rank* and *select* up to the leaves in  $O(\log n)$  time. We have shown in Section 4.2 how to manage inside the leaves in time  $O(f(n) + \log n)$ . The wasted space is  $O(n/f(n))$  across the leaves, and also  $O(n/f(n))$  for the internal tree nodes. By choosing  $f(n) = \log n$  we obtain  $O(\log n)$  time for both query operations, and  $O(n/\log n) = o(n)$  wasted bits of space. This completes the solution for those queries.

The update time within leaves is  $O(\log n + f(n)) = O(\log n)$  according to Section 4.3. Let us now consider the tree adjustments required upon updates.

Inserting and deleting bits requires rewriting the  $p()$  and  $r()$  values from the affected superblock(s) through the root. Creation and deletion of leaves and internal tree nodes is easily handled together with the maintenance of  $r()$  and  $p()$ . We note, however, that although each bit insertion/deletion can produce at most one tree leaf insertion/deletion, it can affect the bits of  $O(f(n))$  leaves, as well as all their  $r()$  and  $p()$  values upwards the root. Yet, we note that those  $O(f(n))$  leaves are contiguous in the tree and therefore their total number of ancestors do not exceed  $O(f(n)) + O(\log n) = O(\log n)$ . It is not hard to organize those updates so as to work  $O(\log n)$  time overall. Thus we achieve  $O(\log n)$  overall time complexity for insertions and deletions as well.

#### 4.5 Changing $\log n$

Our result so far assumes that  $\log n$  stays constant during the operations. This value fixes the superblock/block hierarchy and the global preprocessed tables. This assumption can be removed in two ways: (1) performing a global rebuild whenever  $\log n$  changes; (2) maintaining partial structures ready for values  $(\log n) - 1$ ,  $\log n$ , and  $(\log n) + 1$  (which we call the *previous*, *current*, and *next*).

Approach (1) is easy to implement. We can rebuild all structures in  $O(n)$  time when necessary to accommodate the new value of  $\log n$ . Amortized over all insertions and deletions, this costs only  $O(1)$  time per operation.

Approach (2) is more complex but is inspired on a standard mechanism to convert amortized complexity into worst-case complexity. The idea is to split the current elements among the *previous*, *current*, and *next* structures, so that the first elements are in *previous*, the last are in *next*, and *current* holds the middle elements. It is trivial to run *rank* and *select* queries on this split structure. When  $n$  is zero or a power of 2, all the elements are in *current*, and the other two are empty. Upon an insertion, the size of *next* must grow by 2 unless it is already full, and *previous* must shrink by 1 unless it is already empty; a deletion must cause the opposite effect; and *current* acts as a variable-size buffer.

To achieve this, let us denote  $x \rightarrow y$  or  $x \leftarrow y$  the movement of one element among structures, for  $x, y \in \{p, c, n\}$ , e.g.  $p \leftarrow c$  means moving the first element of *current* to *previous*. If the source structure is empty, the movement is just ignored. Then, we insert (delete) in the proper structure and then, depending on where the insertion (deletion) point lies, we move elements as follows:

—*previous*:  $p \rightarrow c, p \rightarrow c, c \rightarrow n, c \rightarrow n$  ( $c \leftarrow n, c \leftarrow n, p \leftarrow c, p \leftarrow c$ ).

—*current*:  $p \rightarrow c, c \rightarrow n, c \rightarrow n$  ( $c \leftarrow n, c \leftarrow n, p \leftarrow c$ ).

—*next*:  $p \rightarrow c, c \rightarrow n$  ( $c \leftarrow n, p \leftarrow c$ ).

It is easy to see by inspection that, no matter where an insertion (deletion) lies, *previous* will shrink (grow) by one element, and *next* will grow (shrink) by two elements. Shrinking is impossible iff the set is already empty, and, more importantly, growing is impossible iff the other sets are already empty. We recall that we start with all the  $n$  elements in *current* and both *previous* and *next* empty.

Let us now consider a mixed sequence of insertions and deletions. As long as growing is possible for both *previous* and *next*, we have that *next* will hold  $2r$  elements if the number of insertions minus deletions since *current* held all the elements is  $r \geq 0$ . Similarly, *previous* will hold  $r$  elements if the number of deletions minus insertions is  $r \geq 0$ . If, at some point, *next* cannot grow, this means that *previous* and *current* are empty, and thus *next* contains the  $n$  elements originally in *current* plus the  $r$  new elements inserted, thus its size is  $n + r$ . On the other hand, we know that its size must be  $2r$  as there have been  $r$  net insertions. From the equality  $n + r = 2r$  we get that, after  $r = n$  net insertions, *next* will hold all the  $2n$  elements. At this point it becomes *current* and the new *previous* and *next* structures are empty. Similarly, after  $r$  net deletions so that *previous* contains all the (remaining)  $n - r$  elements, it must hold  $n - r = r$  and thus we have that, after  $r = n/2$  deletions, all the elements are in *previous*, which becomes the new *current*.

At those points, precisely,  $n$  is a new power of 2 and  $\log n$  has changed its value. The space requirement is still  $n + o(n)$  bits, even when the tree pointers of *next* require  $\log(2n)$  bits.

Note, however, that we do not have time to build the right  $R$  table corresponding to the new *next* or *previous* structures, as we need it immediately. The solution is to maintain tables  $R$  ready for 5 values of  $\log n$ , from  $(\log n) - 2$  to  $(\log n) + 2$ . Thus, upon a change in  $\log n$ , we have the new  $R$  table we need immediately available.

For example, if  $\log n$  increases, we have ready the 3 tables for  $\log n$ ,  $(\log n) + 1$  and  $(\log n) + 2$ . Yet, to maintain the invariant, we should now build the  $R$  table for  $(\log n) + 3$ , but we have plenty of time to build it in parallel with the new operations: After  $O(\sqrt{n})$  operations we have managed to build the new  $R$  table (as there are  $O(\sqrt{n})$  cells and each is easily built in  $O(\log n)$  time). This is much less than the time necessary for  $\log n$  to change again. If, on the other hand,  $\log n$  shrinks back before we build the new  $R$  table, we can just discard the partial work done.

Finally, note that the trees have only two types of nodes (internal and leaves), so memory can be easily managed in constant time per allocation or deallocation request. We obtain the following result, where we recall that  $w$  is the number of bits in the word of the RAM model.

**THEOREM 1.** *The DYNAMIC BIT VECTOR WITH INDELS problem under RAM model with constraint  $\log n = \Theta(w)$  can be solved using  $n + O(n/\log n)$  bits of space supporting the operations rank, select, insert, and delete, in  $O(\log n)$  worst-case time.*

#### 4.6 Handling the Case $\log n = o(w)$

In Section 4.5 we assumed that  $\log n = \Theta(w)$ . This permits us using the system memory, with  $w$ -bit pointers, and assume that each such pointer takes  $O(\log n)$  bits. This is a common assumption in the literature, but we find it too restrictive in the dynamic setting. In this section we extend our result to the case where  $\log n = o(w)$ . We use different solutions for the internal tree nodes and for the leaves.

Let us start with the tree nodes. Assume our tree has  $u$  internal nodes and holds  $n$  bits. Each time tree *next* becomes *current* in Section 4.5, we allocate  $u' = 2n/(f(2n)\log(2n))$  cells of space for the new *next* tree (those cells require  $(\log u')$ -bit pointers). This ensures that, when the number of bits we handle reaches  $2n$ , all the tree will be in *next*, *next* will become *current*, and thus the new *current* memory area will be completely full. Similarly, when *previous* becomes *current*, we allocate space for  $u'' = (n/2)/(f(n/2)\log(n/2))$  cells in the new *previous* tree. Note also that, according to the rules to move elements in Section 4.5, *current* never increases its size after it is created. Thus, there is no flaw in creating it with exactly the number of cells to hold its current number of elements. Overall, we are using  $O(n/f(n))$  bits of space (more precisely, about  $3.5n/f(n)$  bits) for all the trees, plus  $3w$  bits necessary for the pointers to the three memory areas.

Within each memory area, we must provide a mechanism to manage tree node allocation and freeing in constant time. This is very easy because all the nodes have the same size within each area. An implicit list of free cells, where the list pointers use the same free space they mark, is sufficient for memory management within each area. The above paragraph shows that overflows do not occur within these areas.

We must use a different mechanism for the tree leaves, as they add up  $n$  bits and thus we cannot afford allocating  $3.5n$  bits of space. Let us focus in one individual structure (say, *current*). Each leaf takes  $f(n)\log n$  bits of space. A separate memory area will store all those leaves in array form, so an  $O(\log n)$ -bits index into the array suffices as a pointer from an internal tree node to a leaf. This array must be kept in

compact form upon insertions and deletions, which is easily achieved by storing the parent of each leaf in the array area (this permits moving the last leaf to the hole left by a deletion and updating the tree node that points to that leaf). Therefore allocating or freeing a leaf requires  $O(f(n))$  time on a RAM machine.

The problem is how to allocate memory when this array grows. We divide the whole memory for the leaves into  $\sqrt{n/w}$  chunks of  $\sqrt{nw}$  bits each<sup>8</sup>. All chunks will be full except for the last one, costing us  $O(\sqrt{nw})$  extra bits. As they are all of the same size, it is easy to allocate and deallocate chunks from system memory in constant time. Any pointer to a leaf is composed of two parts: the first  $\frac{1}{2} \log \frac{n}{w}$  bits indicate the chunk number, and the other  $\frac{1}{2} \log(nw)$  bits give the offset within the chunk. To achieve constant-time access from a pointer, we need a global array of the chunk addresses in system memory, requiring  $O(w\sqrt{n/w}) = O(\sqrt{nw})$  bits. The overall extra space for the chunk mechanism is  $O(\sqrt{nw})$ , which is  $o(n)$  unless  $n$  is very small,  $n = O(w)$ . In this case  $\sqrt{nw} = O(w)$ , and as we are already paying  $O(w)$  bits for a constant number of system memory pointers, we can spend the same space as if we had  $w$  bits in our structure (using a single chunk).

Thus we can reexpress the result of Theorem 1 under our weaker model of computation as follows:

**THEOREM 2.** *The DYNAMIC BIT VECTOR WITH INDELS problem under RAM model with constraint  $\log n = O(w)$  can be solved using  $n + O(n/\log n + \sqrt{nw} + w) = n + o(n) + O(w)$  bits of space, supporting the operations *rank*, *select*, *insert*, and *delete*, in  $O(\log n)$  worst-case time.*

We note that the only price we are finally paying is  $O(w)$  extra bits of space. This is asymptotically optimal if we assume that at least it is necessary to have a single system pointer to any dynamic structure.

#### 4.7 Searchable Partial Sums with Indels

Assume now that our sequence  $A = a_1 \dots a_n$  is formed by  $k$ -bit numbers. Now the length of  $A$  in bits is  $kn$ . Over this sequence we wish to support the operations *sum*, *search*, *insert*, and *delete* (*update* can be obtained with *delete* and *insert*).

We can apply essentially the same technique as for bits (case  $k = 1$ ), by choosing blocks of  $t = (\log n)/(2k)$  numbers, so that we can handle  $(\log n)/2$  bits in constant time (let us for now assume  $k \leq (\log n)/2$ ). We maintain  $s = f(n) \log(n)/k$  numbers (that is,  $f(n) \log(n)$  bits) in a superblock, so we can still handle it in  $O(f(n))$  time. Now table  $R$  receives a sequence of  $t$  numbers and delivers their sum. Table  $R$  still requires  $O(\sqrt{n} \text{ polylog}(n))$  bits. The process to compute *sum* within a leaf is completely analogous to *rank* in Section 4.2: we add up  $R$  values until reaching the block, and then add up  $O(\log(n)/k)$   $k$ -bit numbers. Likewise, *search* is analogous to *select*. Both operations are carried out in  $O(f(n) + \log n)$  time within a superblock.

For  $k > (\log n)/2$  (but still  $k = \Theta(\log n)$ ), we do not use table  $R$ , but simply add up the numbers in the leaf one by one. This still solves *sum* and *search* in  $O(f(n))$  time.

<sup>8</sup>This partitioning does not change along the lifetime of the tree under consideration, that is, we have to interpret  $n$  as  $2^{\lceil \log n \rceil}$  here.

The updates and the tree work exactly as for bits, using circular arrays. For the tree, the  $r()$  values in the nodes are now the sums of the numbers in the node subtree, and all the management is exactly analogous. We thus achieve  $O(\log n + f(n))$  time for the operations, and  $O(kn/f(n))$  bits of extra space. We can choose  $f(n) = \log n$  as before. (Note that the  $r()$  values also fit in  $O(\log n)$  bits, as we have  $n$  numbers of  $k$  bits, which add up at most  $n2^k$ , and  $\log(n2^k) = k + \log n = O(\log n)$  bits.)

Changing  $\log n$  can be handled smoothly, just as when dealing with bits. Note that, when  $\log n$  changes, we may switch from using  $R$  to not using it, or vice versa. If we assume that  $\log n = \Theta(w)$ , we have the following result.

**THEOREM 3.** *The SEARCHABLE PARTIAL SUMS WITH INDELS problem under RAM model with constraint  $\log n = \Theta(w)$  with  $k$ -bit numbers,  $k = O(\log n)$ , can be solved using  $kn + o(kn)$  bits of space, supporting the operations *sum*, *search*, *insert*, and *delete*, in  $O(\log n)$  worst-case time.*

The best current result permitting indels [Hon et al. 2003b] works only for  $k = O(1)$ , and also requires  $kn + o(kn)$  bits of space. It needs  $O(\log_b n)$  time for *sum* and *search*, and  $O(b)$  amortized time for *insert* and *delete*, for  $b = \Omega(\log n / \log \log n)^2$ . Within the minimum update time they can achieve  $O(\log n / \log \log n)$  time for the queries. Our result achieves slightly worse complexity for queries and better complexity for updates. In addition, all of our complexities are worst-case and we can handle any  $k = O(\log n)$  value.

A final point is to consider the case  $\log n = O(w)$ . We use the same memory arrangement of Section 4.6 to achieve the same time complexities and only  $O(w)$  extra space (the same analysis applies verbatim with  $kn$  instead of  $n$  bits). Yet, we could like to handle  $k$  values as large as  $k = \Theta(w)$ . In this case the tree requires  $O(w \cdot kn / (f(n) \log n))$  bits of space for the  $r()$  values (the  $r$  values are now sums of  $O(n)$   $w$ -bit numbers, and thus they require  $O(w + \log n) = O(w)$  bits). We must choose  $f(n) = \omega(w / \log n)$  to achieve sublinear extra space. Let us choose  $f(n) = w / \log^{1-\varepsilon} n$ , for any constant  $\varepsilon > 0$ . The time complexity raises to  $O(w / \log^{1-\varepsilon} n + \log n)$ .

**THEOREM 4.** *The SEARCHABLE PARTIAL SUMS WITH INDELS problem under RAM model with constraint  $\log n = O(w)$  with  $k$ -bit numbers,  $k = O(w)$ , can be solved using  $kn + o(kn) + O(w)$  bits of space, supporting the operations *sum*, *search*, *insert*, and *delete*, in  $O(w / \log^{1-\varepsilon} n + \log n)$  worst-case time, for any constant  $\varepsilon > 0$ .*

## 5. DYNAMIC ENTROPY-BOUND STRUCTURES FOR BIT VECTORS

We design two data structures to represent a bit sequence  $A = a_1 \dots a_n$  of binary zero-order entropy  $H_0$ , using essentially  $nH_0$  bits of space and performing operations *rank*, *select*, *insert* and *delete* in  $O(\log n)$  time.

Our two solutions differ in the extra space they achieve on top of the  $nH_0$  bits. Their common parts are as follows. Both use balanced trees with leaves reserving  $s = f(n) \log n$  bits of space, where  $O(\log n)$  bits can be wasted within each leaf. Both share the  $O(\log n)$ -time mechanism for the operations in the tree, differing only in how they manage within the leaves. Both use the mechanism of propagation

to partial leaves to ensure that most leaves are almost full, and at most one out of  $f(n)$  leaves is partial, so  $O(f(n))$  leaves are affected by an insertion or a deletion. Both use precomputed tables to process, in constant time,  $\Theta(\log n)$  bits within leaves.

We note that, since now the sequence is not directly available, we must provide a way to retrieve any bit  $a_i$  from  $A$ . In a binary sequence this is easy, as  $a_i = \text{rank}(A, i) - \text{rank}(A, i - 1)$ , so we can do it also in  $O(\log n)$  time. Actually, in our second solution, we can retrieve an  $O(\log^2 n)$ -bit chunk from  $A$  within the same  $O(\log n)$  time.

## 5.1 Gap Encoding

The first mechanism is suitable for sequences where 1-bits are sparse (the complementary technique can be used when 0-bits are sparse). Let  $\ell$  be the number of 1-bits in  $A$ , then the space we require with this method is essentially  $\ell \log \frac{n}{\ell}(1 + o(1)) + O(\ell) \leq nH_0(1 + o(1)) + O(\ell)$  bits.

Recall from Section 3.4 the structure by Blandford and Blelloch [Blandford and Blelloch 2004]. We show how to improve to 1 the constant multiplying the entropy term of their space requirement.

**5.1.1 Operations Inside a Superblock.** We maintain as many complete gaps ( $\delta$ -codes) as possible within each leaf of  $s$  bits, and assume that the 1-bit after the last encoded gap belongs to the leaf (thus the last gap of  $A$  requires special treatment). This representation may leave up to  $\log n + O(\log \log n)$  unused bits at the end of the leaf because the next  $\delta$ -code does not fit in it, and thus it is written at the next leaf. We also need  $O(\log \log n)$  bits to record the number of bits used in the leaf, as well as to mark the beginning of the circular array. We do not use the concept of block in this solution, just superblocks (that is, leaves) formed by gaps.

Note that a leaf of  $s = f(n) \log n$  bits may represent as many as  $\Theta(s)$  gaps, and as few as  $f(n)(1 + o(1))$ . In order to process a whole leaf in  $O(f(n))$  time, we need universal tables that let us process it by chunks of  $\Theta(\log n)$  bits. Let  $G$  be a table that receives  $t = (\log n)/2$  bits as follows:  $G(x) = (b, r, l)$  indicates that it is possible to decode the first  $l$  bits of  $x$  (that is, the final  $t - l$  bits do not make up a complete  $\delta$ -code), and that in those  $l$  bits there are  $r$  gaps that add up  $b$ . Note that it is possible that a  $\delta$ -code is longer than  $t$  bits. If  $x$  is the prefix of such a  $\delta$ -code, then  $G(x) = (0, 0, 0)$  indicates that  $G$  is unable to decode it. Table  $G$  requires  $O(\sqrt{n} \log n)$  bits of space.

To decode  $\delta$ -codes longer than  $t$  bits we use a different table which decodes only parts ( $a$ ) and ( $b$ ) of the  $\delta$ -code (see Section 3.4.1).  $U(x) = (d, l)$  means that the first  $\delta$ -code represented in  $x$  (or of which  $x$  is a prefix) represents a number of  $d$  bits, and that parts ( $a$ ) and ( $b$ ) of its representation require  $l$  bits. A further access for the next  $d \leq \log n$  bits (once we skip the first  $l$  bits of  $x$ ) completes the decoding of the long gap. Table  $U$  handles entries of  $O(\log \log n)$  bits, and thus it needs  $O(\text{polylog}(n))$  bits of space. Using  $G$  and  $U$  we can, in a constant number of accesses, decode at least one gap and at least  $t = \Theta(\log n)$  bits from the leaf.

A  $\text{rank}(S, i)$  query inside a leaf is handled in  $O(f(n) + \log n)$  time by decoding successive gaps using  $G$  (and occasionally  $U$ ) and adding the  $r$  values (that is, 1-bits), as long as the sum of the  $b + 1$  values (that is, gap lengths plus their

terminating 1-bit) does not exceed  $i$ . The  $l$  values delivered by  $G$  are used to advance in the  $\delta$ -encoded sequence. Once the next  $G$  access exceeds  $i$ , we reread those bits code-wise, using  $U$  (even for short codes) and adding 1 per gap to the result, until we read the gap where position  $i$  is exceeded. Overall we spend  $O(f(n))$  time with  $G$  and  $O(\log n)$  time with  $U$ .

Similarly  $select(S, j)$  is solved by adding values  $b + 1$  until the sum of the  $r$  values exceeds  $j$ , and then rereading the last argument of  $G$  code-wise until  $j$  is reached. For  $select_0(S, j)$  we must add values  $b + 1$  until the sum of the  $b$  values exceeds  $j$ , and the rest is straightforward.

To insert a bit  $a$  preceding position  $i$  in  $S$ , we sequentially look for the gap where  $a$  should be inserted (using  $G$  and  $U$  as before). Say that  $a$  must be inserted at relative position  $i'$  within  $0^{g_k}1$ ,  $1 \leq i' \leq g_k + 1$ . If  $a = 1$  we must replace  $\delta(g_k)$  by  $\delta(i' - 1)\delta(g_k - i' + 1)$  (see Section 3.4.2). Otherwise, if  $a = 0$  we must replace  $\delta(g_k)$  by  $\delta(g_k + 1)$ . All the  $\delta$ -codes that follow must be shifted to make room for the new code. The replacement can be easily done in  $O(\log n)$  time and the shifting can be carried out in  $O(f(n))$  time as in Section 4.2. Deletion of a bit is analogous.

We note that insertion of a new bit can expand the code sequence within the leaf by  $O(\log n)$  bits, which may overflow and require that (other)  $O(\log n)$  bits formed by whole overflowing codes be moved to the next leaf. This propagation is identical to that of Section 4. The fact that we move  $O(\log n)$  bits instead of one bit changes nothing under the RAM model: To copy  $O(\log n)$  bits from the previous leaf, one first makes room for them by taking out  $O(\log n)$  bits from the end of the circular array; then the desired bits are copied just before the beginning of the circular array; and the bits that were moved out overflow to the next leaf. All this is handled in a constant amount of  $\Theta(\log n)$ -bit moves. Thus the insertion of a bit is handled in  $O(\log n + f(n))$  time. Deletion is analogous.

There is, however, the problem of accessing in constant time the last whole codewords of a leaf (those occupying the last  $O(\log n)$  bits), as  $\delta$ -codes cannot be read backwards. There are several possible solutions to this problem. Probably the shortest to describe is that we can use an alternative gap encoding that modifies that of Section 3.4.1. Recall parts (a), (b) and (c) of the classical  $\delta(x)$  encoding. In our encoding,  $\delta'(x)$ , we represent parts (a), (b), (c), then (b) reversed, and finally (a) reversed. It is easy to see that  $\delta'(x)$  can be read in either direction in constant time, and that it requires  $|\delta'(x)| = \log x + 4 \log \log x + O(1) = \log x + o(\log x)$  bits of space. All the asymptotic analysis remains just as when using the original  $\delta(x)$  codes, and now the last whole codewords of a leaf can be easily identified.

**5.1.2 Changing  $\log n$  and Handling the Case  $\log n = o(w)$ .** The case of varying  $\log n$  can be treated analogously to Sections 4.5 and 4.6. We must have tables  $G$  and  $U$  ready for 5 values of  $\log n$ , from  $(\log n) - 2$  to  $(\log n) + 2$ , and have plenty of time to build them for the next change of  $\log n$ . The way to handle the case  $\log n = o(w)$  is analogous too.

We note that the extra space of the tree and partially full leaves adds up  $n'/f(n)$ , not  $n/f(n)$  (recall that  $n'$  is the number of bits in the encoded sequence). Also, the  $\sqrt{nw}$  space complexity of Section 4.6 is actually  $\sqrt{n'w}$ . By choosing again  $f(n) = \log n$  we achieve the following result.

**THEOREM 5.** *The DYNAMIC BIT VECTOR WITH INDELS problem under RAM*

model with constraint  $\log n = O(w)$  can be solved using  $nH_0(1 + o(1)) + O(\ell + \sqrt{n} \text{polylog}(n) + w) = nH_0 + o(n) + O(\ell + w)$  bits of space supporting the operations rank, select, insert, and delete, in  $O(\log n)$  worst-case time. Here,  $H_0 \leq 1$  is the empirical zero-order entropy of the sequence and  $\ell$  the number of bits set.

To compare the extra space against the (static) structure of [Gupta et al. 2006a], we rewrite the part related to the sequence representation,  $n' = nH_0 + O(\ell)$ , into the more precise form  $n' = \ell \log \frac{n}{\ell} + O(\ell \log \log \frac{n}{\ell})$ . Our  $1 + o(1)$  is actually  $1 + O(1/\log n)$ , so the product of both is still  $\ell \log \frac{n}{\ell} + O(\ell \log \log \frac{n}{\ell})$ , just as in [Gupta et al. 2006a]. In addition we have a dependence on the uncompressed stream size, yet this is mild,  $O(\sqrt{n} \text{polylog}(n))$ . In Section 5.2 this dependence becomes stronger, but in exchange the extra space  $O(\ell)$  is removed. This is relevant if  $\ell = \Theta(n)$ .

**5.1.3 Searchable Partial Sums Revisited.** Consider again the problem of managing a sequence  $A$  of  $k$ -bit positive numbers  $a_i$ ,  $1 \leq a_i \leq 2^k$ . Assume we represent it as a binary sequence  $A'$  of  $\sum_{i=1}^n a_i$  bits. In  $A'$  we set bits  $\text{sum}(A, i)$  for all  $i$ . Then, it holds  $\text{sum}(A, i) = \text{select}(A', i)$  and  $\text{search}(A, j) = 1 + \text{rank}(A', j - 1)$ . Sequence  $A'$  will be represented using gap encoding. Inserting a number  $a$  into  $A$  is equivalent to inserting a whole gap  $\delta(a - 1)$  into  $A'$ . This can be done in a form completely analogous as how we inserted individual bits (and even slightly simpler). The same holds for deletions.

Thus all the operations are supported in  $O(\log n)$  time. As for the space, calling  $n' = \sum_{i=1}^n \log a_i \leq kn$ , the number of bits in  $A'$  is upper bounded by  $n' + o(n') + O(n)$ . The following result is immediate (for the details related to  $w$  recall Section 4.7).

**THEOREM 6.** *The SEARCHABLE PARTIAL SUMS WITH INDELS problem under RAM model with constraint  $\log n = O(w)$  with  $k$ -bit positive numbers can be solved using  $n' + o(n') + O(n + w)$  bits of space, where  $n' \leq kn$  adds up the exact number of bits needed to represent each number in the sequence. This representation supports the operations sum, search, insert, and delete, in  $O(w/\log^{1-\varepsilon} n + \log n)$  worst-case time for any constant  $\varepsilon > 0$ . In particular, this is  $O(\log n)$  if  $\log n = \Theta(w)$ .*

Note that this result is similar to that of Theorem 4 if all the numbers are at least  $2^{k-1}$ . Yet, when there are small and large numbers together, this theorem achieves a more compact representation.

## 5.2 Block Identifier Encoding

The second mechanism to compress bit sequences is slightly more complex, yet it removes the  $O(\ell)$  term from the space complexity. This is important when the sequence is sufficiently dense of 1-bits.

The solution in this section uses a scheme close to the one described in Section 3.1, albeit simplified because we do not need to achieve constant time within a leaf. We divide  $A$  into blocks and superblocks, where superblocks (the tree leaves) reserve  $s = f(n) \log n$  bits of space and maintain as many complete blocks as possible. Each block represents  $t = (\log n)/2$  bits, but it is stored in fewer bits using its  $(c, o)$  identifier. We do not represent the  $L$  and  $Q$  sequences of Section 3.1, just the  $D$



sequence of block identifiers. Each leaf has at most  $t + O(\log \log n)$  wasted bits, for the unused space at the end and to store the exact length of the  $D$  sequence within the block. This amounts to  $O(n/f(n))$  wasted bits overall.

**5.2.1 Queries Inside a Superblock.** A table  $G$ , similar in spirit to that of Section 5.1, is used to decode  $\Theta(\log n)$  bits from the leaf in constant time.  $G(x) = (b, r, l)$  indicates that it could decode up to  $l \leq t$  bits from  $x$  (since the rest did not encode a whole block), where it found  $b$  encoded blocks, adding up  $r$  1-bits overall. When  $G(x) = (0, 0, 0)$ , we are in presence of a long code (of length  $> t$ ), which is decoded in constant time as follows. We first read the  $O(\log \log n)$  bits of  $c$  in constant time. Then, a small universal table  $C(c) = \lceil \log \binom{t}{c} \rceil$  tells us the number of bits of the  $o$  entry. We read in constant time the next  $C(c)$  bits, which gives us  $o$ . Finally, a table  $U(c, o) = (x, r)$  gives us the explicit  $t$ -bit content  $x$  of the block encoded as  $(c, o)$ , and its total number  $r$  of 1-bits. Thus, in constant time we decode  $\Theta(\log n)$  bits from the leaf, and at least one entry. Tables  $G$  and  $U$  require  $O(\sqrt{n} \text{ polylog}(n))$  bits of space, whereas  $C$  requires  $O(\text{polylog}(n))$  bits.

A  $\text{rank}(S, i)$  query inside a leaf is handled in  $O(f(n) + \log n)$  time by decoding successive blocks using  $G$  and adding up the  $r$  values (that is, 1-bits), as long as the sum of the  $t \cdot b$  values (that is, processed block lengths) does not exceed  $i$ . The  $l$  values delivered by  $G$  are used to advance in the encoded sequence. Once the sum of  $tb$  values exceeds  $i$  after a  $G$  access, we reread those bits block-wise using  $C$  and  $U$  (even for short codes), and add up the  $r$  values given by  $U$ , until we read the block that contains position  $i$ . This last block is reprocessed bitwise using the  $x$  value given by  $U$ . Overall we spend  $O(f(n) + \log n)$  time.

Similarly  $\text{select}(S, j)$  is solved by adding values  $tb$  until the sum of the  $r$  values exceeds  $j$ , then rereading the last argument of  $G$  block-wise until  $j$  is exceeded again, and finally processing the last block bit-wise. For  $\text{select}_0(S, j)$  we must add values  $tb$  until the sum of  $tb - r$  values exceeds  $j$ .

**5.2.2 Inserting and Deleting Bits.** To insert a bit  $a$  preceding position  $i$  in  $S$  we sequentially find, using  $G$ ,  $C$ , and  $U$ , the block  $b$  where the insertion is to take place,  $b = 1 + \lfloor (i - 1)/t \rfloor$ . All the  $D(1 \dots b - 1)$  entries are directly copied into a new memory area where the updated representation of  $S$  is to be built. On a RAM machine this copying can be done in  $O(f(n))$  time.

The block  $D(b) = (c, o)$  to modify is obtained in constant time with tables  $C$  and  $U$ . Let  $B = a_1 \dots a_t$  be the bits of this block, and let  $i' = i - (b - 1)t$  be the position to insert the bit  $a$  within  $B$ . Thus we compute  $B' = a_1 \dots a_{i'-1} a a_{i'} \dots a_{t-1}$  and save  $a_t$  for later. Again,  $B'$  can be computed in constant time using bit shifts. To compress  $B'$  we use a universal table  $H$ , which given a  $t$ -bit block gives its  $(c, o)$  representation.  $H$  requires  $O(\sqrt{n} \log n)$  bits, and gives  $H(B') = (c', o')$  in constant time. This description  $D(b)' = (c', o')$  is appended at the updated copy of  $S$  we are constructing.

We must now take care of the remaining blocks to the right. We have a bit  $a_t$  that fell off  $B$ . To perform all this propagation in  $O(f(n))$  time, we use yet another universal table  $J(a, x)$ , where  $a$  is a bit to insert at the beginning of the next block and  $x$  is the sequence of the first (compressed)  $t$  bits of  $D(b + 1 \dots)$ .  $J(a, x) = (D', a')$  means that, if we decode from  $x$  as many integral blocks as we

can, append bit  $a$  at the beginning, and reencode them, we obtain sequence  $D'$  and bit  $a'$  falls off at the end of  $D'$ . Another table  $V(x) = r$  tells us how many bits we could use from  $x$ , so we can advance in the processing of sequence  $D$  by  $r$  bits after having copied  $D'$  to the new version of  $S$  we are constructing. If  $V(x) = 0$ , this means that  $x$  starts a long block (that is, whose compressed representation occupies more than  $t$  bits). In this case we treat the block individually: We decode it using  $C$  and  $U$ , insert bit  $a$  at its beginning, call  $a'$  the bit that overflows at its end, and recompress it using  $H$ . Therefore, in constant time we process  $\Theta(\log n)$  bits from the leaf, and at least one entry. The process continues until we complete the leaf and then replace  $S$  by its updated version. Note we still have one overflown bit.

Tables  $J$  and  $V$  require  $O(\sqrt{n} \text{ polylog}(n))$  bits. With the occasional help of  $C$ ,  $U$ , and  $H$ , they process the leaf in  $O(f(n))$  time, plus the time necessary to write the modified leaf by  $\Theta(\log n)$ -bit chunks.

Let us consider how much can the superblock grow by the insertion of a single bit. If a new block is started (which can occur only in a partial leaf), we need  $O(\log \log n)$  more bits. In addition, the  $D$  entry of a block may grow because its  $(c, o)$  descriptor changes. The maximum value of  $\lceil \log \binom{t}{c+1} \rceil - \lceil \log \binom{t}{c} \rceil$  is  $\lceil \log t \rceil$ , achieved when  $c = 0$ . Propagated over at most  $O(f(n) \log n / \log \log n)$  blocks, the sequence of  $D$  values might be increased by  $\Theta(f(n) \log n)$  bits. This is as large as a whole superblock. Indeed, a single bit insertion might double the size of the superblock in some extreme cases. For example, if the sequence is  $(0^t 1^t)^r$ , all the  $c$  values will be 0 or  $t$ , and the  $o$  indexes will be empty, thus we will store  $f(n) \log n / (\log \log(n) - 1)$  blocks in the superblock. If we now insert a 1 at the beginning of the sequence, each  $o$  descriptor becomes  $\log t = (\log \log n) - 1$  bits wide, which adds up  $f(n) \log n$  extra bits. Still, the new superblock is also  $O(f(n) \log n)$  size and can be output using  $J$  and  $V$  in  $O(f(n))$  time.

**5.2.3 Overflow to the Next Superblock.** At the end of the operation, it might be that the new sequence does not fit within the  $s$  bits allocated to the leaf. If so, we take out as many blocks as necessary from the end of the leaf, so as to move them to the beginning of the next leaf. We have seen that we might have to move up to  $\Theta(f(n) \log n)$  bits. In addition we must insert the excess bit at the next leaf (after the blocks we are moving, if any).

The circular array mechanism is not useful this time. The process completely rewrites the next leaf  $S'$ . We move the overflowing  $D$  entries to the beginning of  $S'$ . Then we must insert the carry bit at the beginning of the original entries of  $S'$ . This can be carried out in  $O(f(n))$  time using tables  $J$  and  $V$ . Yet, this bit insertion may produce another  $O(f(n) \log n)$ -bits overflow, in addition to the original  $O(f(n) \log n)$  bits. We can create a new leaf as soon as we have enough overflown bits. This ensures that at most  $s$  bits are ever propagated to the next leaf. The propagation can thus be carried out in  $O(f(n))$  time per leaf rewritten/created. Moreover, as a leaf of  $s$  bits can grow up to size  $O(s)$ , each leaf can trigger the creation of  $O(1)$  further leaves. The mechanism of partial leaves (Section 4.3) limits the propagation among leaves: only  $O(f(n))$  leaves are rewritten or created in the process.

For deletions we proceed similarly, using a table  $J'$  very similar to  $J$ :  $J'(x, a)$

deletes the first bit of the blocks represented by  $x$  and adds bit  $a$  at their end. The bit  $a$  we give to  $J'$  is obtained in constant time using  $C$  and  $U$ , as the first bit of the block encoded in  $D$  at offset  $V(x)$  from the current position. Also, we ensure that leaves are as full as possible. If some space is left at the end of the leaf, we check that the first blocks from the next leaf can be moved back, and propagate the underflow similarly as the overflows. Partial leaves are handled as before upon deletions. Note that, just as whole leaves can be created due to an insertion, up to  $\Theta(f(n))$  whole leaves can disappear due to a deletion (as their contents can shrink so as to be packed within fewer leaves).

Note that, because of the changes in  $|o|$  widths, an insertion can actually produce an underflow and a deletion can produce an overflow. This is not problematic. Overall (still not considering how to manage tree nodes), we have  $O(1/f(n))$  extra space per bit and  $O(f(n)^2)$  insertion/deletion time.

**5.2.4 Final Global Aspects.** Creation and deletion of leaves and internal tree nodes is easily handled together with the maintenance of  $r()$  and  $p()$  in the tree, as in Section 4.4. We note, however, that we permit that a single update affects  $O(f(n))$  leaves, and it creates/deletes  $O(f(n))$  leaves. At this point, we opt for a red-black tree as our balanced tree structure. Once the leaf to be inserted or deleted is located, the red-black tree needs constant time to rebalance, so this adds up  $O(f(n))$  time per insertion or deletion. As for propagating the red-black coloring upwards the root, the same reasoning used for blocked  $r()$  and  $p()$  updates (Section 4.4) applies. Thus the total work in the tree is  $O(\log n)$ .

Handling changes in  $\log n$  is totally analogous to Section 5.1.2. We must have tables  $G, U, C, H, J, J'$  and  $V$  ready for 5 values of  $\log n$ , and we have plenty of time to build them for the next change of  $\log n$ . The way to handle the case  $\log n = o(w)$  is analogous too.

In this case it is also true that the extra space of the tree and partially full leaves adds up  $n'/f(n)$ , not  $n/f(n)$  (where  $n'$  is the compressed sequence length). Since now times are up to  $O(f(n)^2)$ , we have to choose  $f(n) = \sqrt{\log n}$  to obtain  $O(\log n)$  time and  $O(n'/\sqrt{\log n})$  space.

**THEOREM 7.** *The DYNAMIC BIT VECTOR WITH INDELS problem under RAM model with constraint  $\log n = O(w)$  can be solved using  $nH_0(1+o(1))+O(n \log \log n / \log n + w) = nH_0 + o(n) + O(w)$  bits of space supporting the operations rank, select, insert, and delete, in  $O(\log n)$  worst-case time. Here,  $H_0 \leq 1$  is the empirical zero-order entropy of the sequence.*

## 6. HANDLING SEQUENCES OF SYMBOLS

We show now how our last result on bit sequences (Section 5.2) can be extended to sequences of symbols over a general alphabet  $[1, q]$ . Note that this looks similar to the  $k$ -bit version of Section 4.7, but the operations to support here are quite different.

### 6.1 Queries Inside a Superblock

We use the general scheme of Section 5.2, adapting it to handle larger alphabets. We use superblocks of  $s = f(n) \log n$  bits (or  $f(n) \log_q n$  symbols). Blocks are of  $t = (\log_q n)/2$  symbols, and thus span  $(\log n)/2$  bits. We use the encoding of

[Ferragina et al. 2007], where the  $(c, o)$  pairs of Section 5.2 are extended to handle non-binary sequences (recall the end of Section 3.1).

A table  $G$  similar to that of Section 5.2 decodes  $\Theta(\log n)$  bits from the leaf in constant time.  $G(x) = (b, r_1, \dots, r_q, l)$  indicates that it could decode up to  $l \leq t \log q$  bits from  $x$  (since the rest did not encode a whole block), where it found  $b$  encoded blocks, adding up  $r_a$  occurrences of each symbol  $a \in [1, q]$  overall. When  $b = 0$ , we are in presence of a long code (of more than  $t \log q$  bits), which is decoded in constant time as follows. We first read the bits of  $c$  in constant time (those are at most  $(\log n)/2$  bits according to the second bound at the end of Section 3.1). Then, the rest is handled with tables analogous to  $C$  and  $U$  of Section 5.2.1: now  $C(c)$  tells the length of  $o$  entries of class  $c$ , and  $U(c, o) = (x, r_1, \dots, r_q)$  gives the explicit  $t$ -symbol content of class  $(c, o)$  and all the  $r_a$  values within the block. Thus, in constant time we decode  $\Theta(\log n)$  bits from the leaf, and at least one entry. All these tables require  $O(\sqrt{n} (\log n + q \log \log n))$  bits of space.

Queries  $\text{rank}_a(S, i)$  and  $\text{select}_a(S, j)$  inside a leaf are handled in  $O(f(n) + \log n)$  time just as in Section 5.2, adding up  $r_a$  values. We can also retrieve  $a_i$  within the same complexity, by just locating the right block and using  $U$  to obtain the explicit symbols of it. Actually we can obtain, within the same time complexity, any chunk of  $\log^2 n / \log q$  consecutive symbols. This is the best that can be obtained within that time.

## 6.2 Inserting and Deleting Symbols

The mechanism is totally analogous to Section 5.2. Table  $H$  that recompresses the new block  $B'$  in constant time still requires  $O(\sqrt{n} (\log n + q \log \log n))$  bits. Propagation of symbols to next leaves is analogous as well. Tables  $J$  and  $V$  also require  $O(\sqrt{n} \text{polylog}(n))$  bits.

Let us consider how much can the superblock grow by the insertion of a single bit. If a new block is started in a partial leaf, we need  $O(\log n)$  bits for its  $c$  entry<sup>9</sup>. On the other hand, the growth of the  $D$  entries is still limited by  $\log t = O(\log \log n)$  bits. The upper bound of  $f(n) \log n / \log \log n$  blocks per superblock still holds, as the entries  $c$  require at least of  $\Theta(\log \log n)$  bits as in the binary case. Thus, added over all possible blocks, we have that the block expansion is limited by  $O(f(n) \log n)$ , of the same order of the current block size. All the rest on handling overflows is as in Section 5.2. For deletions we use table  $J'$ , of the same size of  $J$ .

Overall (not yet considering how to manage the tree) we have, on top of  $nH_0$ ,  $O((n \log q)/f(n) + (n \log q)q \log \log n / \log n)$  extra space. The first term is the sequence length divided by the overhead due to partial leaves and unused space at full leaves. The second is due to the  $c$  entries,  $(n/t)q \log \log n$  (first bound at the end of Section 3.1). All the operations are handled within  $O(\log n + f(n)^2)$  time. We can choose, as before,  $f(n) = \sqrt{\log n}$  to obtain  $O(\log n)$  time and  $O(n \log q / \sqrt{\log n})$  extra space for the first term. The second term is  $o(n \log q)$  for  $q = o(\log n / \log \log n)$ .

<sup>9</sup>This comes from the second bound at the end of Section 3.1. A natural question is which is the point of compressing  $b = (\log n)/2$  bits into  $(c, o)$  if just  $c$  takes so much space. Yet, this is just a brutal bound that is sometimes convenient. The other bound we are using is  $O(q \log \log n)$  bits.

### 6.3 Managing the Tree

In each internal node of the tree we must now store the total occurrences of each symbol within the node subtree,  $r_a()$ . This requires  $O(qnH_0/f(n))$  additional bits of space. This is  $o(n \log q)$  as long as  $q = o(\sqrt{\log n})$ .

The search time within the tree is still  $O(\log n)$ , as in all cases only one  $r_a()$  value is involved. Updates, however, are more complicated. A single symbol insertion/deletion may involve moving many symbols to the next leaf, and this in turn involves updating many  $r_a()$  values upwards. Although those movements to the next leaf cancel out at their lowest common ancestor, it is not hard to build examples where we need to update  $\Theta(q \log n)$  values of  $r_a()$  (imagine moving a block from the last leaf of the left root child to the first leaf of the right root child: since  $q < t$  we can have  $q$  updates whose common ancestor is  $\log n$  nodes away). Therefore, insertions and deletions cost  $O(q \log n)$ .

### 6.4 Changing $\log n$

This is analogous to Section 5.2.4. We must have tables  $G, U, C, H, J, J'$  and  $V$  ready for 5 values of  $\log n$ . As long as those tables take sublinear space, we have time to build them for the next change of  $\log n$  (as this requires  $\Theta(n)$  operations, among which we can deamortize the construction of the small tables). The way to handle the case  $\log n = o(w)$  is analogous too.

In this case it holds again that the extra space of the tree and partially full leaves adds up  $O(nH_0/f(n))$ , not  $O(n/f(n))$ . By choosing again  $f(n) = \sqrt{\log n}$  we achieve the following result.

**THEOREM 8.** *The DYNAMIC SEQUENCE WITH INDELS problem under RAM model with constraint  $\log n = O(w)$  and symbols in  $[1, q]$ , for  $q = o(\sqrt{\log n})$ , can be solved using  $nH_0 + o(n \log q) + O(w)$  bits of space, supporting the operations rank and select in  $O(\log n)$  worst-case time, and insert and delete in  $O(q \log n)$  worst-case time. Here,  $H_0 \leq \log q$  is the empirical zero-order entropy of the sequence.*

### 6.5 Handling Larger Alphabets

We now extend the result of the previous section to alphabets of size  $\sigma$ , larger than  $q = o(\sqrt{\log n})$ . The idea is to build a wavelet tree [Grossi et al. 2003] (recall Section 3.2) over sequences represented using Theorem 8 [Ferragina et al. 2007].

Let us assume we represent the sequence for each wavelet tree level using the dynamic solution of Theorem 8. We have the restriction  $q = o(\sqrt{\log n})$ . The wavelet tree has  $O(\log_q \sigma)$  levels. Time complexities for the operations is the number of levels times the cost per level. This is  $O(\log n \log_q \sigma)$  for the query operations, and  $O(q \log n \log_q \sigma)$  for the update operations.

Changes in  $\log n$  occur simultaneously in all symbol sequences, and they are smoothly encapsulated within the tree of each level. Handling the case  $\log n = o(w)$  is also analogous. We share a single memory area for all the  $O(\log \sigma)$  sequences, so that we still need only  $O(1)$   $w$ -bit pointers.

**THEOREM 9.** *The DYNAMIC SEQUENCE WITH INDELS problem under RAM model with constraint  $\log n = O(w)$  and symbols in  $[1, \sigma]$  can be solved using  $nH_0 + o(n \log \sigma) + O(w)$  bits of space, supporting operations rank and select in  $O(\log n \log_q \sigma)$  worst-case time, and insert and delete in  $O(q \log n \log_q \sigma)$  worst-case time, for any*

$q = o(\sqrt{\log n})$ . Here,  $H_0 \leq \log \sigma$  is the empirical zero-order entropy of the sequence, and we assume  $\sigma = o(n)$ .

The case  $q = 2$  corresponds to bit sequences, thus the wavelet tree is built directly over the representation of Section 5.2. In this case the wavelet tree has  $O(\log \sigma)$  levels and all the operations cost  $O(\log n \log \sigma)$ . Another interesting choice is  $q = \log^\epsilon n$ , for  $0 < \epsilon < \frac{1}{2}$ . The height of the wavelet tree is  $O(\frac{1}{\epsilon} \log \sigma / \log \log n)$ . The time complexities are  $O(\frac{1}{\epsilon} \log n \log \sigma / \log \log n)$  for the query operations, and  $O(\frac{1}{\epsilon} \log^{1+\epsilon} n \log \sigma / \log \log n)$  for the update operations.

## 7. DYNAMIC FULL-TEXT INDEXES

In this section we extend the result of Chan, Hon, and Lam [Chan et al. 2004] for the DYNAMIC TEXT COLLECTION problem (recall Section 3.6) in several aspects. The most important is a considerable reduction in space. We also change the model of operation and simplify some structures.

### 7.1 Reducing Space and Increasing Time

The most immediate improvement to [Chan et al. 2004] is to replace their COUNT structure by the one of Theorem 9 (with  $q = 2$  in principle, although other tradeoffs could be interesting too). This converts the time for *rank*, *insert* and *delete* into  $O(\log n \log \sigma)$ , and requires only  $nH_0 + o(n \log \sigma)$  bits of space. Note that  $H_0$  refers to the zero-order entropy of  $A = \text{bwt}(\mathcal{C})$ , but this coincides with the zero-order entropy of  $\mathcal{C}$  as they are a permutation of each other. We immediately obtain an entropy-bound index for counting pattern occurrences on a dynamic collection of texts.

To locate occurrences and display text substrings, we can use their same MARK structure, yet sampling one out of  $\log_\sigma n \log \log n$  text positions, so as to have  $o(n \log \sigma)$  extra bits of space for it. With this sampling and our *rank* structure we can report each occurrence in time  $O(\log^2 n \log \log n)$ . Displaying a text substring of length  $\ell$  can be carried out in time  $O(\log n(\ell \log \sigma + \log n \log \log n))$ .

We obtain, almost automatically, the following compressed version of their structure.

**THEOREM 10.** *The DYNAMIC TEXT COLLECTION problem can be solved with a data structure of size  $nH_0(\mathcal{C}) + o(n \log \sigma) + O(w)$  bits supporting counting of occurrences of a pattern  $P$  in  $O(|P| \log n \log \sigma)$  time, and inserting and deleting a text  $T$  in  $O(|T| \log n \log \sigma)$  time. After counting, any occurrence can be located in time  $O(\log^2 n \log \log n)$ . Any substring of length  $\ell$  from any  $T$  in the collection can be displayed in time  $O(\log n(\ell \log \sigma + \log n \log \log n))$ . Deletion and displaying times assume that we know the lexicographic position of  $T$  within the other texts in the collection. Here  $n$  is the length of the concatenation  $\mathcal{C} = 0 T_1 0 T_2 \cdots 0 T_m$  of the  $m$  texts, and we assume  $\sigma = o(n)$ .*

Note that we have already used the fact, pointed out in Section 3.6, that knowing the lexicographic position  $j$  of  $T$  within the others is sufficient to locate its last character in  $A[j]$ . We also point out that the restriction  $\sigma = o(n)$  comes from structure  $\mathcal{C}$ , which needs  $O(\sigma \log n)$  bits. This is  $o(n \log \sigma)$  as long as  $\sigma = o(n)$ .

It is good time to give simple descriptions for  $\mathcal{C}$  and MARK. We note that  $C[c]$  is the number of occurrences of characters smaller than  $c$  in  $\mathcal{C}$ . Let us consider  $K[c]$  as the number of occurrences of  $c$  in  $\mathcal{C}$ , and build a SEARCHABLE PARTIAL SUMS structure for it. Now  $C[c] = \sum_{c' < c} K[c']$  is a *sum* query, and upon text insertions/deletions we must increase/decrease by 1 some entry in  $K$ . Using Theorem 3, we require  $\sigma \log n(1+o(1))$  bits for  $K$  and can answer  $C[c]$  and perform the updates in  $O(\log n)$  time. This does not affect the given time complexities.

For MARK, we will maintain a text sampling so that the distance between consecutive samples is at most  $\log_\sigma n \log \log n$ , and no three consecutive distances add up less than  $\log_\sigma n \log \log n$ . This ensures  $\Theta(n/(\log_\sigma n \log \log n))$  samples in the collection and  $o(n \log \sigma)$  extra space for these structures. It will also ensure  $O(\log^2 n \log \log n)$  time to carry out the operations, even using a wavelet tree requiring  $O(\log n \log \sigma)$  time per step.

There are two queries to handle using this structure. The first is, given a position in  $\mathcal{A}$ , know whether it is sampled or not, and if it is, know the corresponding value. We maintain an array  $S_A$  with the differences between consecutive sampled positions in  $\mathcal{A}$  (starting with an artificial 1), and a SEARCHABLE PARTIAL SUMS WITH INDELS structure over  $S_A$ . To know whether  $\mathcal{A}[i]$  is sampled, we ask whether  $\text{sum}(S_A, \text{search}(S_A, i)) = i$ . If it is, then it is the  $\text{search}(S_A, i)$ -th sample in the set. The sampled  $\mathcal{A}[i]$  values are stored in a balanced tree in order of increasing  $i$  so that they can easily be found by position.

The second query to handle is for displaying. Given a text position, we wish to know which is the nearest sampled position following it in  $\mathcal{C}$ . For this sake we maintain array  $S_C$ , storing differences between consecutive samples in  $\mathcal{C}$ , and also processed for SEARCHABLE PARTIAL SUMS WITH INDELS. The text sample following  $j$  is thus  $\text{sum}(S_C, \text{search}(S_C, j))$ . We also use  $\text{search}(S_C, j)$  to access a balanced tree storing the samples in text position order and storing the corresponding  $\mathcal{A}$  position.

Upon a character insertion in  $\mathcal{A}[i]$ , corresponding to position  $j$  in  $\mathcal{C}$ , we must increase  $S_A[\text{search}(S_A, i)]$  and  $S_C[\text{search}(S_C, j)]$  by 1. If the latter entry exceeds  $\log_\sigma n \log \log n$ , we insert a new sample for  $i$  and  $j$  at the positions we have found in  $S_A$  and  $S_C$ , respectively. In either case, this corresponds to deleting the current entry and inserting 2 new entries replacing it. The balanced trees that contain the samples are updated too. Note that one of the two new entries might be rather short. To ensure that any three consecutive entries add up more than the minimum limit, we consider the first (second) of the two new entries and the one preceding (following) it, merging them if they add up less than  $\log_\sigma n \log \log n$  (two merge, we remove two entries and insert a new one replacing them). Similarly, when a character is removed from the collection, entries are decremented and merged with a neighboring one if necessary.

All these operations are carried out in  $O(\log n)$  time and  $o(n \log \sigma)$  bits of space, which does not affect time nor space complexities. This description for MARK is simpler than the one in [Chan et al. 2004].

## 7.2 An Improved Model for Handling the Collection

We find that asking the users of the data structure to know the lexicographic position of their texts within the collection is delegating a problem the same data

structure should solve. In this section we introduce a more friendly model that maintains the same time complexities and requires  $O(m \log n)$  additional bits of space. This extra space should be irrelevant unless the texts are very short<sup>10</sup>.

When the user inserts a new text  $T$  into  $\mathcal{C}$ , we return a *handle* for it. The handle is a  $\log m$ -bits number, where  $m$  is the current number of texts in the collection. To delete  $T$  later, we only require its original handle. Pattern occurrences are given in the form  $(i, j)$ , where  $i$  is the handle of the text where the occurrence lies and  $j$  is the position within that text. Finally, to retrieve text substrings we only need the handle of the text to display and the positions within it.

To implement this we store a balanced tree **HANDLE**, where the handles are the keys and are stored at the leaves, and another balanced tree **LEX** where the leaves store the same handles in lexicographical order of their corresponding texts. Associated to each key in **HANDLE** we store a pointer to the leaf corresponding to it in **LEX**. Each internal node in **LEX** contains the size of its subtree, which together with parent pointers, easily permits discovering the lexicographic position of a leaf in **LEX** by an upwards traversal (adding up the size of left subtrees of parent nodes we arrive at from the right child).

Together, both trees permit determining the lexicographic position of a text given its handle in  $O(\log m) = O(\log n)$  time, and require  $O(m \log n)$  bits of space.

After a new text  $T$  is inserted in the collection, we must determine its lexicographical position among the other texts, so as to insert a new corresponding leaf at the correct position in **LEX**. This is easy, as the lexicographic position corresponds to the position where  $t_{|T|}$  was inserted in  $A$ . Once we do this, we insert the new text handle in **HANDLE** and point to the newly created **LEX** leaf. All the operations in **LEX** take  $O(\log m)$  time.

Let us now switch our attention to **HANDLE**. Upon a text insertion we find the smallest unused handle number (this guarantees that the handle will require  $\log m$  bits as promised). This is easily achieved by storing in each internal node of **HANDLE** the subtree size and the maximum handle value stored within: When the numbers in the left subtree differ (as there are holes in there) we descend to the left, otherwise to the right. These data at internal nodes are easy to maintain upon tree updates, all in  $O(\log m)$  time.

Still, note that there is a potential problem with the handle numbers we manage. It is possible to insert  $m$  collections and then delete all but the last one, so that we have only one collection but maintain a handle of  $\log m$  bits. The only way to fix this is to permit the structure to modify handles upon deletions, even those for texts that do not participate in the operations. That is, upon a deletion, we should inform that the largest existing handle has been renamed to use the value of the deleted handle. This ensures that all handles are within  $[1, m]$ . Otherwise the space required by the structure is  $O(m \log(n + M))$ , where  $M$  is the largest number of texts in the collection we have ever had.

The only missing piece is how to report occurrence positions in the format (*han-*

---

<sup>10</sup>Removing this term was the explicit goal in [Chan et al. 2004], precisely for the case of many short texts, so we are addressing different goals. Yet, as mentioned in Section 3.6, they could have converted their  $O(\log^2 n)$  deletion time into  $O(\log n)$  within their model, and without resorting to  $\Psi$ .



dle, local position) instead of absolute position in the collection. A new balanced tree POS stores the handles in text position order. Each handle stores the distance in  $\mathcal{C}$  to the previous leaf (that is, the length of the corresponding text), and internal nodes accumulate these distances. Then it is immediate that the absolute position obtained using MARK can be converted into its handle plus relative position in a root-to-leaf traversal on the tree. Similarly, a display request for  $T_j[l, r]$  is converted using POS into a display request for  $\mathcal{C}[l', r']$ , by having a pointer from HANDLE leaves to their corresponding leaves in POS, and traversing POS from the leaf to the root. Tree POS can easily be maintained in  $O(\log m)$  time for each text insertion and deletion.

Finally, we must consider the case of  $\log n$  changing. In this case we can use a standard method: We maintain three copies of all the extra structures (MARK, HANDLE, LEX, POS,  $C$ ,  $S_A$ ,  $S_C$ ), for  $(\log n) - 1$ ,  $\log n$ , and  $(\log n) + 1$ . When the change occurs we switch to the new structures, and have sufficient time to build new structures for  $(\log n) + 1$  or  $(\log n) - 1$  before  $\log n$  changes again. The case  $\log n = o(w)$  is handled essentially as for the tree in Section 4.6, as we can afford a constant factor penalty in the space overhead of these structures.

**THEOREM 11.** *The DYNAMIC TEXT COLLECTION problem can be solved with a data structure of size  $nH_0(\mathcal{C}) + o(n \log \sigma) + O(m \log n + w)$  bits supporting counting of occurrences of a pattern  $P$  in  $O(|P| \log n \log \sigma)$  time, and inserting and deleting a text  $T$  in  $O(|T| \log n \log \sigma)$  time. After counting, any occurrence can be located in time  $O(\log^2 n \log \log n)$ . Any substring of length  $\ell$  from any  $T$  in the collection can be displayed in time  $O(\log n(\ell \log \sigma + \log n \log \log n))$ . Here  $n$  is the length of the concatenation  $\mathcal{C} = 0 T_1 0 T_2 \cdots 0 T_m$ , and we assume  $\sigma = o(n)$ .*

Note that we have assumed that the method of modifying the handles is acceptable, otherwise the  $O(m \log n)$  extra space is  $O(m \log(n + M))$  as explained. The time complexities do not change.

### 7.3 Space is Actually $h$ -th Order Entropy

Recall the partitioning of  $L$  into  $\ell$  pieces  $L^1 L^2 \cdots L^\ell$  according to the  $h$ -contexts (end of Section 3.5.1): It is sufficient to achieve zero-order entropy within each partition to obtain  $h$ -th order entropy overall. Previous work [Ferragina et al. 2007] on a static setting made use of this property by building wavelet trees over the partitions, so as to obtain  $h$ -th order entropy from the sum of zero-order entropies of the wavelet trees. Trying to maintain such an optimal partitioning under a dynamic setting seems to be very difficult because the partitioning can change abruptly due to a single character insertion. It is still possible to maintain a dynamic partition for a given fixed context length  $h$ , by keeping a trie of the existing contexts with a local wavelet tree at each trie leaf. In this section, however, we prove a much more striking result: We show that the solution obtained in the previous section, with just a single wavelet tree for all the text, *is* indeed an  $nH_h$ -bits space solution.

For the proof, let us first state formally a couple of results reviewed earlier in the paper. The first is mentioned in Section 3.2, and the second in Section 3.5.1.

**LEMMA 1** [GROSSI ET AL. 2003]. *Let  $L$  be a string and  $B_v$  the corresponding binary sequence for each node  $v$  of the wavelet tree of  $L$ . Then  $\sum_v |B_v| H_0(B_v) = |L| H_0(L)$ .*

LEMMA 2 [MANZINI 2001]. *Let  $L = L^1 L^2 \dots L^\ell$  be a partition of  $L = \text{bwt}(T)$ , according to contexts of length  $h$  in  $\mathcal{M}$ . Then  $\sum_{1 \leq i \leq \ell} |L^i| H_0(L^i) = n H_h(T)$ .*

We are now ready to prove our main Lemma.

LEMMA 3. *Let  $L = L^1 L^2 \dots L^\ell$  be any partition of  $L = \text{bwt}(T)$ . If the bitmaps are compressed using [Raman et al. 2002], then the number of bits used by a partition  $L^j$  in the wavelet tree of  $L$  is upper bounded by  $|L^j| H_0(L^j) + O(|L^j| \log \sigma \log \log n / \log n + \sigma \log n)$ .*

PROOF. The bits corresponding to  $L^j$  form a substring of the bit vectors at each node of the wavelet tree, as their positions are mapped to the left and right child using  $\text{rank}_0$  or  $\text{rank}_1$ , thus order is preserved. Let us consider a particular node of the wavelet tree and call  $B$  its bit sequence. Let us also call  $B^j$  the substring of  $B$  corresponding to partition  $L^j$ , and assume  $B^j$  has  $l^j$  bits set. Consider the blocks of  $t$  bits that compose  $B$ , according to the partitioning of [Raman et al. 2002] (Section 3.1). Let  $B_{blk}^j = B_1^j B_2^j \dots B_b^j$  be the concatenation of those bit blocks that are *fully contained* in  $B^j$ , so that  $B_{blk}^j$  is a substring of  $B^j$  of length  $t \cdot b$ . Assume  $B_i^j$  has  $l_i^j$  bits set, so that  $B_{blk}^j$  has  $l_1^j + \dots + l_b^j \leq l^j$  bits set. The space the  $o$  fields of the  $(c, o)$  representations of blocks  $B_i^j$  take in the compressed  $B_{blk}^j$  is

$$\sum_{i=1}^b \left\lceil \log \binom{t}{l_i^j} \right\rceil \leq \log \binom{t \cdot b}{l_1^j + \dots + l_b^j} + b \leq \log \binom{|B^j|}{l^j} + b \leq |B^j| H_0(B^j) + b$$

where all the inequalities hold by simple combinatorial arguments [Pagh 1999] and have been reviewed in Section 3.1.

Note that those  $B^j$  bit vectors are precisely those that would result if we built the wavelet tree just for  $L^j$ . According to Lemma 1, adding up those  $|B^j| H_0(B^j)$  over all the  $O(\sigma)$  wavelet tree nodes gives  $|L^j| H_0(L^j)$ . To this we must add three space overheads. The first is the extra  $b$  bits above, which add up  $O(|L^j| \log \sigma / \log n)$  over the whole wavelet tree because  $t \cdot b \leq |B^j|$  and the  $|B^j|$  lengths add up  $|L^j|$  at each wavelet tree level. The second overhead is the space of the blocks that overlap with  $B^j$  and thus were not counted: As  $B^j$  is a substring of  $B$ , there can be at most 2 such blocks per wavelet tree node. At worst they can take  $O(\log n)$  bits each, adding up  $O(\sigma \log n)$  bits over the whole wavelet tree. The third overhead is that of the  $c$  fields, which add up  $O(|L^j| \log \sigma \log \log n / \log n)$ .  $\square$

The above lemma lets us split the wavelet tree “horizontally” into pieces. Let us add up all the zero-order entropies for the pieces. If we partition  $L$  according to contexts of length  $h$  in  $\mathcal{M}$ , and add up all the space due to all partitions in the wavelet tree, we get  $\sum_{1 \leq j \leq \ell} |L^j| H_0(L^j) = n H_h(T)$  (Lemma 2). To this we must add (i)  $O(|L^j| \log \sigma / \log n)$ , which sums up to  $O(n \log \sigma / \log n) = o(n \log \sigma)$  bits over all the partitions; (ii)  $O(\sigma \log n)$  bits per partition, which gives  $O(\ell \sigma \log n)$ ; and (iii)  $O(|L^j| \log \sigma \log \log n / \log n)$ , which sums up to  $O(n \log \sigma \log \log n / \log n) = o(n \log \sigma)$ . In the partitioning we have chosen we have  $\ell \leq \sigma^h$ , thus the upper bound  $n H_h + o(n \log \sigma) + O(\sigma^{h+1} \log n)$  holds for the total number of bits spent in the wavelet tree. The next theorem immediately follows.

THEOREM 12. *The space required by the wavelet tree of  $L = \text{bwt}(T)$ , if the bitmaps are compressed using the technique of [Raman et al. 2002], is  $n H_h(T) +$*

$o(n \log \sigma) + O(\sigma^{h+1} \log n)$  bits for any  $h \geq 0$ . This is  $nH_h(T) + o(n \log \sigma)$  bits for any  $h \leq (\alpha \log_\sigma n) - 1$  and any constant  $0 < \alpha < 1$ . Here  $n$  is the length of  $T$  and  $\sigma$  its alphabet size.

Note that this holds *automatically and simultaneously* for any  $h$ , and we do not even have to care about  $h$  in the index. The next improvement over Theorem 11 is now immediate.

**THEOREM 13.** *The DYNAMIC TEXT COLLECTION problem can be solved with a data structure of size  $nH_h(\mathcal{C}) + o(n \log \sigma) + O(\sigma^{h+1} \log n + m \log n + w)$  bits, simultaneously for all  $h$ . It supports all the operations of Theorem 11 with the same time complexities. For  $h \leq (\alpha \log_\sigma n) - 1$ , for any constant  $0 < \alpha < 1$ , the space complexity simplifies to  $nH_h(\mathcal{C}) + o(n \log \sigma) + O(m \log n + w)$  bits.*

Something striking about the above result is that it holds for the static full-text self-indexes in the literature that build on the wavelet tree of the BWT of the text [Mäkinen and Navarro 2005; Ferragina et al. 2007], but this has gone unnoticed and more complicated arrangements have been made to reach  $h$ -th order entropy. In [Mäkinen and Navarro 2005], they first run-length compress the BWT in order to reduce its length to  $O(nH_h)$  and then apply the BWT. In [Ferragina et al. 2007] they explicitly cut the BWT into pieces  $L^j$  so that the sum of  $nH_0$  sizes of the pieces adds up  $nH_h$ . In both papers, the simpler version they build on (just the wavelet tree of the BWT) would have been sufficient. Thus, we have achieved a significant simplification in the design of static full-text indexes as well. (There are other results in those papers, some of which we have used here.)

In [Ferragina and Manzini 2004] they propose an algorithm to cut  $A$  optimally, so as to minimize the sum of local zero-order entropies plus the overheads of maintaining the separate structures. The optimum partitioning might not correspond to any fixed  $h$  value, but rather use longer contexts in some parts of  $A$  and shorter in others. What we have shown is that the space produced by *any* splitting of  $A$  into pieces is achieved in the simple arrangement having just one wavelet tree, without the need of finding such an optimal partitioning. Their technique, on the other hand, is more general as it works for any zero-order compressor.

Also the paper where the wavelet tree is originally proposed [Grossi et al. 2003] as an internal tool to design one of the most space-efficient compressed full-text indexes, would benefit from our simplification. They cut  $A$  into a table of *lists* (columns) and *contexts* (rows). All the entries across a row correspond to a contiguous piece of  $A$ , that is, some context  $L^j$ . A wavelet tree is built over each table row so as to ensure, again, that the sum of zero-order entropies over the rows adds up to global  $h$ -th order entropy. Our finding implies that all rows could have been concatenated into a single wavelet tree and the same space would have been achieved. This would greatly simplify the original arrangement and possibly expose the deep relationship with the BWT-based approaches [Navarro and Mäkinen 2007]. Interestingly, in [Grossi et al. 2004] they find out that, if they use gap encoding over the successive values along a *column*, and they then concatenate all the columns, the total space is  $O(nH_h)$  without any table partitioning as well. Both findings share the same source: the sum of zero-order entropies of the table cells, no matter the order, adds up to  $nH_h$ .

Finally, it is interesting to point out that, in a recent paper [Ferragina et al. 2006], the possibility of achieving  $h$ -th order compression when applying wavelet trees over the BWT is explored (among many other results), yet they resort to run-length compression to achieve this. Once more, our finding is that this is not really necessary to achieve  $h$ -th order compression if the levels of the wavelet tree are represented using the technique of block identifier encoding [Raman et al. 2001].

Another consequence of our result is that we obtain an  $O(n \log n \log \sigma)$  time construction algorithm for a compressed self-index requiring  $nH_h + o(n \log \sigma)$  bits *working space* during construction: This is obtained by just inserting text  $T$  into an empty collection. This index can be easily converted into a more efficient static self-index, where a static wavelet tree requires the same space and reduces the  $O(\log n \log \sigma)$  time complexities to  $O(\lceil \log \sigma / \log \log n \rceil)$  [Ferragina et al. 2007].

Therefore, we have obtained the *first* compressed self-index with space essentially equal to the  $h$ -th order empirical entropy of the text collection, which in addition can be built within this working space. Alternative dynamic indexes or constructions of self-indexes [Ferragina and Manzini 2000; Hon et al. 2003a; Arroyuelo and Navarro 2005; Chan et al. 2007] achieve at best  $O(nH_h)$  bits of space (with constants larger than 4), and in many cases worse time complexities, as explained in the Introduction.

Note also that, from the dynamic index just built, it is very easy to obtain the BWT of  $T$ . It is a matter of finding the characters of  $A$  one by one. This takes  $O(n \log n \log \sigma)$  time, just as the construction, and gives an algorithm to build the BWT of a text (achieving  $h$ -th order compression) within entropy bounds. The best result we know of, in terms of space complexity [Kärkkäinen 2004], achieves  $O(n \log^2 n)$  time ( $O(n \log n)$  on average) using  $O(n)$  bits in addition to the  $n \log \sigma$  bits of the text.

## 8. FINAL REMARKS

We have introduced a technique to maintain a dynamic bit sequence of length  $n$  using  $nH_0 + o(n)$  bits of space, where  $0 \leq H_0 \leq 1$  is the zero-order entropy of the sequence. The structure answers *rank* and *select* queries and permits insertions and deletions of bits, in worst-case logarithmic time. This is the first dynamic data structure achieving this space. Closely related lower bounds [Patrascu and Demaine 2004] suggest that the time complexities are optimal, yet a proof is missing for this particular set of operations.

From this central result we have uncovered many connections with other problems and derived a surprising number of results in their dynamic setups, using less space and/or time compared to the best existing solutions. We have obtained improved update times and slightly better space for searchable partial sums with indels; the first results on dynamic sequences over alphabets of size  $\sigma$  (achieving zero-order entropy space with times of the form  $O(\log n \log \sigma)$  per operation); compressed dynamic full-text self-indexes and compressed construction of full-text self-indexes (achieving high-order entropy space and  $O(\log n \log \sigma)$  worst-case time per character in all the operations).

All our results are worst-case and support varying  $\log n$ , being valid even for the case  $\log n = o(w)$ , where  $w$  is the size of the machine word. The traditional

techniques to support varying  $\log n$  cannot be directly adapted because we cannot afford the extra space to maintain several copies of the data structure.

Some of the results we have achieved match existing lower bounds. Yet, others seem to be improvable. In particular, it should be possible to improve the  $O(\log n \log \sigma)$  time complexity for accessing and updating dynamic wavelet trees, perhaps with a fractional cascading mechanism. This would immediately affect several other results we achieved.<sup>11</sup> Alternatively, it could be possible to dynamize other static methods not based on wavelet trees, which achieve  $O(\log \log \sigma)$  time for queries [Golynski et al. 2006].

On the other hand, we have achieved zero-order dynamic representations for the data sequence itself. There exist static high-order representations [Sadakane and Grossi 2006; González and Navarro 2006; Ferragina and Venturini 2007] that can be composed with extra data for computing *rank* and *select*. If dynamized, such compressed representations would immediately yield high-order dynamic compressed sequences.

Note that we have indeed achieved high-order compression for dynamic sequences, yet for performing text searching operations on them, not *rank/select* (this is actually intriguing). More precisely, we have shown that wavelet trees, when built over the BWT of a text, automatically achieve high-order entropy. This translates into a significant simplification to many existing self-indexes that achieve high-order entropy (e.g., [Mäkinen and Navarro 2005; Ferragina et al. 2007]), by showing that the base technique they build on naturally achieves the result without need of any further engineering. Our finding also impacts several other works that use this technique in one form or another [Grossi et al. 2003; Ferragina and Manzini 2004; Ferragina et al. 2006].

Still, the results in [Mäkinen and Navarro 2005; Ferragina et al. 2007] have practical value. In their actual implementation (<http://pizzachili.dcc.uchile.cl> or <http://pizzachili.di.unipi.it>), zero-order entropy is achieved by using *uncompressed* bit streams over a *Huffman-shaped* wavelet tree, as this requires less space overhead and implementation effort than using the technique of [Raman et al. 2002] over balanced wavelet trees. In this case the locality property does not hold, and  $h$ -th order entropy would not be achieved if just the simple wavelet tree of the BWT was used. Now, our findings suggest that implementing the technique of [Raman et al. 2002] over a balanced wavelet tree is indeed promising, as in exchange for its (sublinear) space overhead and implementation effort, it would need no extra data structure to achieve higher order compression. Thus we expect it to constitute a simple and competitive alternative in practice.

This, in particular, would immediately derive into a simple and powerful practical algorithm to build, within entropy bounds, different compressed self-indexes, the BWT of a text, and so on. Moreover, this can have a noticeable practical impact on difficult real-life problems such as building indexes for texts that do not fit in main

---

<sup>11</sup>Indeed, in a very recent paper [Lee and Park 2007], they build over our scheme to achieve  $n \log \sigma (1 + o(1))$  bits of space and  $O(\log n (1 + \frac{\log \sigma}{\log \log n}))$  time. Using the same space, another very recent result [Gupta et al. 2006b] achieves  $O(\frac{1}{\epsilon} \log \log n)$  time for queries in exchange for  $O(\frac{1}{\epsilon} n^\epsilon)$  time for updates, for any constant  $\epsilon > 0$ . Both update times are amortized, and the schemes require asymptotically no extra data on top of the sequence, yet they do not achieve compression.

memory, even compressed. In practice, one of the best algorithms for this problem [Crauser and Ferragina 2002] is still a multi-pass technique with I/O complexity  $O(n^2/M)$  [Gonnet et al. 1992], where  $M$  is the maximum text size that can be indexed in main memory. The use of our compressed construction technique on main memory translates into much larger values of  $M$ , and thus fewer passes over the disk, in exchange for (much less important) higher CPU times within each pass.

This is connected with possibly the most important current challenge for compressed data structures, and for compressed full-text self-indexes in particular. Compressed data structures mainly aim at avoiding the use of disk whenever possible, usually in exchange for slower operation in main memory. This pays off by far because main memory is much faster than secondary memory (and this happens at any level of the memory hierarchy, e.g., one can achieve better cache usage just because more compressed data fit in the cache, even if the access patterns are not particularly cache-friendly). Yet, when the data does not fit in main memory anyway, one wishes to have a compressed data structure with good locality of reference, so as to minimize the I/O complexity. This is challenging because better space usage does not automatically translate into fewer block accesses. Indeed, many of the existing solutions suffer from poor locality of reference.

## REFERENCES

- APOSTOLICO, A. 1985. The myriad virtues of subword trees. In *Combinatorial Algorithms on Words*. NATO ISI Series. Springer-Verlag, 85–96.
- ARROYUELO, D. AND NAVARRO, G. 2005. Space-efficient construction of LZ-index. In *Proc. ISAAC'05*. LNCS 3827. 1143–1152.
- BELL, T., CLEARY, J., AND WITTEN, I. 1990. *Text Compression*. Prentice Hall, Englewood Cliffs, New Jersey.
- BLANDFORD, D. AND BLELLOCH, G. 2004. Compact representations of ordered sets. In *Proc. 15th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 11–19.
- BURROWS, M. AND WHEELER, D. 1994. A block sorting lossless data compression algorithm. Tech. Rep. Technical Report 124, Digital Equipment Corporation.
- CHAN, H.-L., HON, W.-K., AND LAM, T.-W. 2004. Compressed index for a dynamic collection of texts. In *Proc. 15th Annual Symposium on Combinatorial Pattern Matching (CPM)*. LNCS 3109. 445–456.
- CHAN, H.-L., HON, W.-K., LAM, T.-W., AND SADAKANE, K. 2007. Compressed indexes for dynamic text collections. *ACM Transactions on Algorithms* 3, 2, article 21.
- CRAUSER, A. AND FERRAGINA, P. 2002. A theoretical and experimental study on the construction of suffix arrays in external memory. *Algorithmica* 32, 1, 1–35.
- DIETZ, P. 1989. Optimal algorithms for list indexing and subset rank. In *Proc. WADS'89*. 39–46.
- ELIAS, P. 1975. Universal codeword sets and representation of the integers. *IEEE Transactions on Information Theory* 21, 2, 194–20.
- FERRAGINA, P., GIANCARLO, R., AND MANZINI, G. 2006. The myriad virtues of wavelet trees. In *Proc. 33rd International Colloquium on Automata, Languages and Programming (ICALP)*. 560–571.
- FERRAGINA, P. AND MANZINI, G. 2000. Opportunistic data structures with applications. In *Proc. FOCS'00*. 390–398.
- FERRAGINA, P. AND MANZINI, G. 2004. Compression boosting in optimal linear time using the Burrows-Wheeler transform. In *Proc. 15th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 655–663.
- FERRAGINA, P. AND MANZINI, G. 2005. Indexing compressed texts. *Journal of the ACM* 52, 4, 552–581.

- FERRAGINA, P., MANZINI, G., MÄKINEN, V., AND NAVARRO, G. 2007. Compressed representation of sequences and full-text indexes. *ACM Transactions on Algorithms* 3, 2, article 20.
- FERRAGINA, P. AND VENTURINI, R. 2007. A simple storage scheme for strings achieving entropy bounds. *Theoretical Computer Science* 372, 1, 115–121.
- GOLYNSKI, A., MUNRO, I., AND RAO, S. 2006. Rank/select operations on large alphabets: a tool for text indexing. In *Proc. 17th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 368–373.
- GONNET, G., BAEZA-YATES, R., AND SNIDER, T. 1992. *Information Retrieval: Data Structures and Algorithms*. Prentice-Hall, Chapter 3: New indices for text: Pat trees and Pat arrays, 66–82.
- GONZÁLEZ, R. AND NAVARRO, G. 2006. Statistical encoding of succinct data structures. In *Proceedings of the 17th Annual Symposium on Combinatorial Pattern Matching (CPM 2006)*. LNCS 4009. 295–306.
- GROSSI, R., GUPTA, A., AND VITTER, J. 2003. High-order entropy-compressed text indexes. In *Proc. SODA'03*. 841–850.
- GROSSI, R., GUPTA, A., AND VITTER, J. 2004. When indexing equals compression: Experiments with compressing suffix arrays and applications. In *Proc. 15th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 636–645.
- GROSSI, R. AND VITTER, J. 2006. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. *SIAM Journal on Computing* 35, 2, 378–407.
- GUPTA, A., HON, W.-K., SHAH, R., AND VITTER, J. 2006a. Compressed data structures: dictionaries and data-aware measures. In *Proc. 5th International Workshop on Experimental Algorithms (WEA)*. 158–169.
- GUPTA, A., HON, W.-K., SHAH, R., AND VITTER, J. 2006b. Dynamic rank/select dictionaries with applications to XML indexing. Tech. Rep. CSD TR #06-014, Purdue University. July.
- HON, W.-K., LAM, T.-W., SADAKANE, K., AND SUNG, W.-K. 2003. Constructing compressed suffix arrays with large alphabets. In *Proc. 14th Annual International Symposium on Algorithms and Computation (ISAAC)*. 240–249.
- HON, W.-K., SADAKANE, K., AND SUNG, W.-K. 2003a. Breaking a time-and-space barrier in constructing full-text indexes. In *Proc. FOCS'03*. 251–260.
- HON, W.-K., SADAKANE, K., AND SUNG, W.-K. 2003b. Succinct data structures for searchable partial sums. In *Proc. ISAAC'03*. LNCS 2906. 505–516.
- KÄRKKÄINEN, J. 2004. Fast BWT in small space by blockwise suffix sorting. In *Proc. DIMACS Working Group on the Burrows-Wheeler Transform: Ten Years Later*.
- LEE, S. AND PARK, K. 2007. Static and dynamic rank-select dictionaries for run-length encoded texts. In *Proc. 18th Annual Symposium on Combinatorial Pattern Matching (CPM)*. LNCS 4580. To appear.
- MÄKINEN, V. AND NAVARRO, G. 2005. Succinct suffix arrays based on run-length encoding. *Nordic Journal of Computing* 12, 1, 40–66.
- MÄKINEN, V. AND NAVARRO, G. 2007. Rank and select revisited and extended. *Theoretical Computer Science*. To appear.
- MANBER, U. AND MYERS, G. 1993. Suffix arrays: a new method for on-line string searches. *SIAM Journal on Computing* 22, 5, 935–948.
- MANZINI, G. 2001. An analysis of the Burrows-Wheeler transform. *Journal of the ACM* 48, 3, 407–430.
- NAVARRO, G. AND MÄKINEN, V. 2007. Compressed full-text indexes. *ACM Computing Surveys* 39, 1, article 2.
- PAGH, R. 1999. Low redundancy in dictionaries with  $O(1)$  worst case lookup time. In *Proc. 26th International Colloquium on Automata, Languages and Programming (ICALP)*. 595–604.
- PATRASCU, M. AND DEMAINE, E. 2004. Tight bounds for the partial-sums problem. In *Proc. 15th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 20–29.
- RAMAN, R., RAMAN, V., AND RAO, S. S. 2001. Succinct dynamic data structures. In *Proc. WADS'01*. 426–437.

- RAMAN, R., RAMAN, V., AND RAO, S. S. 2002. Succinct indexable dictionaries with applications to encoding  $k$ -ary trees and multisets. In *Proc. SODA '02*. 233–242.
- SADAKANE, K. AND GROSSI, R. 2006. Squeezing succinct data structures into entropy bounds. In *Proc. 17th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 1230–1239.