

# Off-Line Compression by Greedy Textual Substitution

ALBERTO APOSTOLICO, SENIOR MEMBER, IEEE, AND STEFANO LONARDI, STUDENT MEMBER, IEEE

*Greedy off-line textual substitution refers to the following approach to compression or structural inference. Given a long textstring  $x$ , a substring  $w$  is identified such that replacing all instances of  $w$  in  $x$  except one by a suitable pair of pointers yields the highest possible contraction of  $x$ ; the process is then repeated on the contracted textstring until substrings capable of producing contractions can no longer be found. This paper examines computational issues arising in the implementation of this paradigm and describes some applications and experiments.*

**Keywords**—Augmented suffix tree, compression of biological sequences, dynamic text compression, grammatical inference, off-line textual substitution, substring statistics.

## I. INTRODUCTION

In data compression by textual substitution (see, e.g., [1]–[3]), substrings with multiple occurrences in a textstring are replaced by a suitable set of pointers to a unique common copy [for instance, by giving: 1) a textstring position starting from which the substring can be recopied and 2) the length of that substring]. Disparate conventions, regarding issues such as the location of the common copy and the mechanics of the encoding-decoding process, give rise to various *macro schemes* of compression. In general, the relative performance of such schemes depends on many factors, including the often subtle interplay between pointer sizes and dictionary parameters (say, number of entries and average length). Partly in response to this fact, techniques were devised for the compact encoding of integers in an unbounded domain (see, e.g., [4]–[6]). Unfortunately, however, the optimal implementation of the majority of macro schemes translates into  $\mathcal{NP}$ -complete problems [3], even before the problem of encoding of pointers is taken into account. One noteworthy exception to this rule is represented by the well-known

Lempel–Ziv (LZ) schemes [7]–[9], which attain asymptotic optimality both in terms of compression achieved and algorithmic complexity. In LZ data compression, data can be processed *on-line* as it is read, a feature that nicely fits the standard paradigm of sequential transmission. The unidirectional or “polar” nature of pointers is crucial in determining the computational efficiency inherent to this scheme.

In some applications, like, for instance, in the production of a CD-ROM or magnetic disk for massive data dissemination, one could afford to perform the compression *off-line*, in particular, to issue pointers in either direction if this brings an increase in compression. Off-line heuristics may be expected to introduce extra time by whatever sequential implementation, but their possible implementation on parallel, perhaps dedicated architectures (see, e.g., [10], [11]), may be expected to achieve sufficient speed to process streams of large consecutive textfile windows consecutively in real-time for any practical purpose. Within the realm of sequential computation, investing more time in the compression may be desirable and feasible for information destined to be massively distributed, as long as the decompression can be still carried out fast and on-line [12]. In other situations, such as, e.g., in backup archiving, the odds of having to restore the data might be feeble enough that even the requirement that this phase be on-line could be forfeited. Finally, as we briefly illustrate at the end of this paper, the study and implementation of macro schemes of the kind considered here may be of some interest in the germane field of inference of hierarchical structures or grammars for sequences (see, e.g., [13]–[15]).

The idea that some of the polarity or greediness inherent to LZ schemes could be traded in for increased compression is intuitively appealing and not new. In [16]–[18], for instance, the authors discuss variations such as, e.g., relaxing the longest-match criterion in determining the next phrase within an LZ parse. The underlying goal is to try and converge faster to the entropy of the source. In view of the intractability of optimal off-line macro schemes, we concentrate here on the implementation of approximate methods such as one of the simplest possible *steepest descent* paradigm. This will consist of performing repeated stages in each one of which we identify a substring of the current version of the text yielding the maximum compression, and then replace

Manuscript received February 15, 2000; revised June 30, 2000. This work was supported in part by NSF Grants CCR-9201078 and CCR-9700276, by NATO Grant CRG 900293, by British Engineering and Physical Sciences Research Council Grant GR/L19362, by Purdue Research Foundation Grant 690-1398-3145, and by the Italian Ministry of University and Research. This work was published in part in the Proceedings of the Annual Data Compression Conference of the IEEE, Snowbird, Utah, 1998 and 2000.

The authors are with the Department of Computer Sciences, Purdue University, West Lafayette, IN 47907-1398 USA (e-mail: axa@cs.purdue.edu; stelo@cs.purdue.edu).

Publisher Item Identifier S 0018-9219(00)09987-4.

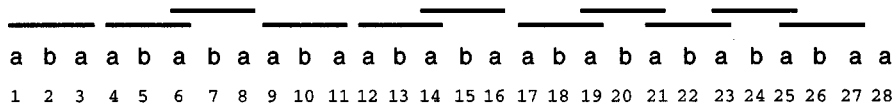


Fig. 1. Overlapping and nonoverlapping occurrences.

all those occurrences except one with a pair of pointers to the untouched occurrence. This is somewhat dual with respect to the bottom up offline scheme introduced by Rubin [19] and recently revived by [20]. As we shall see, this simple scheme already poses some interesting algorithmic problems, some of which we discuss in detail. However, the main issue that we try to address here is that of whether and to what extent a greedy use of bidirectional pointers can yield good compression. As it turns out, the method does outperform all current LZ implementations in most of the cases. More interestingly, it performs quite well on biological sequences and sequence families, where it beats all other generic compression methods, and approaches the performance of methods specifically built around some peculiar regularities of DNA sequences, such as tandem repeats and palindromes, that are neither distinguished nor treated selectively here. The most interesting performances, however, are obtained in the compression of entire groups of genetic sequences forming *families* with similar characteristics. This is becoming a standard and useful way to group sequences in a growing number of important specialized databases. On such inputs, the approach presented here yields scores that are not only better than those of any other method, but also improve increasingly with increasing input size. This is to be attributed to a certain ability to capture distant relationships among the sequences in a family, a feature the merits of which were dramatically exposed in the recent paper [21].

Biological sequences, specially DNA, have been long recognized among important classes of data for which the two tasks of compression and interpretation are often and subtly intertwined. (see, e.g., [22]). The deoxyribonucleic acid (DNA) constitutes the physical medium in which all properties of living organisms are encoded. The knowledge of its sequence is fundamental in molecular biology. Important molecular biology databases (e.g., EMBL, Genbank, DDJB, Entrez, SwissProt, etc.) have been developed to collect hundreds of thousand of sequences of nucleotides and amino-acids from biological laboratories all over the world. The size of these databases, that is currently in the order of thousands of gigabytes, grows at an exponential rate. DNA compression by standard methods such as, e.g., the LZ family of schemes does not seem to fully exploit the redundancies inherent to those sequences. The design of *ad hoc* methods for the compression of genetic sequences constitutes, therefore, an interesting and worthwhile task. Along these lines, a corpus of specialized approaches to DNA compression has been developed in the recent past. As highlighted above, pendant notions of information content and structure have been associated with the compressibility of a sequence. From such a perspective, the amount of compression achievable on genetic sequences has been used in the

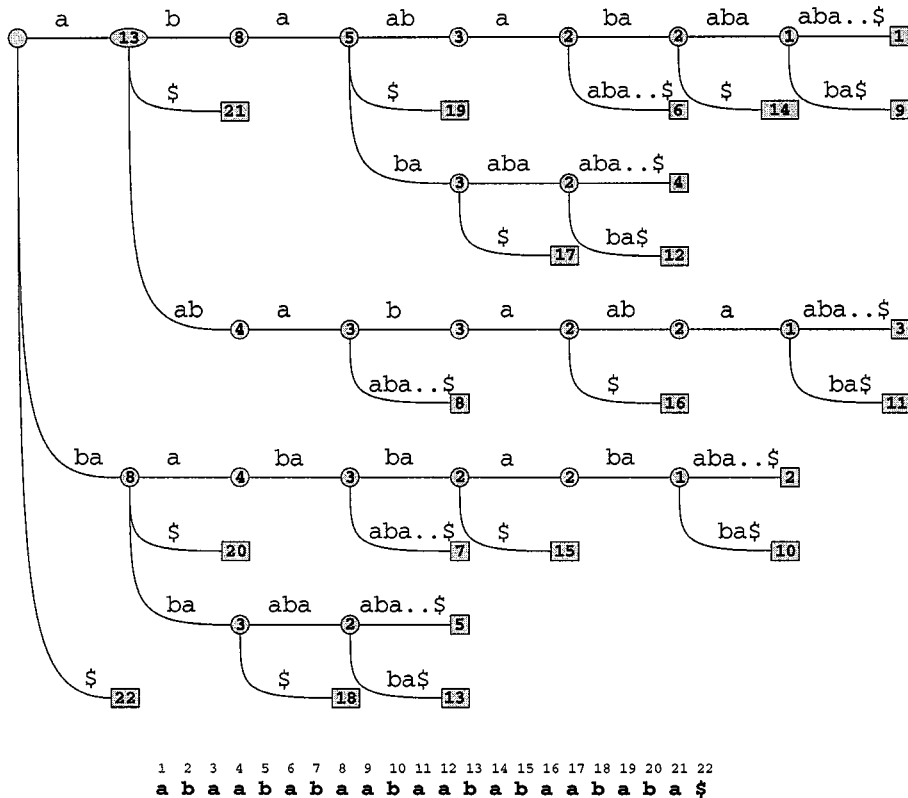
detection of fragments carrying biological significance, or in assessing the relatedness of fragments and sequences. We refer to, e.g., [23]–[32] and references therein for a sampler of the rich literature existing on these subjects.

The structure of the paper is as follows. In Sections II and III, we give some background and notation, and describe a data structure used to gather the statistics of the text. The overall design is presented in Section IV, which is followed by a presentation of the main experimental results in Section V. A discussion of finer implementation details and some final remarks conclude the paper.

## II. SUBSTRING STATISTICS WITHOUT OVERLAP

We use  $\Sigma$  to denote an *alphabet* of symbols. For a string  $x$  over  $\Sigma$ , the number of consecutive symbols in  $x$  is the *length*  $|x|$  of  $x$ , and we write  $x[i]$ ,  $1 \leq i \leq |x|$  to indicate the  $i$ th symbol in  $x$ . In the following, we assume  $|x| = n$ . We use  $x[i, j]$  shorthand for the *substring*  $w$  of  $x$  composed by  $x[i] \cdot x[i + 1] \cdot \dots \cdot x[j]$  where  $1 \leq i \leq j \leq |x|$ , and  $x[i, i] = x[i]$ . Finally, substrings in the form  $x[1, j]$  are called *prefixes* of  $x$ , and substrings in the form  $x[i, |x|]$  are called *suffixes* of  $x$ . For any substring  $w$  of  $x$ , we denote by  $f_w$  the number of *nonoverlapping* occurrences of  $w$  in  $x$ . Clearly,  $f_w$  may be different from the total number of occurrences of  $w$ . For example,  $w = aba$  occurs 11 times in  $x = abaababaabaababaabababababaa$ , with starting positions in the set  $\{1, 4, 6, 9, 12, 14, 17, 19, 21, 23, 25\}$  (see Fig. 1). However, occurrences starting at positions 4 and 6, or 12 and 14, etc., overlap with each other. We can have no more than seven occurrences of  $w$  in  $x$  so that no two of them overlap. For instance, we could take those with starting positions in  $\{1, 4, 9, 12, 17, 21, 25\}$ . Thus,  $f_{aba} = 7$ . To understand our interest in the count of nonoverlapping occurrences, assume that a substring  $w$  appears repeatedly in  $x$ . Then, replacing all occurrences of  $w$  except one with a pointer to the unique reference copy might yield a more compact description of  $x$ . If  $f_w$  is known, then it is also possible to assess beforehand the contraction in length that  $x$  would undergo following such an encoding. If, now, we were asked to identify the one substring  $w$  inducing the highest contraction on  $x$ , we could clearly do so based on the  $f$ -value and length of the individual substrings. Choosing instead on the basis of the total number of occurrences would neither guarantee nor allow us to precompute the best contraction.

The computation of the statistics of all substrings of a string  $x$  is an easy application of the *suffix tree*  $T_x$  of  $x$ . As is well known, the latter is a trie (digital search tree) collecting all the suffixes of  $x\$$ , where  $\$$  is a special symbol not included in  $\Sigma$ . The tree in compact form is built by iterated insertion of consecutive suffixes in  $\Theta(n^2)$  worst case time and  $O(n \log n)$  expected time (see, e.g., [33]). A number of



**Fig. 2.** The minimal augmented suffix tree for `abaababaabaababaababa$` shows node additions and weights needed to change a suffix tree into an index for statistics without overlaps. The label of a leaf is the starting position of the suffix corresponding to that leaf. The label of a node is the number of occurrences without overlap of the string on the path ending at that node or anywhere in the middle of the preceding arc.

more clever constructions are available achieving linear time for finite alphabets (see, e.g., [34]). The number of occurrences (with overlap) of a string  $w$  of  $x$  is trivially given by the number of leaves reachable from the node closest to the locus of  $w$  in  $T_x$ , irrespective of whether or not  $w$  ends in the middle of an arc. Thus, labeling every internal node  $\alpha$  of  $T_x$  with the number  $c(\alpha)$  of the leaves in the subtree rooted at  $\alpha$  yields this statistics for all substrings of  $x$ .

The problem becomes more involved if we wanted to build a similar index for the statistics without overlap. A perusal of Fig. 2 shows that this transition induces a twofold change in our structure: on the one hand, the weight in each node does no longer necessarily coincide with the number of leaves; on the other, extra nodes must be now introduced to account for changes in the statistics that occur in the middle of arcs. The efficient construction of this augmented index in minimal form (i.e., with the minimum possible number of unary nodes) is quite elaborate [35]. For a string  $x$ , the resulting structure is denoted  $\hat{T}_x$  and called the *Minimal Augmented Suffix Tree* of  $x$ . It is not difficult to build  $\hat{T}_x$  in  $O(n^2)$  time and space by embedding the necessary weighting as part of the iterated suffix insertion procedure, hence, at an expected cost of  $O(n \log n)$  [33]. The time required by the construction given in [35] is instead  $O(n \log^2 n)$  in the worst case. The number of auxiliary nodes was bounded by  $O(n \log n)$  in [35]. A tighter  $O(n)$  bound is implied by recent developments in [36].

### III. IMPLEMENTING THE DATA STRUCTURES

When it comes to the actual allocation in memory of a suffix tree, one faces a number of design choices, prominent among which those pertaining to the implementation of nodes. There are three main possibilities in this regard.

- The node is implemented as an array of size  $|\Sigma|$ . This yields fast searches, but is likely to introduce an unbearable amount of waste even for small alphabets.
- The node is implemented as a linked list (or, better, as a balanced search tree). This keeps space to a minimum, but introduces an overhead on the search.
- The adjacency of a node is realized as part of a global hash coding. This yields expected constant time search within overall  $\Theta(n \log n)$  space.

In our case, the space is of high practical concern, so that we use the linked list. Fig. 2 displays the minimal augmented suffix tree of our example textstring. As is customary, the substrings representing edge labels are not stored explicitly in the nodes but rather encoded each by an ordered pair of integers to a unique common copy of  $x$ , so as to achieve overall linear space. However, even linear space can be problematic: at 20 bytes per node and with a number of nodes 1.5 times the number of symbols in the input string, as typically featured in our experiments, a text of size  $n$  needs approximately  $30n$  bytes of storage space. In general, although the

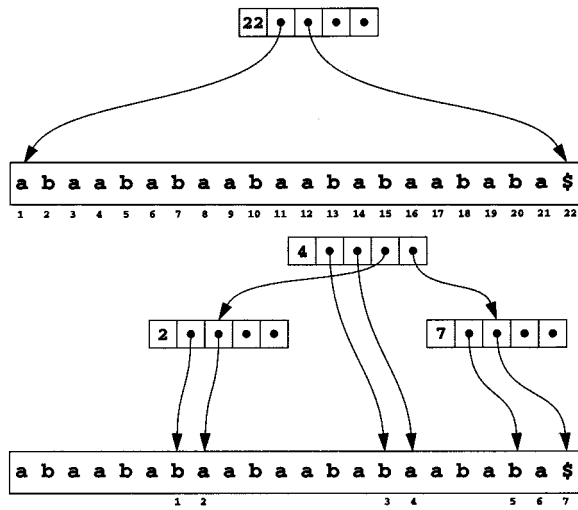


Fig. 3. The data structure allocating textstring `abaababaabaabaabaaba$` prior to and after the removal of `aba`.

size of the suffix tree depends on the particular implementation, one might expect it to be never lower than 20 bytes per input symbol (or *bps*) in the worst case. We refer to [37] for a comparative study of various space-efficient allocations.

In general, these space savings are achieved at the expense of higher complexity in either construction, or searching, or both; thus, for instance, the suffix array and the PAT tree need  $O(n \log n)$  time for the construction [ $O(n)$  on average for the array] and  $O(|w| + \log n)$  when searching for a string  $w$ .

We use  $\langle w \rangle$  to denote the node, if it exists, precisely at the end of the path in  $\hat{T}_x$  labeled by the string  $w$ . If, instead,  $w$  ends in the middle of an arc, then  $\langle w \rangle$  denotes the node corresponding to the shortest extension of  $w$  that ends in a node. In our realization,  $\langle w \rangle$  contains the following items:

- two indices  $[i, j]$  identifying an occurrence of  $w$  in  $x$ , i.e., such that  $w = x[i, j]$ ;
- one pointer to the list of children and one to the list of siblings of  $\langle w \rangle$ ;
- one counter to store the number of nonoverlapping occurrences of  $w$  in  $x$ .

The data structure allocating the textstring  $x$  should support somewhat contrasting primitives such as, for instance, efficient string searching and repeated substring deletions. To accommodate the repeated contractions of  $x$ , the latter is maintained in a linked list of dynamic arrays, as follows. At the beginning, the text is read from the source into a single array of length  $n$ . Subsequently, the removal of the occurrences of a substring  $w = \text{aba}$  will partition the array into linked fragments, as shown in Fig. 3. These arrangements are complemented by refresh cycles that will recombine the text in a single array, from time to time, to counteract excessive fragmentations.

Repeatedly building the suffix tree at each stage exacts a considerable toll irrespective of the method adopted. Ideally, one would like to build the tree once and then maintain it, together with updated statistics, following every substring selection and removal. Linear time algorithms for dynamically maintaining the tree under deletion of a string were originally

proposed by McCreight together with his construction. Similar problems have been studied by Fiala and Green [17] in the context of sliding window compression. More recently, Larsson [38] showed that the algorithm by Ukkonen can be easily extended to accommodate the sliding window update of the suffix tree in amortized linear time. Gu *et al.* [39] introduced a new data structure for dynamic text indexing that supports insertion and deletion of a single character in  $O(\log n)$  time and the  $i$  updates involving a substring  $w$  that occurs  $occ_w$  times in  $O(|w| + occ_w \log i + i \log |w|)$ . Several recent efforts address the dynamic maintenance of tries of various nature. However, we did not find an existing satisfactory solution to the problem of quickly modifying our statistical index so as to reflect the deletion from the corresponding textstring of *all* the occurrences of a given substring. In our experiments, every new version of the suffix tree was built from scratch. In a later section, we present some heuristics designed to alleviate such computation efforts.

#### IV. CHOOSING AND COMPUTING A GAIN MEASURE

By “gain measure,” we refer here to the function  $G$  that drives, at each stage, the selection of the substring that yields the highest compression. In practice, it is not easy to define precisely such a measure, as we explain below.

The main difficulty is due to the fact that at the time when we need to compute the contraction that would be induced by a particular substring, we lack some important costs such as those associated with the optimal encodings of pointers or integers, which can be computed precisely only at the outset. Letting  $l(i)$  represent the number of bits needed to encode integer  $i$ , we assume for simplicity  $l(i) = \lceil \log i \rceil$  at the time the gain is computed. Note that this choice does not affect the appraisal of final compression, the latter being based on purely empirical measures. Along the same lines, one could choose an expression for  $l$  that reflects more accurately the efficient encoding of integers in an unknown range [4]–[6]. However, as long as the ultimate encoding of the compressed string is not based on those representations, but rather on some statistical treatment (e.g., Huffman encoding), there is hardly any sense in resorting to them and hardly any way to compute  $l(i)$  accurately at this stage.

With this choice made, we describe now in succession three possible measures of gain. Let us denote by  $w$  the word that maximizes  $G$  at some iteration and, therefore, is selected to be replaced by pointers.

In Scheme 1, we assume that *all* the  $f_w$  occurrences of the string  $w$  are removed from the text, while  $w$  itself is saved in an auxiliary data structure that contains:

- the string  $w$ , that is  $Bm_w$  bits long, where  $m_w = |w|$  and  $B = \log |\Sigma|$ ;
- the length  $m_w$  of  $w$ , at a cost of  $l(m_w)$  bits;
- the value of  $f_w$ , at a cost of  $l(f_w)$  bits;
- the  $f_w$  positions of  $w$  in  $x$ , at a global cost bounded by  $f_w l(n)$  bits.

Fig. 4 shows the original and compressed representations for the textstring and the corresponding associated costs. The expression underneath the top figure represents the original

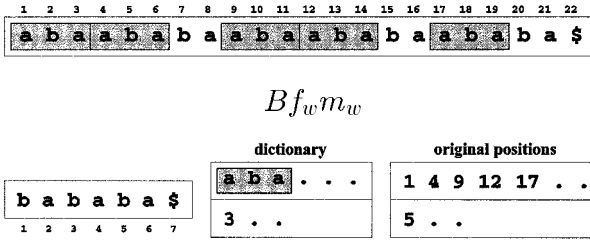


Fig. 4. Illustrating Scheme 1 with **aba**.

Table 1  
Statistical Encoders Available to OFF-LINE

id	encoding type
0	Plain
1	Huffman
2	Arithmetic
3	Deflate (ZLIB)
4	RLE
5	RLE + Huffman
6	RLE + Arithmetic
7	RLE + Deflate (ZLIB)

cost of the occurrences of  $w$  as plain text: with  $|w| = m_w$ , the  $f_w$  copies of  $w$  require  $Bf_w m_w$  bits in the plain text. In practice, the value of  $B$  is appraised based on the empirical entropy of the source: the plain text is encoded with the best statistical encoder from Table 1, and then  $B$  is set to the average length of a symbol.

The expression at the bottom of the figure is the cost of representing the same occurrences of  $w$  based on the external dictionary and auxiliary structures. The difference between the two expressions defines  $G_1(w)$  as follows:

$$G_1(w) = Bf_w m_w - Bm_w - l(m_w) - l(f_w) - f_w l(n) \\ = (f_w - 1)Bm_w - l(m_w) - l(f_w) - f_w l(n).$$

In Scheme 2, we assume that one of the  $f_w$  copies of  $w$  is kept in the original text, marked by a “literal identification” bit, while the remaining  $f_w - 1$  copies are encoded by pointers, each pointer being preceded by a suitable identification bit.

Because of the additional bit, the plain text representation of all the occurrences of  $w$  requires  $(B + 1)f_w m_w$  bits. The pointer-based representation costs are as follows (see Fig. 5):

- $(B + 1)m_w$  bits for the original copy of  $w$ ;
- $(f_w - 1)(l(n) + l(m_w) + 1)$  bits for the  $f_w - 1$  pointers.

The difference of these expressions defines  $G_2$  as follows:

$$G_2(w) = (B + 1)f_w m_w - (B + 1)m_w - (f_w - 1) \\ \cdot (l(n) + l(m_w) + 1) \\ = (f_w - 1)(B + 1)m_w - (f_w - 1)(l(n) + l(m_w) + 1) \\ = (f_w - 1)((B + 1)m_w - l(n) - l(m_w) - 1).$$

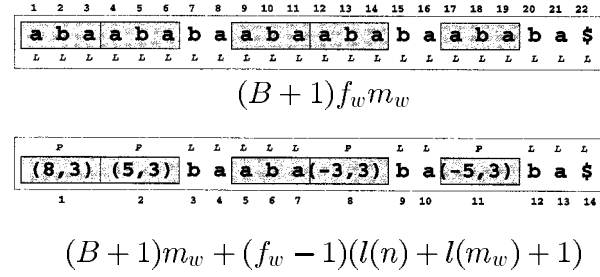


Fig. 5. Illustrating Scheme 2 for **aba**;  $L$  and  $P$  are used to mark literals and pointers, respectively; the pair  $(p, l)$  denotes a pointer to a reference copy that has length  $l$  and that is  $p$  symbols away.

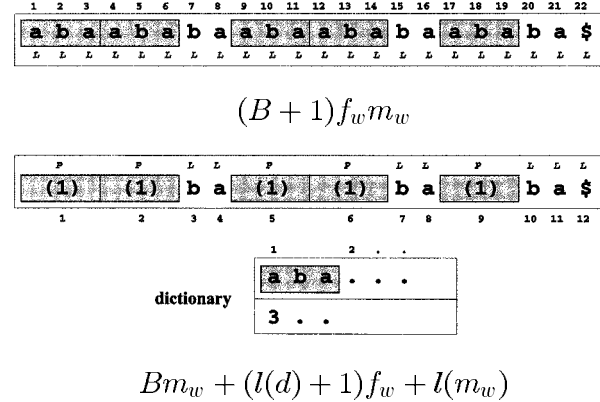


Fig. 6. Illustrating Scheme 3 for  $w = \mathbf{aba}$ ;  $d$  is the size of the dictionary, “(1)” denotes a pointer to the first entry in the dictionary,  $L$  and  $P$  mark literals and pointers, respectively.

In Scheme 3, words in the textfile are replaced by pointers to their corresponding entries in an external dictionary. Thus, following the selection of  $w$  at the generic iteration,  $w$  is added as a new entry into the dictionary and *all* of its occurrences become pointers to that entry. Also in this case, an auxiliary bit-vector is required in general in order to distinguish between pointers and literals at the outset, both in the text and in the dictionary. However, if pointer recursion is forbidden as we assume in our construction, then the words in the dictionary cannot contain pointers, and a bit-vector is not needed there.

The plain text representation of all the occurrences of  $w$  requires  $(B + 1)f_w m_w$  bits. The costs of the pointer-based representation are (see Fig. 6):

- $Bm_w$  bits for the string  $w$  in the dictionary;
- $l(m_w)$  to store the length  $m_w$ ;
- $l(d)f_w$  for the  $f_w$  pointers inside the text, where  $d$  is the size of the dictionary.

The corresponding expression for  $G_3(w)$  is then

$$G_3(w) = (B + 1)f_w m_w - Bm_w - (l(d) + 1)f_w - l(m_w) \\ = B(f_w - 1)m_w + f_w m_w - (l(d) + 1)f_w - l(m_w).$$

We point out that, for any of the above specifications of  $G$  and any word  $w$  in  $x$ ,  $G(w)$  is a monotone increasing function of  $m_w$ . Moreover, the maximum number of nonoverlapping occurrences of  $w$  in  $x$  does not change in the middle of an arc of  $\hat{T}_x$ . Therefore, the word maximizing the gain at each stage *always* ends on a node of  $\hat{T}_x$ . If now  $w$  is this word, then its occurrences are suitably encoded, and the whole process is

```

1: abaababaabaababaabaaba$ Substituted substr: "aba"
2: bababa$ Substituted substr: "ba"
----- Final encoding:

sublen = [3 2]
substr = [ababa]
abspol = [0 0] abspos = [0 0]
relpol = [0 2 0 2 0 0] relpos = [0 0 0 0 0 0]
occurr = [5 3]
text = [$]

```

Fig. 7. A run of OFF-LINE<sub>1</sub> on our example textstring.

repeated until the gain becomes zero or negligible, according to some predetermined threshold.

The three schemes just described were embedded in as many encoders, respectively called OFF-LINE<sub>1</sub>, OFF-LINE<sub>2</sub>, and OFF-LINE<sub>3</sub>. As an example, the iterations of OFF-LINE<sub>1</sub> are highlighted in Fig. 7. The first iteration results in the choice of aba; the second, of ba. The collection of data representing the output encoding appears in the bottom part of the figure. Several implementation details will be given later.

## V. RESULTS

The encoders described in the previous section were coded in C++ using the Standard Template Library (STL) [40] and extensively tested. Table 2 offers a first glance at the performances of the three encoders on two typical inputs, namely, *paper2* from the Calgary Corpus and *mito*, the mitochondrial DNA sequence of the yeast (*Saccharomyces Cerevisiae*). Running times are in the order of 2 or 3 min for files of about 80 KB on a 300-Mhz machine running under Solaris. In terms of compression, the best encoder is OFF-LINE<sub>3</sub>, followed by OFF-LINE<sub>1</sub> and, at some distance, OFF-LINE<sub>2</sub>.

Table 3 compares performances among textual substitution methods over the entire Calgary Corpus. OFF-LINE<sub>3</sub> outperforms the other two encoders on most inputs. As a whole, OFF-LINE encoders perform better than the rest on most inputs, and loose marginally to GZIP where they do. However, a thoroughly faithful comparison to GZIP is made difficult by the many heuristics embedded in that program.

Crossing the boundary of textual substitution methods, the recent block-sorting techniques called BZIP and BZIP2, based on [41], outperform GZIP and OFF-LINE on the whole Calgary Corpus (see Table 4). As seen next, a different scenario is displayed when we turn to biological data sets.

We compare the performance of OFF-LINE encoders with those of standard compression programs in Table 5. The encoder OFF-LINE<sub>3</sub> outperforms each and every general purpose encoder on the fourteen chromosomes and the mitochondrial DNA of the yeast. It should be noted that the actual compressions are very small and sometimes negative.

In fact, raw biological sequences (notably, those coming from coding regions [42]) are known to be hard to compress. However, even comparing our encoders with programs

Table 2

A First Glance at the Three OFF-LINE Encoders' Performances on a 300-MHz Solaris Machine

encoder	paper2 (82,199)		mito (78,521)	
	size	time <sub>[min]</sub>	size	time <sub>[min]</sub>
OFF-LINE <sub>1</sub>	30,848	3.21	16,426	1.66
OFF-LINE <sub>2</sub>	33,757	3.01	17,741	2.24
OFF-LINE <sub>3</sub>	<b>30,219</b>	2.38	<b>16,086</b>	2.38

specifically designed to compress DNA, the difference in performance is not large, as shown in Table 6.

It is worthwhile to highlight such DNA-specific analyzers and compressors. As mentioned, information theoretic analyzes of biological sequences mingle with the very dawn of bioinformatics studies (see, e.g., [22]), but this area has recently known a considerable revival of interest in view of the massive production of genomic sequences of various kinds. In this context, the detection of redundancy serves not only the purpose of achieving more compact descriptors, but also, and perhaps more importantly, may act as a filter of possibly relevant biological functions. The tenet there is that an incompressible string is more random and, thus, less likely than a repetitive one to carry some biological function.

Due to mutations, errors in the sequencing process, and other biological events, a substantial part of the redundancy present in DNA manifests itself in the form of consecutive (*tandem*) repeats of the same word or *motif*, and palindromes. However, such tandem repeats and palindromes are not exact. Rather, they may occur with substitutions, insertions, or deletions of symbols. Moreover, palindromes are actually *complemented*, meaning that in the reverse half of the word the base A is mirrored by a T (and vice-versa), while C is mirrored by a G (and vice-versa).

Among the recent dedicated approaches to DNA compression, the one by Grumbach and Tahi [25], [26], called BIOCOMPRESS2, extends LZ-77 to catch very distant repeats and complementary palindromes.

Loewenstern and Yianilos [27] consider the problem of computing good estimates of the entropy of DNA sequences by building a PPM-like predictive model. With respect to the original PPM, they extend the context model by allowing mismatches. Their algorithm estimates the parameters of the model, called CDNA, via a learning process that tries to optimize a complex objective function. The general problem is known to be  $\mathcal{NP}$ -complete, but they devise more realistic approximation schemes.

Allison *et al.* propose the most computationally intensive approach to DNA compression [43]. They search for both approximate repeats and approximate palindromes. Their primary purpose is not to compress the text, but rather to model the statistical properties of the data as accurately as possible and to find patterns and structures within them. They build a model with parameters such as the probability of repeats, of the length of repeats, and of mismatches within repeats. The parameters of the model are estimated by an

**Table 3**

The Variants of OFF-LINE Against the Other Textual Substitution Compressors, on the Calgary Corpus

<i>File</i>	<i>Size</i> (bytes)	<i>Huffman</i>	<i>LZ-78</i>	<i>LZ-77</i>	OFF-LINE <sub>1</sub>	OFF-LINE <sub>2</sub>	OFF-LINE <sub>3</sub>
		PACK	COMPRESS	GZIP			
<b>bib</b>	111,261	72,868	46,528	35,063	36,145	39,226	<b>34,442</b>
<b>book1</b>	768,771	438,487	332,056	313,376	305,185	323,007	<b>298,735</b>
<b>book2</b>	610,856	368,423	250,759	206,687	<b>203,249</b>	216,494	204,703
<b>geo</b>	102,400	72,836	77,777	68,493	<b>68,229</b>	69,983	68,726
<b>news</b>	377,109	246,516	182,121	144,840	<b>141,257</b>	150,462	143,246
<b>obj1</b>	21,504	16,330	14,048	<b>10,323</b>	10,845	11,271	11,088
<b>obj2</b>	246,814	194,378	128,659	<b>81,631</b>	88,179	93,915	87,574
<b>paper1</b>	53,161	33,457	25,077	<b>18,577</b>	19,994	21,607	19,289
<b>paper2</b>	82,199	47,731	36,161	<b>29,753</b>	30,848	33,757	30,219
<b>pic</b>	513,216	106,737	62,215	56,442	52,036	55,427	<b>50,885</b>
<b>progc</b>	39,611	26,030	19,143	<b>13,275</b>	14,758	15,527	14,127
<b>progl</b>	71,646	43,093	27,148	16,273	18,508	18,919	<b>16,153</b>
<b>progp</b>	49,379	30,328	19,209	11,246	12,890	13,282	<b>11,160</b>
<b>trans</b>	93,695	65,343	38,240	<b>18,985</b>	21,170	21,170	19,662

**Table 4**

Comparing OFF-LINE with Context-Sorting Encoders on the Calgary Corpus

<i>File</i>	<i>Size</i> (bytes)	<i>BWT</i>	<i>BWT</i>	OFF-LINE <sub>1</sub>	OFF-LINE <sub>2</sub>	OFF-LINE <sub>3</sub>
		BZIP	BZIP2			
<b>bib</b>	111,261	<b>27,097</b>	27,467	36,145	39,226	34,442
<b>book1</b>	768,771	<b>230,247</b>	232,598	305,185	323,007	298,735
<b>book2</b>	610,856	<b>155,944</b>	157,443	203,249	216,494	204,703
<b>geo</b>	102,400	57,358	<b>56,921</b>	68,229	69,983	68,726
<b>news</b>	377,109	<b>118,112</b>	118,600	141,257	150,462	143,246
<b>obj1</b>	21,504	<b>10,409</b>	10,787	10,845	11,271	11,088
<b>obj2</b>	246,814	<b>76,017</b>	76,441	88,179	93,915	87,574
<b>paper1</b>	53,161	<b>16,360</b>	16,558	19,994	21,607	19,289
<b>paper2</b>	82,199	<b>24,826</b>	25,041	30,848	33,757	30,219
<b>pic</b>	513,216	<b>49,422</b>	49,759	52,036	55,427	50,885
<b>progc</b>	39,611	<b>12,379</b>	12,544	14,758	15,527	14,127
<b>progl</b>	71,646	<b>15,387</b>	15,579	18,508	18,919	16,153
<b>progp</b>	49,379	<b>10,533</b>	10,710	12,890	13,282	11,160
<b>trans</b>	93,695	<b>17,561</b>	17,899	21,170	21,170	19,662

expectation maximization algorithm that takes time  $O(n^2)$  at each iteration. Their results may well be taken to represent the current “state of the art,” but as said the algorithm is extremely slow.

Finally, we run OFF-LINE<sub>3</sub> on families of related and unrelated genetic sequences. Entries in most genetic databases are flat text files containing one or more sequences that are usually functionally related, with some annotations. The fasta format is the most commonly used standard for storing and exchanging genetic files. The generic fasta

file contains one or more blocks. Each block is composed by one or more annotation lines each starting with the symbol  $>$ , followed by the genetic sequence.

Table 7 shows the results of running OFF-LINE<sub>3</sub> on several families of sequences of the yeast genome. The complete dataset is available at <http://www.cs.purdue.edu/homes/stelo/Off-line/>. The file *Spor\_All\_2x.fasta* is artificially obtained by concatenating *Spor\_All.fasta* with itself, in an attempt to probe into extreme cases of intersequence correlation [21]. The last two families (8 and

**Table 5**  
Comparing OFF-LINE with Other Compression Programs on the Chromosomes of the Yeast

File	Size (bytes)	Huffman	LZ-78	LZ-77	BWT	BWT	OFF-LINE <sub>1</sub>	OFF-LINE <sub>2</sub>	OFF-LINE <sub>3</sub>
		PACK	COMPRESS	GZIP	BZIP	BZIP2			
chrI	230,195	63,144	62,935	66,264	61,674	62,373	57,098	58,631	<b>56,915</b>
chrII	813,137	222,597	219,845	236,837	218,463	221,032	201,617	203,456	<b>201,180</b>
chrIII	315,344	86,281	86,009	91,827	84,809	85,705	77,916	78,983	<b>77,764</b>
chrIV	1,522,191	416,516	409,957	440,056	407,799	411,250	371,230	374,413	<b>370,796</b>
chrV	574,860	157,415	155,944	167,749	154,580	155,731	142,364	143,775	<b>141,919</b>
chrVI	270,148	74,077	73,873	78,925	72,838	73,651	67,451	68,151	<b>67,391</b>
chrVII	1,090,936	298,680	294,417	317,282	293,079	296,245	270,051	272,972	<b>269,265</b>
chrVIII	562,638	154,110	152,265	163,135	151,240	152,992	139,588	140,924	<b>139,271</b>
chrIX	439,885	120,669	118,965	127,805	118,182	119,553	109,507	110,871	<b>109,303</b>
chrX	745,443	204,152	201,783	216,148	200,325	202,223	184,709	186,471	<b>184,287</b>
chrXI	666,448	182,377	180,100	194,119	179,306	180,901	165,780	166,752	<b>165,478</b>
chrXII	1,078,171	295,441	291,754	305,653	288,112	290,800	260,172	261,346	<b>259,898</b>
chrXIII	924,430	253,176	249,099	267,127	248,450	250,735	228,233	231,474	<b>227,610</b>
chrXIV	784,328	215,020	212,219	228,757	210,988	212,816	195,291	196,719	<b>194,947</b>
chrXV	1,091,282	298,762	294,921	317,971	293,838	297,279	270,626	273,366	<b>269,921</b>
chrXVI	948,061	286,579	264,113	278,651	254,947	257,590	234,099	237,365	<b>233,150</b>
mito	78,521	18,149	17,890	19,369	17,962	18,094	16,426	17,741	<b>16,086</b>

**Table 6**  
Comparing OFF-LINE with DNA-Specific Compression Programs on the Third Chromosome (chrIII) of the Yeast (315 344 bps). The Parameter *bpc* Represents the Average Number of Bits Per Character in the Compressed Representation (Some Final Sizes are Extrapolated from Table 1 of [43])

encoder	size	bpc
GZIP	91,827	2.33
PACK	86,281	2.19
COMPRESS	86,009	2.18
BZIP2	85,705	2.17
BZIP	84,809	2.15
OFF-LINE <sub>3</sub>	77,764	1.97
CDNA [27]	76,471	1.94
BIOCOMPRESS2 [26]	75,682	1.92
AED [43]	75,407	1.913

9) are a segment of *all* the upstream regions of the yeast and, thus, not strongly related. Table 7 shows that not only the absolute performance of OFF-LINE, but also its relative advantage over the other methods improves as the input size increases. Likewise, as soon as the input files contain sequences not as strongly related, the improvements, while still present, decay immediately, as shown for files 8 and 9 in the table. The ability to capture distant relationships is enhanced in the comparison with GZIP and BZIP2 as we move from their default window sizes (900 Kb in BZIP2) to smaller sizes. The results, shown in Table 8, suggest that the relative advantage of OFF-LINE will increase as it will be applied to larger and larger families.

## VI. FINE TUNING AND OTHER IMPLEMENTATION DETAILS

The most time-consuming activity of the compression phase is the construction of the index trie and its annotation with the values of the gain. We employed three heuristics to overcome the high computational demands of a “full-fledged” version of the compressor.

Table 9 shows the results achieved by one of these heuristics on the basic algorithm, in which more than just one substring selection and substitution is performed between two consecutive updates of the statistical index. Of course, such an approach saves time on one hand, but it risks blurring the perception of the best candidates for substitution. In our implementation, a heap is maintained with the statistical index, containing at each step the  $Q$  best words in terms of  $G$ , for some chosen value of the parameter  $Q$ . Between any two consecutive index reconstructions, the  $Q$  strings in the heap are retrieved and used in succession in a contraction step for the text. It is possible at some point that a string from the heap will no longer be found in the contracted text. In fact, part of the words in the heap turn out to be useless in general. In any case, as soon as all words in the heap have been considered, a new augmented trie is built on the contracted text.

As the table displays, the number of individual substring substitution passes over the text grows with the maximum allowed size of the heap. On the other hand, we spend less and less time building weighted tries. The overall result is, within a wide interval, a considerable speed up with respect to the eager version of OFF-LINE without substantial penalty in compression performance. When the size of the heap becomes too large (approximately  $Q > 100$  in our experiments), only a small subset of the words in the heap is used: most of the computational effort is spent in pattern searching, which results in deterioration of both speed and compression.



**Table 7**

Comparing OFF-LINE<sub>3</sub> With Other Compression Programs on Families of Sequences of the Yeast. The Figures in Parentheses Report Percentage Gains Achieved by OFF-LINE<sub>3</sub>,  $k$  is the Number of Upstream Sequences in Each Family, Individual Sequence Length is 800 bps Except in the Last Two Rows, Where it is 2000. The Alphabet Consists of About 50 Symbols. The Input Strings 1-9 Correspond, in this Order, to the Families of Spor\_EarlyII.fasta, Spor\_EarlyI.fasta, Helden\_GCN.fasta, Spor\_Middle.fasta, Helden\_All.fasta, Spor\_All.fasta, Spor\_All\_2x.fasta, All\_Up\_400k.fasta, All\_Up\_1M.fasta

Family	Total size		Huffman	LZ-78	LZ-77	BWT	OFF-LINE <sub>3</sub>
	(bytes)	$k$	PACK	COMPRESS	GZIP	BZIP2 -9	
(1)	25,008	29	7,996 <sub>(11.0%)</sub>	7,875 <sub>(9.6%)</sub>	8,008 <sub>(11.1%)</sub>	7,300 <sub>(2.5%)</sub>	<b>7,119</b>
(2)	31,039	36	9,937 <sub>(12.5%)</sub>	9,646 <sub>(9.8%)</sub>	9,862 <sub>(11.8%)</sub>	9,045 <sub>(3.8%)</sub>	<b>8,697</b>
(3)	32,871	38	10,590 <sub>(12.2%)</sub>	10,223 <sub>(9.0%)</sub>	10,379 <sub>(10.4%)</sub>	9,530 <sub>(2.4%)</sub>	<b>9,301</b>
(4)	54,325	63	17,295 <sub>(14.6%)</sub>	16,395 <sub>(9.9%)</sub>	16,961 <sub>(12.9%)</sub>	15,490 <sub>(4.6%)</sub>	<b>14,778</b>
(5)	112,507	130	36,172 <sub>(17.7%)</sub>	33,440 <sub>(11.0%)</sub>	33,829 <sub>(12.0%)</sub>	31,793 <sub>(6.4%)</sub>	<b>29,758</b>
(6)	222,453	258	70,755 <sub>(23.2%)</sub>	63,939 <sub>(15.0%)</sub>	68,136 <sub>(20.3%)</sub>	61,674 <sub>(11.9%)</sub>	<b>54,317</b>
(7)	444,906	516	141,431 <sub>(53.4%)</sub>	124,637 <sub>(47.1%)</sub>	135,816 <sub>(51.5%)</sub>	85,142 <sub>(22.6%)</sub>	<b>65,891</b>
(8)	399,615	191	121,700 <sub>(12.3%)</sub>	115,029 <sub>(7.22%)</sub>	115,023 <sub>(7.22%)</sub>	112,363 <sub>(5.0%)</sub>	<b>106,722</b>
(9)	1,001,002	477	305,054 <sub>(11.9%)</sub>	286,971 <sub>(6.4%)</sub>	285,064 <sub>(5.8%)</sub>	280,334 <sub>(4.1%)</sub>	<b>268,612</b>

**Table 8**

Constraining the Competitors to Work on Small Windows Enhances the Gain of OFF-LINE. Here, the Input Strings 6 and 7 Correspond, Respectively, to the Families of Spor\_All.fasta, Spor\_All\_2x.fasta (see Table 7 for their Respective Statistics)

Family	LZ-77	BWT	OFF-LINE <sub>3</sub>
	GZIP -1	BZIP2 -1	
(6)	76,629 <sub>(29.1%)</sub>	63,332 <sub>(14.2%)</sub>	<b>54,317</b>
(7)	153,103 <sub>(57.0%)</sub>	126,314 <sub>(47.8%)</sub>	<b>65,891</b>

**Table 9**

Performances of OFF-LINE<sub>1</sub> for Different Choices of the Size of the Candidates Heap. We Fixed  $\text{min\_occ} = 2$ ,  $\text{min\_length} = 2$ ,  $l = 100$

$Q$	paper2		mito	
	size	$\text{time}_{[\text{min}]}$	size	$\text{time}_{[\text{min}]}$
1	30,773	19.70	16,326	7.06
2	30,780	10.36	16,367	4.06
5	30,785	5.06	16,405	2.24
10	30,787	3.21	16,446	1.66
20	30,826	2.39	16,476	1.36
50	30,904	1.97	16,632	1.28
100	30,923	1.86	16,702	1.37
1,000	30,923	1.98	16,702	1.47

Whenever one can assume it as being highly unlikely that very long words occur frequently in a text, then building the statistics for *all* the substrings can be a waste of resources. Pruning the tree speeds up the implementation considerably and saves large amounts of memory. Pruning the tree does not mean that we could completely miss the word involved in a long substitution. If the current best substitution is a word  $w$  longer than the threshold  $l$ , then the encoder will eventually

**Table 10**

Comparing the Performance of OFF-LINE<sub>1</sub> for Different Choices of the Maximum Allowed Length of a Candidate for Substitution. We Fixed  $\text{min\_occ} = 4$ ,  $l = 4$ ,  $Q = 10$

$l$	paper2		mito	
	size	$\text{time}_{[\text{min}]}$	size	$\text{time}_{[\text{min}]}$
10	30,986	2.58	17,044	0.29
50	30,664	2.62	16,491	1.32
100	30,636	2.68	16,470	1.38
$\infty$	30,636	19.39	16,470	10.34

choose some substring of  $w$  of length  $l$  because that substring occurs without overlap at least as many times as  $w$ . Table 10 shows that the pruned version of OFF-LINE<sub>1</sub> at  $l = 100$  performs almost ten times faster and achieves exactly the same compression as the version that builds the complete tree.

The collective speed-up gained from these heuristics combined is significant: Our original implementation took several hours to compress those files while, afterwards, it would complete in a few minutes. What is even better, the corresponding loss of efficiency in terms of compression is almost negligible.

As documented in some additional tables, a few hundred iterations of the word selection loop of OFF-LINE suffice on inputs of the order of 100 000 symbols. This suggests that dedicated fine-grained parallel architectures of this kind would implement virtually instantaneous encoders for biosequences and general inputs alike. Tables 11 and 12 show the modest number of iterations of the main loop performed by OFF-LINE on our inputs, which would be negligible in a parallel context. Therefore, the most expensive tasks, represented by the tree constructions, can be limited considerably in a parallel implementation, turning the method into an on-line, even real-time, application.

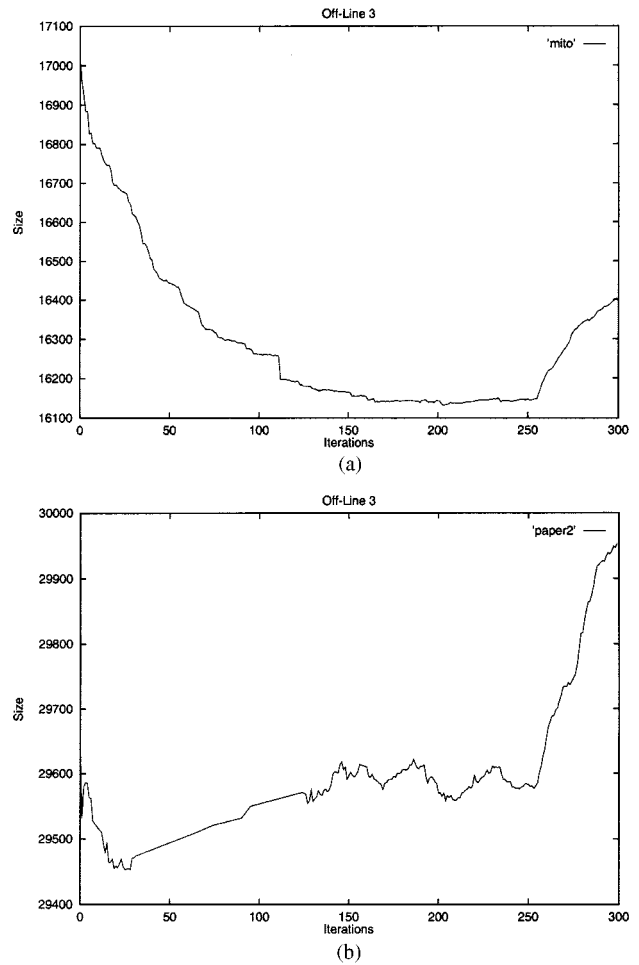
**Table 11**  
Iterations of the Main Loop of OFF-LINE for the Calgary Corpus Files

<i>File</i>	<i>Size</i>	OFF-LINE <sub>1</sub>	OFF-LINE <sub>2</sub>	OFF-LINE <sub>3</sub>
<b>bib</b>	111,261	504	634	465
<b>book1</b>	768,771	2997	2857	2990
<b>book2</b>	610,856	2305	2408	2378
<b>geo</b>	102,400	407	473	503
<b>news</b>	377,109	1789	1634	1619
<b>obj1</b>	21,504	125	111	337
<b>obj2</b>	246,814	1219	1207	1055
<b>paper1</b>	53,161	373	475	342
<b>paper2</b>	82,199	506	717	505
<b>pic</b>	513,216	94	125	222
<b>prog</b>	39,611	255	261	308
<b>prog1</b>	71,646	312	267	273
<b>progp</b>	49,379	208	210	252
<b>trans</b>	93,695	340	253	318

**Table 12**  
Iterations of the Main Loop of OFF-LINE for the Chromosomes of the Yeast

<i>File</i>	<i>Size</i>	OFF-LINE <sub>1</sub>	OFF-LINE <sub>2</sub>	OFF-LINE <sub>3</sub>
<b>chrI</b>	230,195	78	603	80
<b>chrII</b>	813,137	112	474	128
<b>chrIII</b>	315,344	61	309	68
<b>chrIV</b>	1,522,191	383	1297	441
<b>chrV</b>	574,860	109	276	118
<b>chrVI</b>	270,148	22	226	30
<b>chrVII</b>	1,090,936	144	1009	162
<b>chrVIII</b>	562,638	91	264	102
<b>chrIX</b>	439,885	54	543	63
<b>chrX</b>	745,443	108	376	123
<b>chrXI</b>	666,448	49	302	58
<b>chrXII</b>	1,078,171	444	1443	499
<b>chrXIII</b>	924,430	187	706	212
<b>chrXIV</b>	784,328	24	441	72
<b>chrXV</b>	1,091,282	128	924	147
<b>chrXVI</b>	948,061	193	755	217

Since the number of iterations performed determines the size of the vocabulary, whence ultimately of pointers, this generates “quantization” phenomena in the neighborhood of certain values that play critical roles in a computer program. Fig. 8 displays the sensitivity of the current implementations to pointer encodings at the crossing of 1 byte. The two curves



**Fig. 8.** Compressed sizes of (a) *mito* and (b) *paper2* versus number of iterations of OFF-LINE<sub>3</sub>.

plot the sizes of the compressed strings *mito* and *paper2*, respectively, at all consecutive stages of the iterated substitutions performed by OFF-LINE<sub>3</sub>. Following a steady increase until iteration 256, the compression starts decreasing as soon as OFF-LINE<sub>3</sub> must employ more than 1 byte to represent a pointer. In addition to this, the erratic shape of the plot for *paper2* suggests, with its several local minima, that it is hard at run time to pin down precisely the best moment when to stop the iterations.

## VII. CONCLUDING REMARKS

We have presented a small battery of compressors that perform well on all data but especially well on biological data. The basic paradigm is uncluttered, relatively easy to program, and acceptably fast in comparison to *ad hoc*, considerably slower and more involved methods.

Besides the obvious challenge of developing versions specifically tailored to biological sequence data, a number of interesting questions emerged in the course of the experiments that would warrant additional effort. These include possible provisions for variable window sizes, better ways to approximate the gain function  $G$ , the feasibility and usefulness of reiteration of treatment following the first

S → DDC\$  
A → ba  
B → aA  
C → BA  
D → BC

Fig. 9. Hierarchical grammar produced by SEQUITUR for **abaababaabaababaababa\$**.

S → AABAABAB \$  
A → aba  
B → ba

Fig. 10. First layer of grammar produced by OFF-LINE.

application of OFF-LINE, and several issues pertaining to the computational efficiency achievable by sequential and parallel implementations. Among the latter, a prominent concern would be to devise efficient algorithms that avoid building the statistical index from scratch after each word selection, and better storage and matching algorithms for our data structure. In fact, as documented in our tables, a few hundred iterations of the word selection loop of OFF-LINE suffice on inputs of the order of 100 000 symbols. This suggests that dedicated fine-grained parallel architectures of this kind would implement virtually instantaneous encoders for biosequences and general inputs alike.

In view of the discussion in the previous section, it is interesting for a moment to regard OFF-LINE also as a paradigm for inferring hierarchical grammatical structures in sequences. Fig. 9 displays the grammar inferred for our example string by the SEQUITUR algorithm by Nevill–Manning *et al.* [15], which is essentially patterned after an LZ parsing scheme. Except for the one involving the start symbol  $S$ , productions are constrained to have right-hand sides consisting of digrams. A grammar subtended by the strings of Fig. 7 is shown in Fig. 10. Iteration of the treatment would expose productions of the form  $C \rightarrow AAB$  and  $D \rightarrow AB$ , and, finally,  $S \rightarrow CCD$ .

The rationale to build grammar based on some measure of compression can be justified by the “Occam’s razor” principle. Occam’s razor is the *principle of parsimony* in experimental sciences. In machine learning, its expressed goal is to discover the simplest hypothesis (or model) that is consistent with the training data. In this context, the grammar that our encoder is looking for is the shortest “explanation” of the original string in terms of information content.

#### ACKNOWLEDGMENT

The authors are thankful to E. Rivals, J. Storer, and F. Tahiri for helpful discussions.

#### REFERENCES

[1] T. C. Bell, J. G. Cleary, and I. H. Witten, *Text Compression*. Englewood Cliffs, NJ: Prentice-Hall, 1990.  
[2] J. A. Storer, *Data Compression: Methods and Theory*. Rockville, MD: Computer Science, 1988.

[3] J. A. Storer and T. G. Szymanski, “Data compression via textual substitution,” *J. ACM*, vol. 29, pp. 928–951, Oct. 1982.  
[4] A. Apostolico and A. Fraenkel, “Robust transmission of unbounded strings using Fibonacci representations,” *IEEE Trans. Inform. Theory*, vol. 33, no. 2, pp. 238–245, 1987.  
[5] S. Even and M. Rodeh, “Economical encoding of commas between strings,” *Commun. ACM*, vol. 21, pp. 315–317, Apr. 1978.  
[6] P. Elias, “Universal codeword sets and representations of the integers,” *IEEE Trans. Inform. Theory*, vol. IT-21, pp. 194–202, Mar. 1975.  
[7] A. Lempel and J. Ziv, “On the complexity of finite sequences,” *IEEE Trans. Inform. Theory*, vol. IT-22, pp. 75–81, Jan. 1976.  
[8] J. Ziv and A. Lempel, “A universal algorithm for sequential data compression,” *IEEE Trans. Inform. Theory*, vol. IT-23, p. 337, May 1977.  
[9] —, “Compression of individual sequences via variable-rate coding,” *IEEE Trans. Inform. Theory*, vol. IT-24, Sept. 1978.  
[10] M. Crochemore and W. Rytter, “Efficient parallel algorithms to test square-freeness and factorize strings,” *Inform. Process. Lett.*, vol. 38, pp. 57–60, Apr. 1991.  
[11] L. M. Stauffer and D. S. Hirschberg, “PRAM algorithms for static dictionary compression,” in *Proc. 8th Int. Symp. Parallel Processing*, H. J. Siegel, Ed. Los Alamitos, CA: IEEE Computer Society Press, Apr. 1994, pp. 344–348.  
[12] S. De Agostino and J. A. Storer, “On-line versus off-line computation in dynamic text compression,” *Inform. Process. Lett.*, vol. 59, no. 3, pp. 169–174, 1996.  
[13] K. S. Fu and T. L. Booth, “Grammatical inference: Introduction and survey—Part I,” *IEEE Trans. Syst., Man, Cybern.*, vol. 5, pp. 95–111, 1975.  
[14] —, “Grammatical inference: Introduction and survey—Part II,” *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-5, no. 1, pp. 112–127, 1975.  
[15] C. Nevill-Manning, C. Ian H., C. Witten, and D. Mauksby, “Compression by induction of hierarchical grammars,” in *Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT: IEEE Computer Society Press, TCC, 1994, pp. 244–253.  
[16] A. Apostolico and E. Guerrieri, “Linear time universal compression techniques based on pattern matching,” in *Proc. 21st Allerton Conf. Communication, Control and Computing*. Monticello, IL: Univ. Illinois Press, 1983, pp. 70–79.  
[17] E. R. Fiala and D. H. Greene, “Data compression with finite windows,” *Commun. ACM*, vol. 32, pp. 490–505, Apr. 1989.  
[18] R. N. Horspool, “The effect of nongreedy parsing in Ziv–Lempel compression methods,” in *Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT: IEEE Computer Society Press, TCC, 1995, pp. 302–311.  
[19] M. Frank Rubin, “Experiments in text file compression,” *Commun. ACM*, vol. 19, pp. 617–623, Nov. 1976.  
[20] N. J. Larsson and N. J. Alistair Moffat, “Offline dictionary-based compression,” in *Proc. IEEE Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT, Mar. 1999, pp. 296–305.  
[21] J. Bentley and D. McIlroy, “Data compression using long common strings,” in *Proc. IEEE Data Compression Conf.*, Mar. 1999, pp. 287–295.  
[22] L. Gatlin, *Information Theory and the Living Systems*. New York: Columbia Univ. Press, 1972.  
[23] L. Allison and C. N. Yee, “Minimum message length encoding and the comparison of macro-molecules,” *Bull. Math. Biol.*, vol. 52, no. 3, pp. 431–453, 1990.  
[24] M. Farach, M. Noordewier, S. Savari, L. Shepp, A. Wyner, and J. Ziv, “On the entropy of DNA: Algorithms and measurements based on memory and rapid convergence,” in *ACM-SIAM Annu. Symp. Discrete Algorithms*, San Francisco, CA, Jan. 22–24, 1995, pp. 48–57.  
[25] S. Grumbach and F. Tahiri, “Compression of DNA sequences,” in *Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT: IEEE Computer Society Press, TCC, 1993, pp. 340–350.  
[26] —, “A new challenge for compression algorithms: Genetic sequences,” *Inf. Process. Manage.*, vol. 30, no. 6, pp. 875–886, 1994.

- [27] D. M. Loewenstern and P. N. Yianilos, "Significant lower entropy estimates for natural DNA sequences," in *Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT: IEEE Computer Society Press, TCC, 1997, pp. 151–160.
- [28] D. M. Loewenstern, H. M. Berman, and H. Hirsch, "Maximum a posteriori classification of DNA structure from sequence information," presented at the Pacific Symp. Biotech., Jan. 1998.
- [29] D. M. Loewenstern, H. Hirsh, P. Yianilos, and M. Noordewier, "DNA sequence classification using compression-based induction," DIMACS, Tech. Rep. 95-04, Apr. 1995.
- [30] A. Milosavljevic and J. Jurka, "Discovery by minimal length encoding: A case study in molecular evolution.," *Mach. Learn.*, vol. 12, no. 1, 2, 3, pp. 69–87, 1993.
- [31] E. Rivals, J. P. Delahaye, M. Dauchet, and O. Delgrange, "A guaranteed compression scheme for repetitive DNA sequences," in *Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT: IEEE Computer Society Press, TCC, 1996, p. 453.
- [32] E. Rivals, O. Delgrange, J. P. Delahaye, M. Dauchet, M. O. Delorme, A. Henaut, and E. Ollivier, "Detection of significant patterns by compression algorithms: The case of approximate tandem repeats in DNA sequences," *CABIOS*, vol. 13, no. 2, pp. 131–136, 1997.
- [33] A. Apostolico and W. Szpankowski, "Self-alignment in words and their applications," *J. Algorithms*, vol. 13, no. 3, pp. 446–467, 1992.
- [34] E. Ukkonen, "On-line construction of suffix trees," *Algorithmica*, vol. 14, no. 3, pp. 249–260, 1995.
- [35] A. Apostolico and F. P. Preparata, "Data structures and algorithms for the strings statistics problem," *Algorithmica*, vol. 15, pp. 481–494, May 1996.
- [36] A. S. Fraenkel and J. Simpson, "How many squares can a string contain," *J. Combin. Theory Ser. A*, vol. 82, pp. 112–120, 1998.
- [37] S. Kurtz, "Reducing the space requirements of suffix trees," *Softw.-Pract. Exp.*, vol. 29, no. 13, pp. 1149–1171, 1999.
- [38] N. J. Larsson, "Extended application of suffix trees to data compression," in *Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT: IEEE Computer Society Press, TCC, 1996, pp. 190–199.
- [39] M. Gu, M. Farach, and R. Beigel, "An efficient algorithm for dynamic text indexing," in *Proc. 5th Annu. ACM-SIAM Symp. Discrete Algorithms*, Arlington, VA, 1994, pp. 697–704.
- [40] D. R. Musser and A. A. Stepanov, "Algorithm-oriented generic libraries," *Softw.-Pract. Exp.*, vol. 24, pp. 623–642, July 1984.
- [41] M. Burrows and D. J. Wheeler, "A block-sorting lossless data compression algorithm," Digital Equipment Corporation, Tech. Rep. 124, May 1994.
- [42] C. Nevill-Manning and I. H. Witten, "Protein is incompressible," in *Proc. IEEE Data Compression Conf.*, J. A. Storer and M. Cohn, Eds. Snowbird, UT, Mar. 1999, pp. 257–266.
- [43] L. Allison, T. Edgoose, and T. I. Dix, "Compression of strings with approximate repeats," in *Intell. Syst. Mol. Biol.* '98, June 1998, pp. 8–16.
- [44] S. Lonardi, "Off-line data compression by textual substitution," Ph.D. dissertation, Dipartimento di Ingegneria Elettronica e Informatica, Univ. Padova, Feb. 1999.



**Alberto Apostolico** (Senior Member, IEEE) received the Dr. Eng. degree (*cum laude*) in electronics engineering from the University of Naples, Naples, Italy, in 1973, and the diploma of specialization in computer sciences (*cum laude*) from the University of Salerno, Salerno, Italy, in 1976.

A Fulbright Scholar in 1974–1975 at Carnegie-Mellon University, Pittsburgh, PA, he has visited extensively and held positions at institutions in the United States (University of

Illinois at Urbana-Champaign, Rensselaer Polytechnic Institute, Purdue University) and Europe (IASI-CNR in Rome, University of Paris, University of London, King's College London). Prior to joining Purdue University, West Lafayette, IN, in 1983 as an Associate and then Full Professor, he was an Associate Professor with the Department of Computer Science, University of Salerno, Italy. In 1987, he became a Professor of Computer Science in the Italian, University of Padova, Italy, in 1992, where he holds the Chair of Theoretical Computer Science in the School of Engineering. His main research interests are in the area of analysis and design of algorithms. His recent work focuses on pattern matching algorithms and applications, notably, to string searching and comparison, the subject of more than 90 published papers in the major international journals and conferences and several coedited books. He serves on the editorial boards of *Theoretical Computer Science*, *Parallel Processing Letters*, *Journal of Computational Biology*, *Chaos Theory and Applications*, and was a Guest Editor of a special issue of *Algorithmica* devoted to String Algorithmics and its Applications. He has been keynote speaker or lecturer at more than 25 international conferences and advanced schools, and on the program committees of as many international conferences, advanced schools, and workshops.

Dr. Apostolico is the recipient or co-recipient of U.S., French, British, Italian, and International (Fulbright, NATO, ESPRIT) research grants. He also served as referee/reviewer for most TCS journals and major conferences: NSF, Canadian SERC, NATO, Israel Science Council, the Hong Kong and the Finnish Sci. Acad., among others.



**Stefano Lonardi** (Student Member, IEEE) received the "Laurea" degree (*cum laude*) from the University of Pisa, Pisa, Italy, in 1994. He is currently a Ph.D. candidate at the Department of Computer Sciences, Purdue University, West Lafayette, IN.

In 1996, he joined the graduate program at Purdue University. His main research interests include data compression, algorithms on strings, computational molecular biology, and statistical analysis of sequences.

In 2000, Mr. Lonardi received the Student Research Award from the Purdue Chapter of Upsilon Pi Epsilon. He is a Member of the ACM and the honor societies Upsilon Pi Epsilon and Phi Kappa Phi.