

Practical Approaches to Reduce the Space Requirement of Lempel-Ziv-Based Compressed Text Indices

DIEGO ARROYUELO

University of Chile

and

GONZALO NAVARRO

University of Chile

Given a text $T[1..u]$ over an alphabet of size σ , the *full-text search* problem consists in finding the *occ* occurrences of a given pattern $P[1..m]$ in T . The current trend in indexed text searching is that of *compressed full-text self-indices*, which replace the text with a space-efficient representation of it, while at the same time providing indexed access to the text.

The LZ-index of Navarro is a compressed full-text self-index based on the LZ78 compression algorithm. This index requires about 4 times the size of the compressed text, i.e. $4uH_k(T) + o(u \log \sigma)$ bits of space, where $H_k(T)$ is the k -th order empirical entropy of text T . This index has shown to be very competitive in practice for locating pattern occurrences and extracting text snippets. However, the LZ-index is larger than competing schemes, and does not offer space/time tuning options, which limits its applicability in many practical scenarios.

In this paper we study several ways to reduce the space of LZ-index, from a practical point of view and in different application scenarios. The main idea used to reduce the space is to regard the original index as a *navigation scheme* that allows us moving between index components. Then we perform an abstract optimization on this scheme, defining alternative schemes that support the same navigation, yet reducing the original redundancy. We obtain reduced LZ-indices requiring $3uH_k(T) + o(u \log \sigma)$ and $(2 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits of space, for any $0 < \epsilon < 1$. Our LZ-indices have an average locating time of $O(m^2 + \frac{n}{\sigma^{m/2}})$, which is $O(m^2)$ for $m \geq 2 \log_\sigma u$.

We perform extensive experimentation to show that our developments lead to reduced LZ-indices that are competitive with the state of the art in many practical situations, providing interesting space/time trade-offs and allowing in many cases to replace the original LZ-index with a smaller, yet competitive, representation. Given the space that our indices require, they are in most cases the best alternative for the key operations of extracting arbitrary text substrings, as well as searching and then displaying the contexts surrounding the pattern occurrences.

Categories and Subject Descriptors: [] ; [] ; []:

General Terms:

Additional Key Words and Phrases:

1. INTRODUCTION

With the rapid growing of nowadays text repositories, and the need of searching those texts for patterns of interest, the classical problem of *text searching* becomes fundamental in many computer applications. Given a sequence of symbols $T[1..u]$ (the text) over an alphabet $\Sigma = \{1, \dots, \sigma\}$, and given another (short) sequence $P[1..m]$ (the *search pattern*) over Σ , the *full-text search problem* consists in finding all the *occ* occurrences of P in T . There exist three typical kinds of queries, which arise in different types of applications, namely:

- Existential queries*: Operation `exists(P)` tells us whether pattern P exists in T or not.
- Cardinality queries*: Operation `count(P)` counts the number of occurrences of pattern P in T .
- Locating queries*: Operation `locate(P)` reports the starting position of the *occ* occurrences of pattern P in T .

In the case of `locate` queries, one can also be interested in finding only a certain number of pattern occurrences, such that the remaining ones are located while processing the first ones, or on user demand (think for example of the 10 result pages shown by a Web search engine). In other cases we can just need, for example, to find a fixed number K of (arbitrary) occurrences. Let us think again of Web search engines, where the answer to a user-entered query usually consists of a ranking with the (hopefully) most relevant Web pages for the query, plus an arbitrary (short) context where the pattern occurs within every such page; this is displayed in order to help users decide whether a page is relevant or not for their needs. Therefore, in this example we must be able to quickly find just one arbitrary pattern occurrence in the text and display its context. We call *partial locate* queries those that ask to locate K arbitrary occurrences.

Unlike *word-based* text searching, which is the typical scenario in Information Retrieval, we wish to find any *text substring*, not only whole words or phrases. This has applications in texts where the concept of *word* is not clearly defined (e.g. Oriental languages, program code, etc.), or texts where words do not exist at all (e.g., DNA, protein, and MIDI pitch sequences, etc.).

We assume that the text is large and known in advance to queries, and we need to perform several queries on it. Therefore, we can preprocess the text to construct an *index* on it, which is a data structure allowing efficient access to the pattern occurrences, yet increasing the space requirement.

Though classical full-text indices, like *suffix trees* [Apostolico 1985] and *suffix arrays* [Manber and Myers 1993], are very efficient at search time, they have the problem of a high space requirement: they require $O(u \log u)$ and $u \log u$ bits respectively, which in practice is about 10–20 and 4 times the text size respectively, apart from the text itself. Thus, we can have large texts which fit into main memory, but whose corresponding suffix tree (or array) cannot be accommodated in main memory. Using secondary storage for the indices is several orders of magnitude slower, so one looks for ways to reduce their size, with the main motivation of maintaining the indices of very large texts entirely in main memory. Therefore, we seek *to provide fast access to the text using as little space as possible*. The modern trend is to use the compressibility of the text to reduce the space of the index. In recent years there has been much research on *compressed text databases*, focusing on techniques to represent the text and the index using little space, yet permitting efficient text searching [Navarro and Mäkinen 2007].

1.1 Compressed Full-Text Self-Indexing

The track of *compressed full-text self-indices* started some years ago [Grossi and Vitter 2000; Ferragina and Manzini 2000; Sadakane 2000], with the aim of accommodating the indices of very large texts entirely in main memory. A *full-text*

self-index allows one to retrieve any part of the text without storing the text itself, and in addition provides search capabilities on the text. A *full-text compressed self-index* is one whose space requirement is proportional to the compressed text size, e.g. $O(uH_k(T))$ bits, where $H_k(T)$ denotes the k -th order empirical entropy of T [Manzini 2001], which is a standard measure of the compressibility of T (see Section 2.1 for more details). Then a compressed full-text self-index *replaces* the text with a more space-efficient representation of it, which at the same time provides indexed access to the text [Navarro and Mäkinen 2007; Ferragina and Manzini 2005].

Taking space proportional to the compressed text, replacing it, and providing efficient indexed access to it is an unprecedented breakthrough in text indexing and compression. Thanks to the advances in these technologies, it is quite common nowadays that the indices of (very) large texts can be accommodated in the main memory of a desktop computer (e.g., the Human Genome, with about 3 gigabytes).

As compressed full-text self-indices replace the text, we are also interested in operations:

- `display(P, ℓ)`, which displays a context of ℓ symbols surrounding the *occ* occurrences of pattern P in T , and
- `extract(i, j)`, which decompresses the substring $T[i..j]$, for any text positions $i \leq j$.

Thus we can see compressed self-indices as full-text indices compressing the text, or as compressors allowing efficient text extraction and indexed full-text searching.

1.1.1 *Some Applications of Compressed Full-Text Self-Indexes.* In the scenario of compressed full-text self-indices, where we *replace the text* with a representation allowing indexed-search capabilities, being able to efficiently **extract** arbitrary text substrings is one of the most basic and important problems that indices must solve efficiently.

Despite **locate** queries are important in classical full-text indexing (since we have the text at hand to access the occurrences and their contexts, as needed by many applications), they are usually not enough for compressed self-indices, since we obtain just text positions, and no clue about the text surrounding these occurrences. In many applications the context surrounding an occurrence is as important as (and sometimes more important than) the occurrence position itself. For example, a user might be interested in the context surrounding the occurrences to decide whether the answer is interesting or not; think, for example, of the widely-used Unix search tool **grep**, which by default shows the text lines containing the occurrences. The relevance of this information is witnessed by the fact that most modern Web search engines display, along with the answers to a query, a context surrounding a pattern occurrence within each document. Therefore, in our scenario it is usually more important to obtain the contexts surrounding the pattern occurrences (i.e., **display** queries), than just text positions (i.e., **locate** queries). The latter can be interesting in specific cases, for example if one wants to take statistics about the positions of the occurrences (for instance, the average difference between successive occurrences, e.g. for linguistic or data mining applications), but **display** queries are more frequently used in general.

Finally, `count` and `exists` queries have much more specific applications, and they usually conform the internal machinery of more complex tasks. For example, for *approximate pattern matching* and *text categorization* (where a document is assigned a class or category depending on the frequency of appearance of given keywords). However, this is not enough for many other applications. For example, after categorizing the documents, an user might want to see a document, or a portion of it, so we have to ask the index to reproduce it, or may want to search for a given pattern and display the occurrence contexts in order to decide whether the document is of interest or not. Some *pattern discovery* tasks may use the frequency of certain strings to decide that they are important patterns. Another example is *selective dissemination of information*, where user profiles are formed by keywords of interest and the system is interested in the presence or absence of those keywords to send or not the document to the user.

1.1.2 *Families of Compressed Full-Text Self-Indices.* The main families of compressed self-indices [Navarro and Mäkinen 2007] are *Compressed Suffix Arrays* [Grossi and Vitter 2005; Sadakane 2003; Grossi et al. 2003], indices based on *backward search* [Ferragina and Manzini 2005; Mäkinen and Navarro 2005; Ferragina et al. 2007] (which are alternative ways to compress suffix arrays, and known as the *FM-index* family), and the indices based on the *Lempel-Ziv* compression algorithm [Ziv and Lempel 1978] (*LZ-indices* for short) [Kärkkäinen and Ukkonen 1996; Navarro 2004; Ferragina and Manzini 2005; Russo and Oliveira 2008].

In this paper we are interested in LZ-indices, since they have shown to be effective in practice for extracting text, displaying occurrence contexts, and locating the occurrences, outperforming other compressed indices [Navarro 2004; 2008]. What characterizes the particular niche of LZ-indices is the $O(uH_k(T))$ space combined with $O(\log u)$ theoretical worst-case time per located occurrence. Moreover, in practice many pattern occurrences can be actually found in constant time per occurrence, which makes the LZ-indices competitive, for example, when searching for short patterns.

Navarro’s LZ-index [Navarro 2004; 2008] is a compressed full-text self-index based on the Lempel-Ziv 1978 [Ziv and Lempel 1978] (LZ78 for short) parsing of the text. See Section 2.2 for a description of the LZ78 compression algorithm. The LZ-index takes about 4 times the size of the compressed text, i.e. $4uH_k(T) + o(u \log \sigma)$ bits, for any $k = o(\log_\sigma u)$ [Kosaraju and Manzini 1999; Ferragina and Manzini 2005], and answers queries in $O(m^3 \log \sigma + (m + occ) \log u)$ worst-case time. The index also replaces the text (i.e., it is a *self-index*): it can display a text context of length ℓ around an occurrence found (and in fact any sequence of LZ78 phrases) in $O(\ell \log \sigma)$ time, or obtain the whole text in time $O(u \log \sigma)$. The index is built in $O(u \log \sigma)$ time.

However, in practice the space requirement of LZ-index is relatively large compared with competing schemes: 1.2–1.6 times the text size versus 0.6–0.7 and 0.3–0.8 times the text size of *Compressed Suffix Array* (CSA) [Sadakane 2003] and *FM-index* [Ferragina and Manzini 2005], respectively. In addition, the LZ-index does not offer space/time trade-offs, which limits its applicability. Yet, the LZ-index is faster to locate and to display the context of an occurrence. As we said before, fast displaying of text substrings is very important in self-indices, as the

text is not available otherwise. Thus, in this paper the challenge is to reduce the space requirement of LZ-index (providing space/time trade-offs), while retaining the good features of fast locating and fast displaying of text substrings.

The practical implementation and evaluation of compressed self-indices is a very important issue, since many indices are proposed in theory but never implemented. In this respect, the *Pizza&Chilli Corpus* [Ferragina and Navarro 2005] provides practical implementations of compressed indices, as well as some testbed texts. Following this line, in this paper we will design new compressed full-text self-indices from a practical point of view. Our developments are, in some extent, based on the theoretical ideas of [Arroyuelo et al. 2006], yet in this paper we consider much more practical ideas, and perform extensive experimentation.

1.2 Our Contribution

In this paper we study how to reduce the space requirement of Navarro’s LZ-index, using what we call the *navigational scheme* approach. This consists in regarding the original LZ-index data structure as a navigation scheme, which allows us moving among data representations. From this point of view we can study the redundancy introduced by the different data structures that conform the LZ-index. We design new data structures supporting the same navigation, yet reducing the original redundancy.

We define several reduced alternatives requiring $3uH_k(T) + o(u \log \sigma)$ and $(2 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits of space, for $0 < \epsilon < 1$, which is less than the space of the original LZ-index. Our indices cannot provide worst-case guarantees at search time, yet they support `locate` queries for patterns of length m in $O(\frac{m^2}{\epsilon} + \frac{n}{\epsilon \sigma^{m/2}})$ average time, which is $O(\frac{m^2}{\epsilon})$ for $m \geq 2 \log_{\sigma} u$.

We implement and test our indices in several real-life scenarios, and show that they offer very attractive space/time trade-offs in practice for the key operations `extract` and `display`: given the memory space our indices require, they are in most cases the best alternative for these operations. This is very important for compressed self-indices, since we are able to efficiently extract any text substring, as well as efficiently searching the text and obtaining the contexts surrounding the pattern occurrences.

For `locate` queries we test two alternatives: *partial locate* queries, which retrieve a fixed number of (arbitrary) pattern occurrences, as well as *full locate* queries, where we retrieve all of the pattern occurrences. In the former case, we develop a heuristic to get fast access to the first occurrences, avoiding as much as possible the navigation on the tries that compose LZ-index; this navigation becomes expensive if we report only a few occurrences. Our indices offer an interesting alternative, particularly for short patterns and small alphabets (such as DNA data). Though not the best in all cases, our indices also offer an attractive alternative for full `locate` queries, competing against the best existing indices.

We conclude also that our indices can replace the original LZ-index without losing competitiveness. We also demonstrate a very important aspect of our indices: compressed indices based on suffix arrays store extra non-compressible information to efficiently carry out the `locate` and `display` tasks, whereas the extra data stored by LZ-indices is compressible. Therefore, when the texts are highly compressible,

LZ-indices can be smaller and faster than alternative indices; and in other cases they offer attractive space/time trade-offs.

Just as for the original LZ-index, our indices are not competitive for `count` queries, since they must basically locate the occurrences in order to count them. However, we show that our indices are an interesting alternative for `exists` queries, which can be supported much more efficiently than `count` queries. The trick here is to find the first occurrence as fast as possible. We use the novel heuristic defined for partial `locate` queries to quickly test the existence of a pattern in the text, which shows to be effective in practice.

It is also well known that the original LZ-index is more adequate when we need to search for short patterns [Navarro 2008]. This is mainly because we must perform $O(m^2)$ navigation operations on the tries that conform the LZ-index, which in case of long patterns become predominant. We implement a version of our indices based on the *Depth-First Unary Degree Sequence* (DFUDS) trie representation [Benoit et al. 2005]. We adapt this representation for our specific purposes, trading some theoretical (yet unpractical) aspects for more practical ones, obtaining very interesting results, especially for partial `locate` queries where trie navigations are fundamental.

2. BASIC CONCEPTS

2.1 Empirical Entropy

Text compression is a technique to represent a text using less space, profiting from the regularities of non-random texts. A concept related to text compression is that of the k -th order empirical entropy of a sequence of symbols T over an alphabet of size σ , denoted by $H_k(T)$ [Manzini 2001]. The value $uH_k(T)$ provides a lower bound to the number of bits needed to compress T using any compressor that encodes each symbol considering only the context of k symbols that precede it in T . It holds that $0 \leq H_k(T) \leq H_{k-1}(T) \leq \dots \leq H_0(T) \leq \log \sigma$ ($\log x$ means $\lceil \log_2 x \rceil$ in this paper). Formally, we have

Definition 2.1. Given a text $T[1..u]$ over an alphabet Σ , the zero-order empirical entropy of T is defined as

$$H_0(T) = \sum_{c \in \Sigma} \frac{n_c}{u} \log \frac{u}{n_c}$$

where n_c is the number of occurrences of symbol c in T . The sum includes only those symbols c that occur in T , so that $n_c > 0$.

Definition 2.2. Given a text $T[1..u]$ over an alphabet Σ , the k -th order empirical entropy of T is defined as

$$H_k(T) = \sum_{s \in \Sigma^k} \frac{|T^s|}{u} H_0(T^s)$$

where T^s is the subsequence of T formed by all the symbols that occur preceded by the context s . Again, we consider only contexts s that do occur in T .

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
a	l	ab	ar	_	a_	la	_a	lab	ard	a_p	ara	_ap	al	abr	ar1	a\$

Fig. 1. LZ78 phrase decomposition for the running example text $T = \text{“alabar_a_la_alabarda_para_apalabrarla”}$, and the corresponding phrase identifiers.

2.2 Ziv-Lempel Compression

The Ziv-Lempel compression algorithm of 1978 (usually named LZ78 [Ziv and Lempel 1978]) is based on a *dictionary of phrases*, in which we add every new *phrase* computed. At the beginning of the compression, the dictionary contains a single phrase b_0 of length 0 (i.e., the empty string). The current step of the compression is as follows: If we assume that a prefix $T[1..j]$ of T has been already compressed into a sequence of phrases $Z = b_1 \dots b_r$, all of them in the dictionary, then we look for the longest prefix of the rest of the text $T[j + 1..u]$ which is a phrase of the dictionary. Once we have found this phrase, say b_s of length ℓ_s , we construct a new phrase $b_{r+1} = (s, T[j + \ell_s + 1])$, write the pair at the end of the compressed file Z , i.e. $Z = b_1 \dots b_r b_{r+1}$, and add the phrase to the dictionary.

We will call B_i the string represented by phrase b_i , thus $B_{r+1} = B_s T[j + \ell_s + 1]$. In the rest of the paper we assume that the text T has been compressed using the LZ78 algorithm into $n + 1$ phrases, $T = B_0 \dots B_n$, such that $B_0 = \varepsilon$ (the empty string). We say that i is the *phrase identifier* corresponding to B_i , for $0 \leq i \leq n$.

Property 2.3. For all $1 \leq t \leq n$, there exists $\ell < t$ and $c \in \Sigma$ such that $B_t = B_\ell \cdot c$.

That is, every phrase B_t (except B_0) is formed by a previous phrase B_ℓ plus a symbol c at the end. This implies that the set of phrases is *prefix closed*, meaning that any prefix of a phrase B_t is also an element of the dictionary. Therefore, a natural way to represent the set of strings B_0, \dots, B_n is a trie, which we call *LZTrie*. In Fig. 1 we show the LZ78 phrase decomposition for the text $T = \text{“alabar_a_la_alabarda_para_apalabrarla”}$, which will be our running example. We show phrase identifiers above the corresponding phrase in the parsing.

In Fig. 4(a) we show the corresponding *LZTrie*. Inside each *LZTrie* node we show the corresponding phrase identifier.

Property 2.4. Every phrase B_i , $0 \leq i < n$, represents a different text substring.

This property is used in the LZ-index search algorithm (see Section 3). The only exception to this property is the last phrase B_n . We deal with the exception by appending to T a special symbol “\$” $\notin \Sigma$, assumed to be smaller than any other symbol in the alphabet. The last phrase will contain this symbol and thus will be unique too.

Definition 2.5. Let $b_r = (r_1, c_1)$, $b_{r_1} = (r_2, c_2)$, $b_{r_2} = (r_3, c_3)$, and so on until $r_k = 0$ be phrases of the LZ78 parsing of T . The sequence of phrase identifiers r, r_1, r_2, \dots is called the *referencing chain* starting at phrase r .

The referencing chain starting at phrase r reproduces the way phrase b_r is formed from previous phrases and it is obtained by successively moving to the parent in the *LZTrie*. For example, the referencing chain of phrase 9 in Fig. 4(a) is $r = 9$, $r_1 = 7$, $r_2 = 2$, and $r_3 = 0$.

The compression algorithm is $O(u)$ time in the worst case and efficient in practice provided we use the *LZTrie*, which allows rapid searching of the new text prefix (for each symbol of T we move once in the trie). The decompression needs to build the same dictionary (the pair that defines the phrase r is read at the r -th step of the algorithm).

Property 2.6 [Ziv and Lempel 1978]. It holds that $\sqrt{u} \leq n \leq \frac{u}{\log_\sigma u}$. Thus, $n \log u \leq u \log \sigma$ always holds.

LEMMA 2.7 [KOSARAJU AND MANZINI 1999]. *It holds that $n \log n = uH_k(T) + O(u \frac{1+k \log \sigma}{\log_\sigma u})$ for any k .*

In our work we assume $k = o(\log_\sigma u)$ (and hence $\log \sigma = o(\log u)$ to allow for $k > 0$); therefore, $n \log n = uH_k(T) + o(u \log \sigma)$.

2.3 Succinct Representations of Sequences and Permutations

A *succinct data structure* requires space close to the information-theoretic lower bound, while supporting the corresponding operations efficiently. In this section and the next we review some results on succinct data structures, which are necessary to follow our work.

2.3.1 Data Structures for rank and select. Given a bit vector $\mathcal{B}[1..n]$, we define the operation $rank_0(\mathcal{B}, i)$ (similarly $rank_1$) as the number of 0s (1s) occurring up to the i -th position of \mathcal{B} . The operation $select_0(\mathcal{B}, i)$ (similarly $select_1$) is defined as the position of the i -th 0 (i -th 1) in \mathcal{B} . We assume that $select_0(\mathcal{B}, 0)$ always equals 0 (similarly for $select_1$). These operations can be supported in constant time and requiring $n + o(n)$ bits [Munro 1996], or even $nH_0(\mathcal{B}) + o(n)$ bits (including \mathcal{B} itself) [Raman et al. 2002].

There exist a number of practical data structures supporting *rank* and *select*, like the one by González et al. [2005], Kim et al. [2005], Okanohara and Sadakane [2007], etc. Among these, the index of González et al. [2005] is very (perhaps the most) efficient in practice to compute *rank*, requiring little space on top of the sequence itself. *Select* is implemented by binary searching the directory built for operation *rank*, and thus without requiring any extra space for that operation (yet, the time for *select* becomes $O(\log n)$ in the worst case).

Given a sequence $S[1..u]$ over an alphabet Σ , we generalize the above definition for $rank_c(S, i)$ and $select_c(S, i)$ for any $c \in \Sigma$. If $\sigma = O(\text{polylog}(u))$, the solution of Ferragina et al. [2007] allows one to compute both $rank_c$ and $select_c$ in constant time and requiring $uH_0(S) + o(u)$ bits of space. Otherwise the time is $O(\frac{\log \sigma}{\log \log u})$ and the space is $uH_0(S) + o(u \log \sigma)$ bits. The representation of Golynski et al. [2006] requires $n(\log \sigma + o(\log \sigma)) = O(n \log \sigma)$ bits of space [Barbay et al. 2007], allowing us to compute $select_c$ in $O(1)$ time, and $rank_c$ and access to $S[i]$ in $O(\log \log \sigma)$ time.

2.3.2 Succinct Representation of Permutations. The problem here is to represent a permutation π of $\{1, \dots, n\}$, such that we can compute both $\pi(i)$ and its inverse $\pi^{-1}(j)$ in constant time and using as little space as possible. A natural representation for π is to store the values $\pi(i)$, $i = 1, \dots, n$, in an array of $n \log n$ bits. The brute-force solution to the problem computes $\pi^{-1}(j)$ looking for j se-

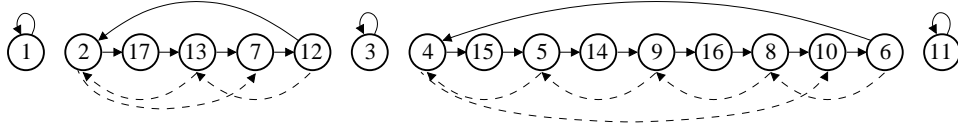


Fig. 2. Cycle representation of permutation ids of Fig. 3(a). Each solid arrow $i \rightarrow j$ in the diagram means $ids(i) = j$. Dashed arrows represent backward pointers.

quentially in the array representing π . If j is stored at position i , i.e. $\pi(i) = j$, then $\pi^{-1}(j) = i$. Although this solution does not require any extra space to compute π^{-1} , it takes $O(n)$ time in the worst case.

A much more efficient solution is based on the cycle notation of a permutation. The cycle for the i -th element of π is formed by elements $i, \pi(i), \pi(\pi(i))$, and so on until the value i is found again. It is important to note that every element occurs in one and only one cycle of π . For example, the cycle notation for permutation ids of Fig. 3(a) is shown in Fig. 2. So, to compute $\pi^{-1}(j)$, instead of looking sequentially for j in π , we only need to look for j in its cycle: $\pi^{-1}(j)$ is just the value “pointing” to j in the diagram of Fig. 2. To compute $ids^{-1}(13)$ in the previous example, we start at position 13, then move to position $ids(13) = 7$, then to position $ids(7) = 12$, then to $ids(12) = 2$, then to $ids(2) = 17$, and as $ids(17) = 13$ we conclude that $ids^{-1}(13) = 17$. The only problem here is that there are no bounds for the size of a cycle, hence this algorithm takes also $O(n)$ time in the worst case. However, it can be improved for a more efficient computation of $\pi^{-1}(j)$.

Given $0 < \epsilon < 1$, we create subcycles of size $O(1/\epsilon)$ by adding a *backward pointer* out of $O(1/\epsilon)$ elements in each cycle of π . Dashed arrows in Fig. 2 show backward pointers for $1/\epsilon = 2$. To compute $ids^{-1}(17)$, we first move to $ids(17) = 13$; as 13 has a backward pointer we follow it and hence we move to position 2. Then, as $ids(2) = 17$ we conclude that $ids^{-1}(17) = 2$, in $O(1/\epsilon)$ worst-case time. We store the backward pointers compactly in an array of $\epsilon n \log n$ bits. We mark the elements having a backward pointer by using a bit vector supporting *rank* queries, which also help us to find the backward pointer corresponding to a given element (see the article by Munro et al. [2003] for details). Overall, this solution requires $(1 + \epsilon)n \log n + n + o(n)$ bits of storage.

2.4 Succinct Representation of Trees

Given a tree with n nodes, there exist a number of succinct representations requiring $2n + o(n)$ bits, which is close to the information-theoretic lower bound of $2n - \Theta(\log n)$ bits. We explain the representations that we will need in our work.

2.4.1 Balanced Parentheses. The problem of representing a sequence of balanced parentheses is highly related to the succinct representation of trees [Munro and Raman 2001]. Given a sequence par of $2n$ balanced parentheses, we want to support the following operations on par :

- $findclose(par, i)$, which given an opening parenthesis at position i , finds the position of the matching closing parenthesis;
- $findopen(par, j)$, which given a closing parenthesis at position j , finds the position of the matching opening parenthesis;

- $excess(par, i)$, which yields the difference between the number of opening and closing parentheses up to position i in the parentheses sequence; and
- $enclose(par, i)$, which given a parentheses pair whose opening parenthesis is at position i , yields the position of the opening parenthesis corresponding to the closest matching parentheses pair enclosing the one at position i .

Munro and Raman [2001] show how to compute all these operations in constant time and requiring $2n + o(n)$ bits of space. They also show one of the main applications of maintaining a sequence of balanced parentheses: the succinct representation of general trees. Among the practical alternatives, we have the representation of Geary et al. [2006] and the one by Navarro [2004, Section 6.1]. The latter has shown to be very effective for representing LZ-indices, and therefore we briefly review it in which follows.

Navarro’s Practical Representation of Balanced Parentheses. To support the operations we could simply precompute and store all the possible answers, requiring $O(n \log n)$ bits overall. However, in many applications (e.g., the representation of trees) matching opening and closing parentheses tend to be close to each other. Profiting from this property, and for instance to support $findclose$, Navarro uses a brute-force approach for these parentheses, sequentially looking for the closing parenthesis within the next few, say 32, parentheses. Actually, this search is performed by using precomputed tables to avoid a bit-per-bit scan.

If the answer cannot be found in this way, Navarro searches a hash table storing the answers for parentheses that are not so close, though not so far away from each other. Say, for example, matching parentheses with a difference of up to 256 positions (parentheses). Instead of storing absolute positions, the difference between positions is stored, and thus we can use 8 bits to code these numbers, which saves space. Finally, if the answer cannot be found in the previous hash table, another table is searched for matching parentheses that are far away from each other (here full numbers are stored, but there are hopefully few entries). A similar approach is used to compute $enclose$ and $findopen$ operations.

The parentheses operations are supported in $O(\log \log n)$ average time [Navarro 2008]. However, this representation does not provide theoretical worst-case guarantees in the space requirement, since in the worst case almost every opening parenthesis has its matching parenthesis far away, so we have to store its information in the tables. Fortunately these cases are not common in practice.

To compute operation $excess(i)$, we need to support operation $rank$ over the binary sequence of parentheses, since $excess(i) \equiv rank_i(par, i) - rank_i(par, i)$. We use the representation of González et al. [2005] to efficiently support $rank$ and $select$ (which will be needed later) on par , while requiring little space on top of the sequence.

2.4.2 BP Representation of Trees. The *balanced parentheses* (BP) representation of a tree defined by Munro and Raman [2001] is built from a depth-first preorder traversal of the tree, writing an *opening* parenthesis when arriving to a node for the first time, and a *closing* parenthesis when going up (after traversing the subtree of the node). In this way, we get a sequence of balanced parentheses, where each node is represented by a pair of opening and closing parentheses. We identify a tree

2.4.3 *DFUDS Representation of Trees.* To get this representation [Benoit et al. 2005] we perform a preorder traversal on the tree, and for every node reached we write its degree in unary using parentheses. For example, a node of degree 3 reads ‘((()’ under this representation. What we get is almost a balanced parentheses representation: we only need to add a fictitious ‘(’ at the beginning of the sequence. A node of degree d is identified by the position of the first of the $d + 1$ parentheses representing the node.

This representation requires also $2n + o(n)$ bits, and supports operations $parent(x)$, $subtreesize(x)$, $degree(x)$ (which gets the degree, i.e., the number of children, of node x), $childrank(x)$ (which gets the rank of node x within its siblings [Jansson et al. 2007]), $ancestor(x, y)$, and $child(x, i)$, all in $O(1)$ time in the following way, assuming that par represents now the DFUDS sequence of the tree:

$$\begin{aligned}
 parent(x) &\equiv select_1(par, rank_1(par, findopen(par, x - 1))) + 1 \\
 child(x, i) &\equiv findclose(par, select_1(par, rank_1(par, x) + 1) - i) + 1 \\
 subtreesize(x) &\equiv (findclose(par, enclose(par, x)) - x) / 2 + 1 \\
 degree(x) &\equiv select_1(par, rank_1(par, x) + 1) - x \\
 childrank(x) &\equiv select_1(par, rank_1(par, findopen(par, x - 1)) + 1) \\
 &\quad - findopen(par, x - 1) \\
 ancestor(x, y) &\equiv x \leq y \leq findclose(par, enclose(par, x))
 \end{aligned}$$

Operation $depth(x)$ can be also computed in constant time on DFUDS by using the approach of Jansson et al. [2007], requiring $o(n)$ extra bits. It is important to note that, unlike the BP representation, DFUDS needs operation $findopen$ on the parentheses in order to compute operation $parent$ on the tree. In practice, if we build on Navarro’s parentheses data structure, this implies that DFUDS needs more space than BP since we need additional hash tables to support $findopen$.

Given a node in this representation, say at position i , its preorder position can be computed by counting the number of closing parentheses before position i ; in other words, $preorder(x) \equiv rank_1(par, x - 1)$. Given a preorder position p , the corresponding node is computed by $selectnode(p) \equiv select_1(par, p) + 1$.

Representing σ -ary Trees with DFUDS. For cardinal trees (i.e., trees where each node has at most σ children, each child labeled by a symbol in the set $\{1, \dots, \sigma\}$) we use the DFUDS sequence par plus an array $letts[1..n]$ storing the labels according to a DFUDS traversal of the tree: we traverse the tree in depth-first preorder, and every time we reach a node x we write the symbols labeling the children of x . In this way, the labels of the children of a given node are all stored contiguously in $letts$, which will allow us to compute operation $child(x, \alpha)$ (which gets the child of node x with label $\alpha \in \{1, \dots, \sigma\}$) efficiently. In Fig. 3(b) we show the DFUDS representation of *LZTrie* for our running example. Notice the inverse relation between the d opening parentheses defining x and the symbols of the children of x : the label of the i -th child is at position i within the symbols of the children of x , while the corresponding opening parenthesis is at position $(d - i + 1)$ within the definition of x . This shall mean extra work when retrieving the symbol by which a given node descends from its parent.

We support operation $child(x, \alpha)$ as follows. Suppose that node x has position

p within the DFUDS sequence par , and let $p' = rank_{\iota}(par, p) - 1$ be the position in $letts$ for the symbol of the first child of x . Let $n_{\alpha} = rank_{\alpha}(letts, p' - 1)$ be the number of α s up to position $p' - 1$ in $letts$, and let $i = select_{\alpha}(letts, n_{\alpha} + 1)$ be the position of the $(n_{\alpha} + 1)$ -th α in $letts$. If i lies within positions p' and $p' + degree(x) - 1$, then the child we are looking for is $child(x, i - p' + 1)$, which, as we said before, is computed in constant time over par ; otherwise x has not a child labeled α . We can also retrieve the symbol by which x descends from its parent with $letts[rank_{\iota}(par, parent(x)) - 1 + childrank(x) - 1]$, where the first term stands for the position in $letts$ corresponding to the first symbol of the parent of node x . The second term, $childrank(x)$, comes from the inverse relation between symbols and opening parentheses representing a node.

Thus, the time for operation $child(x, \alpha)$ depends on the representation we use for $rank_{\alpha}$ and $select_{\alpha}$ queries. Notice that $child(x, \alpha)$ could be supported in a straightforward way by binary searching the labels of the children of x , in $O(\log \sigma)$ worst-case time and not needing any extra space on top of array $letts$. The access to $letts[\cdot]$ takes constant time.

Alternatively, we can represent $letts$ with the data structure of Ferragina et al. [2007], which requires $n \log \sigma + o(n \log \sigma)$ bits of space, and allows us to compute $child(x, \alpha)$ in $O(1 + \frac{\log \sigma}{\log \log u})$ time. The access to $letts[\cdot]$ also takes $O(1 + \frac{\log \sigma}{\log \log u})$ time. These times are $O(1)$ whenever $\sigma = O(\text{polylog}(u))$ holds. On the other hand, we can use the data structure of Golynski et al. [2006], requiring $O(n \log \sigma)$ bits of space, yet allowing us to compute $child(x, \alpha)$ in $O(\log \log \sigma)$ time, and access to $letts[\cdot]$ also in $O(\log \log \sigma)$ time. We will prefer the representation of Ferragina et al., since it is able to improve its time complexity to $O(1)$ for smaller alphabets.

The scheme we have presented to represent $letts$ is slightly different to the original one [Benoit et al. 2005], which achieves $O(1)$ time for $child(x, \alpha)$ for any σ . However, ours is simpler and allows us to efficiently access $letts[\cdot]$, which will be very important in our indices to extract text substrings.

3. THE LZ-INDEX DATA STRUCTURE

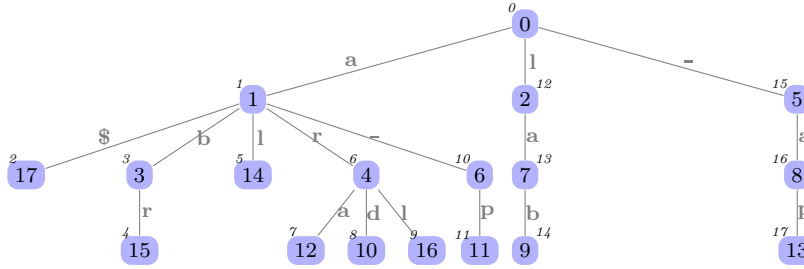
Assume that the text $T[1..u]$ has been compressed using the LZ78 algorithm into $n+1$ phrases $T = B_0 \dots B_n$, as explained in Section 2.2. As we are mainly interested in practical performance in this paper, we describe next a practical representation of the LZ-index and its corresponding search algorithm [Navarro 2008, Section 9].

Hereafter, given a string $S = s_1 \dots s_i$, we will use $S^r = s_i \dots s_1$ to denote its reverse. Also, $S^r[i..j]$ will mean $(S[i..j])^r$.

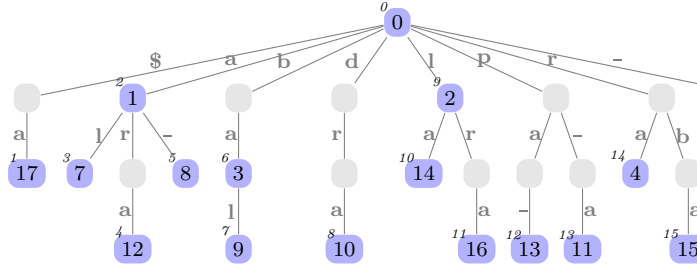
3.1 Original LZ-index Components

The following data structures conform the practical version of LZ-index [Navarro 2004; 2008]:

- (1) *LZTrie*: is the trie formed by all the phrases B_0, \dots, B_n . Given the properties of LZ78 compression, this trie has exactly $n + 1$ nodes, each one corresponding to a phrase.
- (2) *RevTrie*: is the trie formed by all the reverse strings B_0^r, \dots, B_n^r . In this trie there could be internal nodes not representing any phrase. We call these nodes “empty”. Empty unary paths are compressed.



(a) Lempel-Ziv Trie (*LZTrie*) for the running example.



(b) *RevTrie* data structure.

$par: (() (() () ()) (()) (() ()) (() ()) (() ()) (() ()) (()))$
 $B: 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 1$
 $rlets: \$ a l r - b l d l a r p a - r a b - a$
 $rids: 0 17 17 12 8 3 9 10 2 14 16 13 11 4 15 5 6$

(c) Balanced parentheses representation of *RevTrie*, compressing empty unary paths. The bitmap B marks with a 0 the empty non-unary nodes. Notice that array $rids$ stores phrase identifiers only for non-empty nodes.

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
$Node[i]$	0	1	23	4	10	29	18	24	30	25	13	19	11	31	8	5	15	2

(d) *Node* data structure, assuming that the parentheses sequence starts from position zero.

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
$RNode[i]$	0	3	17	11	30	35	36	4	8	12	15	26	6	24	18	32	20	1

(e) *RNode* data structure, assuming that the parentheses sequence starts from position zero.

Fig. 4. LZ-index components for the running example.

(3) *Node*: is a mapping from phrase identifiers to their node in *LZTrie*.

(4) *RNode*: is a mapping from phrase identifiers to their node in *RevTrie*.

Fig. 4 shows the *LZTrie*, *RevTrie*, *Node*, and *RNode* data structures corresponding to our running example. We show preorder numbers, both in *LZTrie* and *RevTrie* (in the latter case only counting non-empty nodes), outside each trie node. In the case of *RevTrie*, empty nodes are shown in light gray.

3.2 Succinct Representation of the LZ-index Components

In the original work [Navarro 2004; 2008], each of the four structures described requires $n \log n + o(u \log \sigma)$ bits of space if they are represented succinctly.

- *LZTrie* is represented using the balanced parentheses representation [Munro and Raman 2001] requiring $2n + o(n)$ bits; plus the sequence *letts* of symbols labeling each trie edge, requiring $n \log \sigma$ bits; and the sequence *ids* of $n \log n$ bits storing the LZ78 phrase identifiers. Both *letts* and *ids* are stored in preorder, so we use $\text{preorder}(x)$ to index them. See Fig. 3(a) for an illustration.
- For *RevTrie*, balanced parentheses are also used to represent the *Patricia* tree [Morrison 1968] structure of the trie, compressing empty unary nodes and so ensuring $n' \leq 2n$ nodes. This requires at most $4n + o(n)$ bits. The *RevTrie*-preorder sequence of identifiers (*rids*) is stored in $n \log n$ bits (i.e., we only store the identifiers for non-empty nodes). Non-empty nodes are marked with bit vector $B[1..n']$, such that $B[j] = 0$ iff node x with preorder j is empty. Thus, the phrase identifier for node x is $\text{rids}[\text{rank}_1(B, j)]$. The symbols labeling the arcs of the trie and the Patricia-tree skips are not stored in this representation, since they can be retrieved by using the connection with *LZTrie*. Therefore, the navigation on *RevTrie* is more expensive than that on *LZTrie*.
- *Node* is just a sequence of n pointers to *LZTrie* nodes. As *LZTrie* is implemented using balanced parentheses, $\text{Node}[i]$ stores the position within the sequence for the opening parenthesis representing the node corresponding to phrase i . As there are $2n$ such positions, we need $n \log 2n = n \log n + n$ bits of storage. See Fig. 4(d) for an illustration.
- Finally, *RNode* is also represented as a sequence of n pointers to *RevTrie* nodes. As *RevTrie* has $n' \leq 2n$ nodes and since we use balanced parentheses for this trie, we need $n \log 4n = n \log n + 2n$ bits of space. See Fig. 4(e) for an illustration.

According to Lemma 2.7, the final size of the LZ-index is $4uH_k(T) + o(u \log \sigma)$ bits for $k = o(\log_\sigma u)$ (and hence $\log \sigma = o(\log u)$ if we want to allow for $k > 0$).

In theory, the succinct trie representations used [Navarro 2004] implement (among others) operations $\text{parent}(x)$ and $\text{child}(x, \alpha)$, both in $O(\log \sigma)$ time for *LZTrie*, and $O(\log \sigma)$ and $O(h \log \sigma)$ time respectively for *RevTrie*, where h is the depth of node x in *RevTrie* (the h in the cost comes from the fact that we must access *LZTrie* to get the label of a *RevTrie* edge). The operation $\text{ancestor}(x, y)$ is implemented in $O(1)$ time both in *LZTrie* and *RevTrie*.

In practice, however, the BP representation [Munro and Raman 2001] for general trees is used. Despite that under this representation operation $\text{child}(x, \alpha)$ is implemented by using operation $\text{child}(x, i)$ in $O(\sigma)$ worst-case time, this has shown to be very effective in practice [Navarro 2004; 2008]. Operation parent is supported in $O(1)$ time under this representation.

3.3 LZ-index Search Algorithm

Let us consider now the search algorithm for a pattern $P[1..m]$ [Navarro 2004; 2008]. For `locate` queries, pattern occurrences are originally reported in the format $\llbracket t, \text{offset} \rrbracket$, where t is the phrase where the occurrence starts, and offset is the distance between the beginning of the occurrence and the end of the phrase. Later,

in Section 5.2, we will show how to map these two values into a simple *text position*. As we deal with an implicit representation of the text (the *LZTrie*), and not the text itself, we distinguish three types of occurrences of P in T , depending on the phrase layout.

3.3.1 Occurrences of Type 1. The occurrence lies inside a single phrase (there are occ_1 occurrences of this type). Given the properties of LZ78, every phrase B_t containing P is formed by a shorter phrase B_ℓ concatenated to a symbol c (Property 2.3). If P does not occur at the end of B_t , then B_ℓ contains P as well. We want to find the shortest possible phrase B_i in the LZ78 referencing chain for B_t that contains the occurrence of P , that is, P is a suffix of B_i . But then P^r is a prefix of B_i^r , and it can be easily found by searching for P^r in *RevTrie* in $O(m^2 \log \sigma)$ time. Say we arrive at node v_r . Any node v'_r descending from v_r in *RevTrie* (including v_r itself) corresponds to a phrase terminated with P . For each such v'_r , we traverse and report the subtree of the corresponding *LZTrie* node v_{lz} (found using *rids* and *Node*). For any node v'_{lz} in the subtree of v_{lz} , we report an occurrence $\llbracket t, m + (\text{depth}(v'_{lz}) - \text{depth}(v_{lz})) \rrbracket$, where t is the phrase identifier (*ids*) of node v'_{lz} . Occurrences of type 1 are located in $O(m^2 \log \sigma + occ_1)$ time, since each occurrence takes constant time in *LZTrie*. For *count* queries we just need to compute the subtree size of each v_{lz} in *LZTrie*.

3.3.2 Occurrences of Type 2. The occurrence spans two consecutive phrases, B_t and B_{t+1} , such that a prefix $P[1..i]$ matches a suffix of B_t and the suffix $P[i+1..m]$ matches a prefix of B_{t+1} (there are occ_2 occurrences of this type). P can be split at any position, so we have to try them all. For every possible split, we search for the reverse pattern prefix $P^r[1..i]$ in *RevTrie* (getting node v_r) and for the pattern suffix $P[i+1..m]$ in *LZTrie* (getting node v_{lz}). The *RevTrie* node v_r for $P^r[1..i]$ is stored in array $C_r[i]$, since it shall be needed later. As in a trie all the strings represented in a subtree form a preorder interval, we have two preorder intervals: one in the space of reversed phrases (phrases finishing with $P[1..i]$) and one in that of the normal phrases (phrases starting with $P[i+1..m]$), and need to find the phrase pairs $(t, t+1)$ such that t is in the *RevTrie* interval and $t+1$ is in the *LZTrie* interval. Then, we check each phrase t in the subtree of v_r and report it if $\text{Node}[t+1]$ descends from v_{lz} . Each such check takes constant time. Yet, if the subtree of v_{lz} has fewer elements, one does the opposite: check phrases from v_{lz} in v_r , using $\text{RNode}[t-1]$. For every pair of consecutive phrases that passes this test, we report an occurrence $\llbracket t, i \rrbracket$.

The time to solve occurrences of type 2 is proportional to the smallest subtree size among v_r and v_{lz} , which can be arbitrarily larger than the number of occurrences reported. That is, we have no worst-case guarantees at search time¹. However, the average search time for occurrences of type 2 is $O(\frac{n}{\sigma^{m/2}})$ [Navarro 2004; 2008]. This is $O(1)$ for long patterns, $m \geq 2 \log_\sigma n$.

3.3.3 Occurrences of Type 3. The occurrence spans three or more phrases, $B_{t-1}, \dots, B_{\ell+1}$, such that $P[i..j] = B_t \dots B_\ell$, $P[1..i-1]$ matches a suffix of B_{t-1} and $P[j+1..m]$ matches a prefix of $B_{\ell+1}$ (there are occ_3 occurrences of this type).

¹The theoretical version [Navarro 2004] uses different structures which do offer such guarantees.

Since the LZ78 algorithm guarantees that every phrase represents a different string (Property 2.4), there is at most one phrase matching $P[i..j]$ for each choice of i and j . Therefore, if we partition P into more than two consecutive substrings, there is at most one pattern occurrence for such partition, which severely limits occ_3 to $O(m^2)$, the number of different partitions of P .

Let us define matrix $C_{lz}[1..m, 1..m]$ and arrays A_i , for $1 \leq i \leq m$, which store information about the search. We first identify the only possible phrase matching each substring $P[i..j]$ by performing a single trie search for each i and increasing j . We record in $C_{lz}[i, j]$ the *LZTrie* node corresponding to $P[i..j]$, and store the pair (id, j) at the end of A_i , such that id is the phrase identifier of the node corresponding to $P[i..j]$. Note that since we search for $P[i..j]$ for increasing j , we get the values of id in increasing order, as the phrase identifier of a node is always larger than that of the parent node. Therefore, the corresponding pairs in A_i are stored by increasing value of id . This process takes $O(m^2 \log \sigma)$ time.

Then we find the $O(m^2)$ maximal concatenations of successive phrases that match contiguous pattern substrings. For $1 \leq i \leq j \leq m$, for increasing j , we try to extend the match of $P[i..j]$ to the right. If id is the phrase identifier for node $C_{lz}[i, j]$, then we have to search for $(id + 1, r)$ in array A_{j+1} , for some r . Array A_{j+1} can be binary searched because it is sorted. If we find $(id + 1, r)$ in A_{j+1} , this means that $B_{id} = P[i..j]$ and $B_{id+1} = P[j + 1..r]$, i.e. the concatenation $B_{id}B_{id+1}$ equals $P[i..r]$. We repeat the process from $j = r$, and stop when the pair $(id + 1, r)$ is not found in the corresponding array (this means that the current concatenation is maximal). As we have to perform $O(m^2)$ binary searches in arrays of size $O(m)$, this procedure takes $O(m^2 \log m)$ worst-case time. In practice, the binary search is replaced by closed hashing schemes, taking $O(m^2)$ time on average [Navarro 2004, Section 6.5].

Once $P[i..j] = B_t \dots B_\ell$ is a maximal concatenation, we check whether phrase $B_{\ell+1}$ starts with $P[j + 1..m]$ by using operation $ancestor(C_{lz}[j + 1, m], Node[\ell + 1])$, in constant time per maximal concatenation. Finally we check whether phrase B_{t-1} ends with $P[1..i - 1]$ by checking whether $ancestor(C_r[i - 1], RNode[t - 1])$ holds in *RevTrie*, in constant time per maximal concatenation. If all these conditions hold, we report an occurrence $\llbracket t - 1, i - 1 \rrbracket$. Overall, occurrences of type 3 are located in $O(m^3 \log \sigma)$ time.

3.3.4 The Search Algorithm in Practice. In practice, the search algorithm proceeds as follows. We first search for every pattern substring $P[i..j]$ in *LZTrie*, and store the corresponding node in $C_{lz}[i, j]$. The search proceeds looking first for $P[1]$, then for $P[1..2]$, then for $P[1..3]$, and so on until we cannot find $P[1..j]$, for some j , or until $j = m$. Then we do the same starting from $P[2]$, and so on. We also store a matrix $C_{id}[1..m, 1..m]$, where $C_{id}[i, j]$ is the LZ78 phrase identifier for phrase $P[i..j]$; if $P[i..j]$ does not exist as an LZ78 phrase, then we store a null value.

The second step consists of searching for every prefix of the reversed pattern $P^r[1..j]$ in *RevTrie*. Recall that we need this for occurrences of type 1 and type 2. Recall also that searching in *RevTrie* is much slower than searching in *LZTrie*, so we try to reduce this work as much as possible. The results already obtained in C_{id} are useful. If we search for $P^r[1..j]$ in *RevTrie*, and $P[1..j]$ exists as a phrase in *LZTrie*, then $RNode[C_{id}[1, j]]$ is the *RevTrie* node we are looking for. Otherwise, $P^r[1..j]$

corresponds to an empty node in *RevTrie*, or to a position in a label between two nodes, and cannot be found with the *LZTrie*. Yet, we can reduce the cost as follows. Let i be the minimum value such that $C_{lz}[i, j]$ is defined, i.e. $P[i..j]$ exists as a phrase in *LZTrie* (and hence $P^r[i..j]$ exists as a reverse phrase in *RevTrie*). Then, we map to *RevTrie* with $RNode[C_{id}[i, j]]$, which corresponds to string $P^r[i..j]$, and from this node try to descend with the string $P^r[1..i - 1]$. This final search has to be done using operation *child* on *RevTrie*. The *RevTrie* node corresponding to $P^r[1..i]$ is stored in $C_r[i]$.

Then we use these results to search for occurrences of type 1, type 2, and type 3 respectively.

3.3.5 Displaying Occurrence Contexts. To display a context of length ℓ surrounding any occurrence reported, if an occurrence starts at phrase i , then we follow the upward path from $Node[i]$ up to the *LZTrie* root, outputting the symbols labeling the upward path. Then we perform the same procedure but now starting from $Node[i + 1]$ (alternatively $Node[i - 1]$) in *LZTrie*, and so on until we display the ℓ desired symbols, taking overall $O(\ell \log \sigma)$ time, because operation *parent* is supported in $O(\log \sigma)$ time in theory [Navarro 2004]. In practice, and using the BP representation for the tries, the time is $O(\ell)$. Finally, we can uncompress the whole text T in $O(u \log \sigma)$ time using the same idea, starting the procedure from the first LZ78 phrase (the time is $O(u)$ in practice).

4. THE LZ-INDEX AS A NAVIGATION SCHEME

4.1 The Original Navigation Scheme

The LZ-index structure can be regarded as a *navigation scheme* that permits us moving back and forth from trie nodes to the corresponding *preorder positions*, both in *LZTrie* and *RevTrie*. The phrase identifiers are common to both tries (arrays *ids* and *rids*) and permit moving from one trie to the other by using *Node* and *RNode* mappings.

Fig. 5 shows the navigation scheme, where solid arrows represent the main data structures of the index. Dashed arrows are asymptotically “for free” in terms of space requirement, since they are followed by applying *preorder* on the corresponding parentheses structure (see Section 2.4). From now on we use the subscript “*lz*” for the operations on *LZTrie*, and subscript “*r*” for *RevTrie*. The four solid arrows in the diagram are in fact the four main components in the space usage of the LZ-index: array of phrase identifiers in *LZTrie* (*ids*) and in *RevTrie* (*rids*), and mapping from phrase identifiers to tree nodes in *LZTrie* (*Node*) and in *RevTrie* (*RNode*). The structure is symmetric and we can move from any point to any other.

The structure, however, is redundant in the sense that the number of arrows is not minimal. Given a graph with t nodes (in our case t index components), t arrows are sufficient to connect them in both directions (actually forming a ring structure). Since nodes and preorder positions in the tries are “connected” using operations *preorder* and *selectnode* over the trie representations (see Section 2.4), we can think that there are only three main index components to connect: *LZTrie* (either nodes or preorder positions), phrase identifiers, and *RevTrie* (either nodes or preorder positions). Next we define more space-efficient representations for LZ-

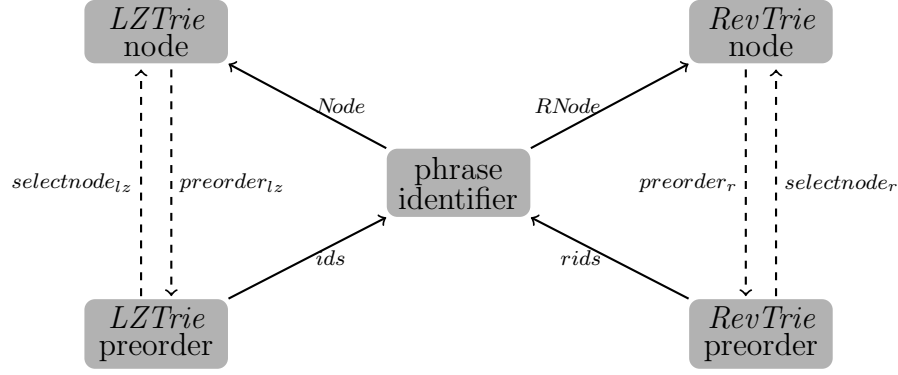


Fig. 5. The original LZ-index navigation structures over index components.

index, trying to reduce the number of arrows in the scheme. Note that, because of Lemma 2.7, we are interested in reducing the number of the index components that require $n \log n = uH_k(T) + o(u \log \sigma)$ bits of storage.

4.2 Schemes Requiring $3uH_k + o(u \log \sigma)$ bits

In this section we present schemes requiring only three solid arrows to connect the LZ-index components, thus forming a ring structure that still allows the same navigation as in the original LZ-index. Different choices yield different efficiencies depending on how often is each type of navigation used during the search.

4.2.1 *Scheme 1.* The following data structures conform this version of LZ-index:

- (1) *LZTrie*: The Lempel-Ziv trie, which is implemented with the following data structures
 - par*[0.. $2n - 1$]: The tree shape of *LZTrie* represented either with balanced parentheses [Munro and Raman 2001] or with DFUDS [Benoit et al. 2005], requiring in any case $2n + o(n)$ bits.
 - letts*[1.. n]: The array of symbols labeling the arcs of *LZTrie*, represented as explained in Section 2.3, depending on the representation used for *par*.
 - ids*[1.. n]: The array of LZ78 phrase identifiers in preorder. Since $ids[0] = 0$, we do not store this value. Note that *ids* is a permutation of $\{1, \dots, n\}$. The space requirement is $n \log n$ bits.
- (2) *RevTrie*: The *Patricia* tree [Morrison 1968] of the reversed LZ78 phrases, which is implemented with the following data structures
 - rpar*[0.. $2n' - 1$]: The *RevTrie* structure, represented either with BP or with DFUDS, compressing empty unary paths and thus ensuring $n' \leq 2n$ nodes, because empty non-unary nodes still exist. Thus, the space requirement is $2n' + o(n')$ bits.
 - rletts*[1.. n']: the array storing the first symbol of each edge label in *RevTrie*, represented as for *LZTrie* and requiring $n' \log \sigma + o(n')$ bits of space.
 - skips*[1.. n']: the *Patricia tree* skips of the nodes in preorder, using $\log \log u$ bits per node and inserting empty unary nodes when the skip exceeds $\log u$.

In this way, one out of u empty unary nodes could be explicitly represented. In the worst case there are $O(u)$ empty unary nodes, of which $O(u/\log u)$ are explicitly represented. This adds $O(u/\log u)$ nodes to n' , which translates into $O((n' + \frac{u}{\log u})(3 + \log \sigma + \log \log u)) = o(u \log \sigma)$ bits overall for the *RevTrie* nodes, symbols, and skips.

- $B[1..n']$: A bit vector supporting *rank* and *select* queries, and requiring $n'(1 + o(1))$ bits of space [Munro 1996]. This bit vector marks the non-empty nodes: The j -th bit of B is 1 iff the node with preorder position j in *rpar* is not empty, otherwise the bit is 0. Given a position i in *rpar* corresponding to a *RevTrie* node, the corresponding bit in B is $B[\text{preorder}_r(i)]$. The preorder of a node p counting only non-empty nodes can be computed as $\text{rank}_1(B, \text{preorder}_r(i))$.
- (3) $RNode[1..n]$: The mapping from phrase identifiers to the corresponding *RevTrie* node. Since we represent nodes as the positions of opening parentheses, and since there are $2n' \leq 4n$ such positions in *RevTrie*, this mapping needs $n \log 4n = n \log n + 2n$ bits. We only store pointers to non-empty nodes.
- (4) $Rev[1..n]$: A mapping from a *RevTrie* preorder position to the corresponding *LZTrie* node, defined as $Rev[i] = Node[\text{rids}[i]]$. Given a position i in *rpar* corresponding to a non-empty *RevTrie* node, the corresponding *Rev* value (i.e., *LZTrie* node) is $Rev[\text{rank}_1(B, \text{preorder}_r(i))]$. The space requirement is $n \log n + n$ bits.

The resulting navigation scheme is shown in Fig. 6(a). The search algorithm remains the same since we can map preorder positions to nodes in the tries and vice versa (see Section 2.3), and also we can simulate the missing arrays $\text{rids}(i) \equiv \text{ids}[\text{preorder}_{l_z}(Rev[i])]$ and $Node(i) \equiv Rev[\text{rank}_1(B, \text{preorder}_r(RNode[i]))]$, all of which take constant time.

We have reduced the space requirement to $3n \log n + 3n \log \sigma + 2n \log \log u + 11n + o(u) = 3n \log n + o(u \log \sigma)$ bits if $\log \sigma = o(\log u)$, which according to Lemma 2.7 is $3uH_k(T) + o(u \log \sigma)$ bits, for any $k = o(\log_\sigma u)$.

The *child* operation on *RevTrie* can now be computed in $O(1)$ time if we use DFUDS, and because we store *rletts* and the skips. Compare to the $O(h \log \sigma)$ time of the original LZ-index [Navarro 2004]. Now, because *RevTrie* is a Patricia tree and the underlying strings are not readily available, it is not obvious how to traverse it. The next lemma addresses this issue.

LEMMA 4.1. *Given a string $s \in \Sigma^*$, we can determine whether it is represented in *RevTrie* or not (finding the corresponding node in the affirmative case) in $O(|s|)$ time.*

PROOF. To find the node corresponding to string s we descend from the *RevTrie* root, using operation $\text{child}(x, \alpha)$ on the first symbol of each edge label, which is stored in *rletts*, and using the skips to compute the next symbol of s to use in the descent. If s cannot be consumed while descending, then we determine that it is not represented in *RevTrie* in $O(|s|)$ time. Otherwise, assume that after consuming string s in this way we arrive at node v_r with preorder j in *RevTrie* (counting only non-empty nodes). The string labeling the root-to- v_r path in *RevTrie* can be computed by accessing the node $v_{l_z} = Rev[j]$ in *LZTrie*, and then extracting the

string labeling the v_{l_z} -to-root path in $LZTrie$. Then we compare that string against s to verify that the node we arrived at corresponds to s , or otherwise that s does not occur in $RevTrie$.

In case node v_r in $RevTrie$ is empty, $Rev[j]$ is undefined. Notice, however, that there must be at least one non-empty node descending from this empty node, since leaves in $RevTrie$ cannot be empty as they always correspond to an LZ78 phrase. Given that the string represented by every non-empty node in the subtree of node v_r has the string s as a prefix, the corresponding strings in $LZTrie$ have s^r as a suffix. So we can use any Rev value within the subtree of node v_r in order to map to the $LZTrie$ and then extract the string it represents. We can use, for example, the value $Rev[rank_1(B, j) + 1]$, which corresponds to the next non-empty node within the subtree of node v_r . We know when to stop extracting, since we know the length of the string we are looking for.

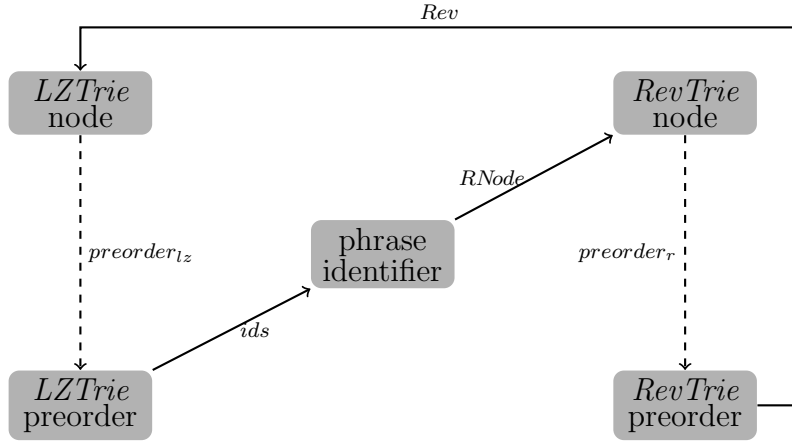
The overall cost for the descending process is therefore $O(|s|)$. \square

Operations *child* and *parent* on $LZTrie$ can be also computed in $O(1)$ time if we use DFUDS on this trie, versus the $O(\log \sigma)$ time (in theory) of the original LZ-index. Hence, occurrences of type 1 are located in $O(m + occ_1)$ time; occurrences of type 2 are located now in $O(\frac{n}{\sigma^{m/2}})$ average time; occurrences of type 3 are located as in the original LZ-index, in $O(m^2)$ average time (by using hashing to find the maximal concatenations of phrases). Therefore, the *occ* occurrences of P can be located in $O(m^2 + \frac{n}{\sigma^{m/2}})$ average time. This time is $O(m^2)$ on average for $m \geq 2 \log_\sigma u$. It has been shown how to achieve worst-case guarantees with reduced schemes of the LZ-index [Arroyuelo et al. 2006]. However, we are focusing in practical performance here.

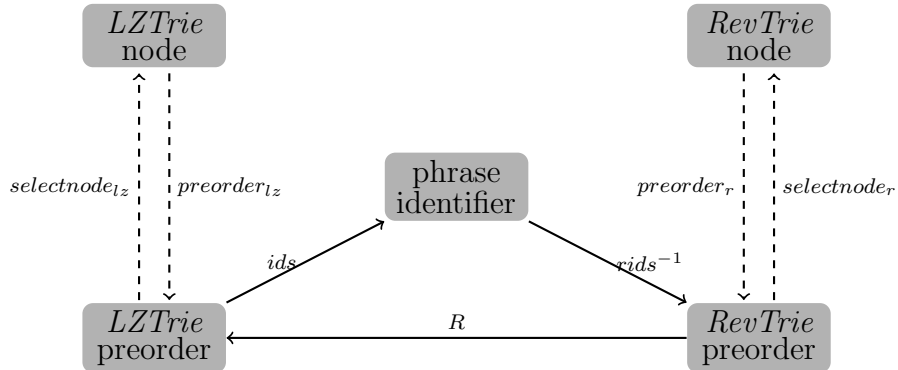
Practical Issues. On the practical side, the access from $RevTrie$ nodes to the corresponding $LZTrie$ node is faster under this scheme, since the direct link *Rev* is faster than the composition of *rids* and *Node* of the original scheme. This is good, for example, for finding occurrences of type 1, which can be dominant for short patterns, as there is a high probability that an occurrence is contained in a single phrase. However, sometimes we must follow longer navigation paths in the search process: for example, when finding occurrences of type 2, we can choose to traverse the subtree in $RevTrie$, and for each phrase identifier id in such subtree apply $Node(id + 1)$ to check whether it descends from the appropriate subtree in $LZTrie$. As now we have to simulate *Node*, this is more expensive (in practice, even if not asymptotically) than in the original scheme. Even worse, since array *rids* is not stored in $RevTrie$, we must simulate $rids(i)$ to get the phrase identifier id . Therefore, the search time could be increased in practice, depending on the number of each type of occurrence.

Let us study occurrences of type 2 in more detail, since they seem to be critical under this scheme. Suppose that for a given partition $P[1..i]$ and $P[i + 1..m]$ of P we get nodes v_{l_z} and v_r in $LZTrie$ and $RevTrie$ respectively. If we choose to traverse the subtree of v_r in $RevTrie$, then for each node v'_r in this subtree we get the corresponding phrase identifier $id = ids[preorder_{l_z}(Rev[p_r])]$, where p_r is set initially as $p_r = rank_1(B, preorder_r(v_r))$, and it is incremented by one with each node in a preorder traversal of the subtree. We then check whether the node $Rev[rank_1(B, preorder_r(RNode[id + 1]))]$ descends from v_{l_z} in $LZTrie$. If, on the

other hand, we choose to traverse the subtree of v_{lz} in $LZTrie$, then for each node v'_{lz} in this subtree we get the phrase identifier as $id = ids[p_{lz}]$, where p_{lz} is set initially as $p_{lz} = preorder_{lz}(v_{lz})$, and it is incremented by one with each node in a preorder traversal. We then check whether the node $RNode[id - 1]$ descends from v_r in $RevTrie$. Empirically, a check from $RevTrie$ to $LZTrie$ is about 3 times as expensive as in the opposite direction, and thus we choose to traverse the subtree of v_r whenever its size is less than $1/3$ the size of the subtree of v_{lz} .



(a) Scheme 1.



(b) Scheme 2.

Fig. 6. Reduced navigation schemes over LZ-index components, requiring $3uH_k + o(u \log \sigma)$ bits.

4.2.2 *Scheme 2.* This scheme tries to reduce the space requirement while also reducing the average path length in the navigation scheme:

- (1) *LZTrie* The Lempel-Ziv trie, defined just as for Scheme 1.

- (2) *RevTrie*: The *Patricia* tree of the reversed LZ78 phrases, defined just as for Scheme 1, but now we add
 — $rids^{-1}[1..n]$: The explicit representation of the inverse of permutation $rids$ of the original LZ-index definition, requiring $n \log n$ bits.
- (3) $R[1..n]$: A mapping from *RevTrie* preorder positions to *LZTrie* preorder positions defined as $R[i] = ids^{-1}(rids[i])$ and requiring $n \log n$ bits. Given a position i in *rpar* corresponding to a non-empty *RevTrie* node, the corresponding R value (i.e., preorder in *LZTrie*) can be computed as $R[rank_1(B, preorder_r(i))]$.

The resulting navigation scheme is shown in Fig. 6(b). We can compute $rids(i) \equiv ids[R[i]]$, $RNode(i) \equiv selectnode_r(rids^{-1}[i])$, and $Node(i) \equiv selectnode_{lz}(R[rids^{-1}[i]])$, all in constant time.

The space requirement is $3n \log n + 3n \log \sigma + 2n \log \log u + 8n + o(u) = 3n \log n + o(u \log \sigma)$ bits if $\log \sigma = o(\log u)$, which according to Lemma 2.7 is $3uH_k(T) + o(u \log \sigma)$ bits, for any $k = o(\log_\sigma u)$.

If we use DFUDS to represent both *LZTrie* and *RevTrie*, then we can locate the occ occurrences of pattern P in $O(m^2 + \frac{n}{\sigma^{m/2}})$ average time, which is $O(m^2)$ for $m \geq 2 \log_\sigma u$.

Practical Issues. It is interesting to note that the average path length of this scheme is shorter than that of Scheme 1, which can translate into a more efficient navigation among index components. In this scheme, for occurrences of type 1 we have direct access to a *LZTrie* preorder by using R , and then we have to apply *selectnode* to get the node whose subtree contains the pattern occurrences. This can be slightly slower than for Scheme 1, where we have direct access to the corresponding node.

However, for occurrences of type 2 we have to follow shorter paths than for Scheme 1. Suppose that for a given partition $P[1..i]$ and $P[i+1..m]$ of P we get nodes v_{lz} and v_r in *LZTrie* and *RevTrie* respectively. If we choose to traverse the subtree of v_r in *RevTrie*, for each node in this subtree we get the corresponding phrase identifier id by using both the R mapping and then ids . Then we use $R[rids^{-1}[id+1]]$ to get the corresponding *LZTrie* preorder, and then we check whether this preorder lies within the subtree of v_{lz} . If, on the other hand, we choose to traverse the subtree of v_{lz} in *LZTrie*, for every node in this subtree we get the corresponding phrase identifier id using ids , and then we check whether the preorder $rids^{-1}[id-1]$ lies within the subtree of v_r in *RevTrie*. Thus, a check from *RevTrie* to *LZTrie* is twice as expensive as in the opposite direction, and thus we choose to traverse the subtree of v_r whenever its size is less than half the size of the subtree of v_{lz} .

Fortunately, the checks of type 2 can be carried out directly on the preorders of both tries, avoiding the use of (the usually expensive) *selectnode* to get the corresponding trie node: if we choose to traverse the subtree of v_r , for example, we compute the preorder interval for the subtree of v_{lz} as $[preorder_{lz}(v_{lz})..preorder_{lz}(v_{lz}) + subtreesize_{lz}(v_{lz}) - 1]$ (recall that *preorder* is computed by means of *rank*), and then we check whether the *LZTrie* preorders we get from the nodes in the subtree of v_r lie within the preorder interval of v_{lz} . In this way, we compute just one *rank* per partition to get the interval, and then we check the *LZTrie* preorder of the

candidates by using just this interval, rather than computing *selectnode* for every possible candidate in that partition. This introduces very important savings in the practical search time.

There are many other possible schemes that achieve $3uH_k(T) + o(u \log \sigma)$ bits of space. We have focused on the two most promising ones. For example, consider a scheme where we only replace the R mapping of Scheme 2 by the *Rev* mapping of Scheme 1. We have again direct access for occurrences of type 1, but occurrences of type 2 now introduce the computation of *rank* in *LZTrie* for every possible candidate, which is expensive.

4.3 Schemes Requiring $(2 + \epsilon)uH_k + o(u \log \sigma)$ bits

In Section 4.2 we have used the minimal number of arrows to connect the three main components of LZ-index, forming a ring structure. It seems that we cannot reduce even more the space requirement of the index by using our navigation-scheme approach. However, many of the data structures of the LZ-index are just permutations, and so the corresponding arrows can be made bidirectional by means of the data structure for permutations [Munro et al. 2003] described in Section 2.3, using just $(1 + \epsilon)n \log n + n + o(n)$ bits for both arrows. This opens several new possibilities.

4.3.1 *Scheme 3.* This scheme represents the following data:

- (1) *LZTrie*: The Lempel-Ziv trie, defined as for Scheme 1, except that now we use the representation of Munro et al. [2003] for *ids* such that the inverse permutation ids^{-1} can be computed in $O(1/\epsilon)$ time, requiring $(1 + \epsilon)n \log n + n + o(n)$ bits for any $0 < \epsilon < 1$.
- (2) *RevTrie*: The *Patricia* tree of the reversed LZ78 phrases, defined as in Scheme 1, but now we add the array *rids* represented using [Munro et al. 2003], so as to be able to compute $rids^{-1}$ efficiently.

The resulting navigation scheme is shown in Fig. 7(a). We can simulate the missing arrays $Node(i) \equiv selectnode_{lz}(ids^{-1}(i))$ and $RNode(i) \equiv selectnode_r(rids^{-1}(i))$, all in $O(1/\epsilon)$ time.

The space requirement is $(2 + \epsilon)n \log n + 3n \log \sigma + 2n \log \log u + 10n + o(u) = (2 + \epsilon)n \log n + o(u \log \sigma)$ bits, which according to Lemma 2.7 is $(2 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits, for any $k = o(\log_\sigma u)$.

Once again, if we use DFUDS to represent both tries of the LZ-index, then occurrences of type 1 can be located in $O(m + \frac{occ_1}{\epsilon})$ time, because we must use ids^{-1} to access to *LZTrie*; occurrences of type 2 are located in $O(\frac{n}{\epsilon \sigma^{m/2}})$ because we must use inverse permutations to move between tries; occurrences of type 3 are located in $O(\frac{m^2}{\epsilon})$ time, since we need to use *Node* and *RNode* to check every possible candidate, in $O(\frac{m^2}{\epsilon})$ time overall. Thus, the *occ* occurrences of P can be located in $O(\frac{m^2}{\epsilon} + \frac{n}{\epsilon \sigma^{m/2}})$ average time, for $0 < \epsilon < 1$. This is $O(\frac{m^2}{\epsilon})$ for $m \geq 2 \log_\sigma u$.

Practical Issues. This scheme stores the phrase identifiers for both tries which, as we have seen for the previous schemes, is very convenient for occurrences of type 2: recall that when traversing the *RevTrie* subtree we have to get the phrase identifier

of each node in the subtree; if we do not store the *RevTrie* identifiers, we have to access *LZTrie* to get them (as is the case of Schemes 1 and 2) and then we have to access *LZTrie* again to perform the check. This is not the case for Scheme 3. However, now we have paths including inverse permutations, which introduce an extra time overhead in practice.

Notice also that this scheme is symmetric in the sense that the checks for occurrences of type 2 cost the same in any direction we choose.

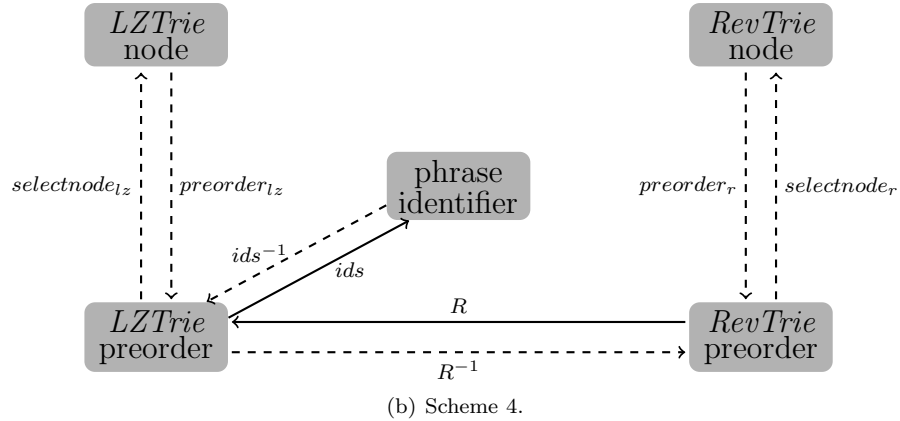
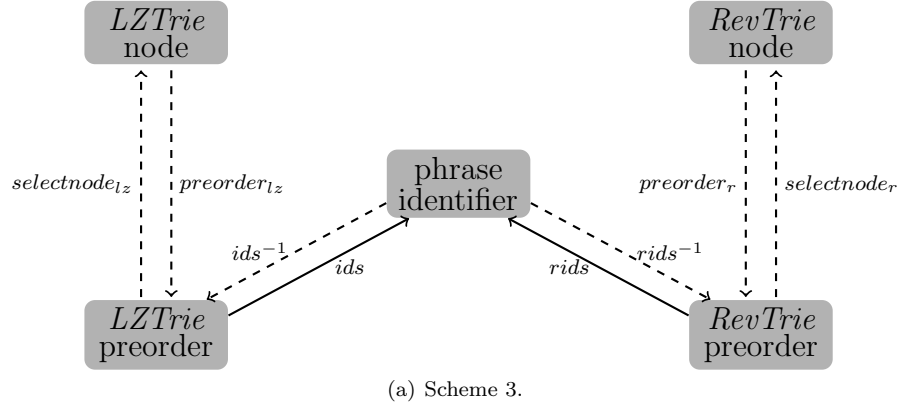


Fig. 7. Reduced navigation schemes over LZ-index components, requiring $(2 + \epsilon)uH_k + o(u \log \sigma)$ bits.

4.3.2 *Scheme 4*. This scheme represents the following data:

- (1) *LZTrie*: The Lempel-Ziv trie, defined just as in Scheme 3.
- (2) *RevTrie*: The *Patricia* tree of the reversed LZ78 phrases, defined just as in Scheme 1.

- (3) $R[1..n]$: The mapping from *RevTrie* preorder positions to *LZTrie* preorder positions, as defined in Scheme 2. This time R is implemented using the succinct data structure for permutations of Munro et al. [2003], requiring $(1 + \epsilon)n \log n + n + o(n)$ bits to represent R and compute R^{-1} in $O(1/\epsilon)$ worst-case time.

In Fig. 7(b) we draw the navigation scheme. We can simulate the missing arrays $rids(i) \equiv ids[R[i]]$, $RNode(i) \equiv selectnode_r(R^{-1}(ids^{-1}(i)))$, and $Node(i) \equiv selectnode_{lz}(ids^{-1}(i))$, all of which take $O(1/\epsilon)$ time.

The space requirement is $(2 + \epsilon)n \log n + 3n \log \sigma + 2n \log \log u + 10n + o(u) = (2 + \epsilon)n \log n + o(u \log \sigma)$ bits, which according to Lemma 2.7 is $(2 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits, for any $k = o(\log_\sigma u)$. Just as for the previous scheme, the *occ* occurrences of P can be located in $O(\frac{m^2}{\epsilon} + \frac{n}{\epsilon \sigma^{m/2}})$ average time, for $0 < \epsilon < 1$.

Practical Issues. This scheme has more efficient access between tries than Scheme 3, as we have to use R in the *RevTrie*-to-*LZTrie* direction, and R^{-1} in the opposite way. However, since we store the phrase identifiers only in *LZTrie*, retrieving the identifier of a *RevTrie* node requires to access two arrays. For occurrences of type 2, the checks from *RevTrie* to *LZTrie* require to access R , then ids , and finally ids^{-1} , while in the opposite way we need to use ids , then ids^{-1} , and finally R^{-1} . The latter case can be more expensive since we have to compute two inverse permutations. Note also that ids^{-1} is used in both directions for occurrences of type 2, which means that this inverse permutation is the most used at search time. Hence, given an amount of space we are able to use, we should use a denser sampling for ids^{-1} than for R^{-1} .

Again, we have focused on the most promising schemes requiring $(2 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits, although there are many other choices.

5. SOME IMPLEMENTATION DETAILS

We describe in this section the most important details in the implementation of our indices. We followed the API interface specification provided in the *Pizza&Chili* Corpus [Ferragina and Navarro 2005], and made the source code available at <http://pizzachili.dcc.uchile.cl/indexes/LZ-index/>.

5.1 Representing the Tries

We defined our indices in Section 4 in such a way that we can use almost any suitable representation for the tries that compose the indices. We just need to define accordingly the *preorder* and *selectnode* operations for the chosen representation. Having this into account, we implement our reduced LZ-indices in two different ways:

- (1) Using the BP representation [Munro and Raman 2001] for both *LZTrie* and *RevTrie*;
- (2) using the DFUDS representation [Benoit et al. 2005] for *LZTrie*, and the BP representation for *RevTrie*.

Note we do not use DFUDS for *RevTrie*, as it requires more space. Moreover, just as for the original LZ-index, we do not store the symbols labeling the *RevTrie* edges (i.e., the first symbol of each string labeling an edge) nor the Patricia-tree skips.

This is in order to save space in practice, since these can be computed by using the connection with the *LZTrie*: previous experiments [Navarro 2004; 2008] showed that this is sufficient for *RevTrie* as most navigations on it are supported by using the *LZTrie* (and those that are not are usually deep in the trie, where the arity is very low and the attractive of DFUDS vanishes). We needed to use these arrays in theory in order to guarantee the average-time complexity of our data structures.

We describe these implementations in which follows.

5.1.1 Using BP Representation. We implemented the trie operations on top of the data structure for balanced parentheses of Navarro [2004, Section 6.1]. The trie operations are supported as explained in Section 2.3. Recall that in order to support the operations on BP we must support operations *findclose*, *excess*, and *enclose* on the sequence of balanced parentheses. Therefore, we do not store information to support *findopen*, thus saving space. Moreover, we do not need to support operation *parent* on *RevTrie*, and thus we do not store information to support *enclose* on the parentheses representing this trie.

Operation *child*(x, α) is supported by using *child*(x, i), for $i = 1, 2, \dots$, until finding the child labeled α . This is because the symbols labeling the children of x are scattered throughout array *letts* and must be found one by one using the operations on the parentheses. Whenever we need to support *rank* and *select* queries (this is, on top of the parenthesis sequences, to represent bitmap B in *RevTrie*, and for the permutation data structures), we use the data structure of González et al. [2005].

5.1.2 Using DFUDS Representation. The main idea of using the DFUDS representation for *LZTrie* is to reduce the time overhead for computing operation *child*(x, α) incurred by BP.

As done for BP, we represent the DFUDS sequence of *LZTrie* on top of the data structure for balanced parentheses of Navarro [2004]. Note that the DFUDS representation of a trie tends to have far matching parentheses, since every node is formed by a number of opening parentheses (indicating the degree of the node), and only a few of these parentheses have the corresponding closing parentheses close enough so as not to be stored in the hash tables. Thus the hash tables tend to require more space under DFUDS. We also use the data structure of González et al. [2005] to support *rank* and *select* queries.

We study the way in which the trie operations are used by the LZ-index search algorithm in order to make this representation more efficient. For example, DFUDS introduces a heavier use of *rank* and *select* in its operations (see Section 2.4), according to the original definitions [Benoit et al. 2005]. However, many of these are redundant in the LZ-index (sometimes they are repeated twice in a given sequence of operations, as we will see below), and many others can be replaced by sequential scans over *par*, since the DFUDS position we are looking for should be not so far away from the current position. A list with the most important implementation details follows:

Supporting Operation findopen. Unlike the representation of LZ-index based exclusively in BP, with DFUDS we need to provide operation *findopen* over the parentheses sequence. This is necessary, for example, to compute operation *parent*

symbols, and then we map again to the DFUDS sequence to find the corresponding child (this involves only one *findclose* operation).

Avoiding Operation depth on DFUDS. The original DFUDS representation [Benoit et al. 2005] does not provide operation *depth*, which is latter supported by Jansson et al. [Jansson et al. 2007] in $O(1)$ time and requiring $o(n)$ extra bits of space. However, in the LZ-index we need to use *depth* in a very limited way, which helps us implement this operation simply and efficiently. We choose not to store any extra depth information, and thus we save space. The depth of a node could be computed by brute force, by successively going to the parent in *LZTrie* up to reaching the trie root. Therefore, we seek to avoid the computation of *depth* in the LZ-index search algorithm.

In the LZ-index search algorithm, we only need to use operation *depth* when locating occurrences of type 1 (see Section 3.3.1). Recall that in order to find occurrences of type 1 we must first search for P^r in *RevTrie*, getting node v_r . Then, for each node in the subtree of v_r we map to the corresponding node v_{lz} in *LZTrie*, and then we traverse the subtree of v_{lz} to report occurrences of type 1. The problem here is how to compute the offset of every occurrence within the corresponding phrase (in Section 3.3.1 we use operation *depth* to do that). However, note that the offset for node v_{lz} is m , since the phrase ends with P . Note also that the offset for the children of v_{lz} is $m + 1$, and in general the offset of node v'_{lz} within the subtree of v_{lz} is $m + d$, where d is the difference of depths between v_{lz} and v'_{lz} .

So, instead of computing the depth of the nodes within the subtree of v_{lz} , which is expensive in our representation, we compute the offset of the nodes. This problem can be solved in a straightforward way in BP, since we perform a preorder traversal from v_{lz} , with initial offset m . We then increment the offset every time we enter a new subtree (which is indicated by ‘(’ in BP), and decrement it when leaving a subtree (which is marked by ‘)’ in BP). The preorder traversal is carried out by sequentially traversing the BP sequence and array *ids*, and not by using operation *child* on *LZTrie*. However, this is not so easy in DFUDS, since, for example, there is no clear marker of the end of a subtree (i.e., in a sequential scan on DFUDS there is no direct way to know whether a node is the last one in a given subtree).

Therefore, in the DFUDS representation we start a preorder traversal from node v_{lz} , and store the degree (number of children) of v_{lz} and its offset m in an initially empty stack. We then continue with the next node in preorder. At every step, the offset for the current node is computed as the offset of the parent node (which is the one at the top of the stack) plus one. Every time the degree stored in the top is (or becomes) 0 (leaves are a particular case), we pop it, and then decrement the degree of the node in the top, indicating that a new child of this node has been fully processed (this can eventually produce more pops when all of the children of the top node have been processed). As it can be seen, we are using the stack to know when the subtree of a node has been completely processed. This procedure ends when the stack becomes empty.

As in the case of BP, the preorder traversal is performed by a sequential scan on the DFUDS sequence, and the degree of every node x is computed by counting the number of opening parentheses in the representation of x .

This procedure could be also used to compute the depth of all nodes within the

subtree of v_{lz} , by initializing the stack with the depth of v_{lz} , instead of storing the initial offset m . Note also we do not need to explicitly store depths or offsets, as they correspond to the stack height plus a constant.

Implementation of Operation degree. In our implementation we do not use the original definition for operation *degree*, which is based on operation *select* in order to find the next closing parenthesis which finishes the definition of the node. This allows us to count the number of opening parentheses defining the current node. In practice, in most cases this closing parenthesis is not so far away from the current position in DFUDS, except perhaps for the trie root. This is because in practice the tries tend to have high degrees only in the first levels. Thus, we explicitly store the degree of the root node, and for the rest of the nodes we perform a sequential scan on the DFUDS sequence. To avoid looking at every parenthesis in the process, we advance by machine words (in our experiments this shall mean 32 bits), until finding the first word containing a closing parenthesis. We represent opening parentheses with a 0, and closing ones with a 1. Thus, we advance while the corresponding word represents a 0, i.e. it stores only opening parentheses. However, in the case of very large alphabets the original definition of *degree* could be better, or we could replace the sequential search by an exponential search using the rank subdirectory.

5.2 Computing Text Positions

We add to our indices a data structure in order to transform pattern occurrences in the format $\llbracket t, offset \rrbracket$ into real text positions. We define array $TPos$ storing the absolute starting position (in the text) for the LZ78 phrases with identifier $i \cdot b$, for $i \geq 0$, for a total of $\frac{n}{b} \log u$ bits. We also define array $Offset$, storing in $Offset[i \cdot b + j]$, for $j \geq 0$, the offset (in number of text symbols) of phrase $i \cdot b + j$ with respect to phrase $i \cdot b$, requiring $n \log M$ extra bits of space, where M is the maximum number of text positions between two consecutive sampled phrases. If $LZTrie$ has height $h = O(\log n)$ (which is true in practice with high probability [Knessl and Szpankowski 2000]), then $M = O(b \log n)$, and thus array $Offset$ requires $O(n(\log b + \log \log n))$ bits. By choosing $b = \log u$ the total space requirement for both arrays is $n + O(n \log \log u) = o(u \log \sigma)$ bits.

Given an occurrence $\llbracket t, offset \rrbracket$, the real text position for that occurrence can be computed in constant time as $TPos[\lfloor (t+1)/b \rfloor] + Offset[t+1] - offset$. In our experiments, we choose $b = 32$ (in a 32-bit machine) such that the division by b and taking the floor can be carried out efficiently by using shifts on machine words. For **extract** queries, we are given a text position from where to extract and we want to know the corresponding $LZTrie$ node from where to start uncompressing the text. Therefore, given a text position pos we can obtain the phrase containing pos by first binary searching $TPos$, finding the greatest phrase $i \cdot b$ such that its position $TPos[i]$ is smaller or equal to pos . Then, we sequentially look in the corresponding segment of $Offset[i \cdot b..(i+1) \cdot b]$ for the greatest phrase t whose starting position does not exceed pos . Thus, the text position pos belongs to phrase t . The time for this operation is $O(\log u)$, because of the binary search on $TPos$ and the sequential search in the segment of $Offset$.

In Section 6.2 we will show the experimental space requirement of this data structure for a set of real-life texts we have used.

5.3 Supporting Partial locate Queries

In many applications it is quite common that we do not need to find all of the pattern occurrences, but just a few (arbitrary) of them. For this kind of applications, we design an algorithm to answer *partial locate* queries, where we are interested in locating just K arbitrary pattern occurrences. Our algorithm, which profits from the properties of the LZ-index in order to support fast searches, is as follows:

- (1) Given a search pattern P , notice that a particular occurrence of it can be found by searching for P in $LZTrie$. In other words, P equals an LZ78 phrase, which is a particular case of occurrence of type 1, and can be found very fast in practice since this is better than using the slower *RevTrie* [Navarro 2008]. If P exists as a phrase, say corresponding to node v_{lz} in $LZTrie$, then all of the nodes descending from v_{lz} in the trie also correspond to occurrences of P , and can be used to answer the query. Thus, we traverse the subtree of v_{lz} and report every node found, as done for occurrences of type 1 (see Section 3.3.1). We stop the procedure as soon as we find K pattern occurrences.
- (2) If the previous step was not enough to answer the query, and in case that P exists in $LZTrie$, then we map to the node corresponding to P^r in *RevTrie* (which exist for sure since P is an LZ78 phrase), and go on to locate the rest of occurrences of type 1 as usual. Otherwise, we delay occurrences of type 1 for a further step and go to the next step.
- (3) In case P does not exist as an LZ78 phrase, we proceed with occurrences of type 2, trying to reuse as much as possible the work already done in step (1). We are thus delaying occurrences of type 1 for a further step. Let $P[1..i]$ be the longest proper prefix of the pattern that exists as an LZ78 phrase. Hence, $P^r[1..i]$ exists as a reverse phrase in *RevTrie*. Because of Property 2.3, every prefix of $P[1..i]$ also exists as a phrase in LZ78 (and the corresponding reverses exist in *RevTrie*). Then, to reuse the work already done when searching for P in $LZTrie$, we map to the *RevTrie* node corresponding to $P^r[1..i]$, which gives us node v_r , and then search for $P[i + 1..m]$ in $LZTrie$, to get node v_{lz} . We then search for occurrences of type 2 corresponding to the partition $P[1..i]$ and $P[i + 1..m]$, using the nodes v_{lz} and v_r , in the usual way and stopping as soon as we find K occurrences. Note that by choosing the longest prefix $P[1..i]$ that exists in $LZTrie$, we are reducing as much as possible the length of the suffix $P[i + 1..m]$ to be searched in $LZTrie$. If this was not enough to answer the query, we repeat the process for the occurrences of $P[1..i - 1]$ and $P[i..m]$, in a similar way, and so on.
- (4) We search for the remaining occurrences of type 2 (i.e., using those partitions of the pattern that were not tried in the previous step)
- (5) Then, we continue with occurrences of type 1, as usual and just if P does not exist in $LZTrie$ (i.e., these were not tried before).
- (6) Finally we try occurrences of type 3.

We call *level 0* of the search to the step of searching for P in $LZTrie$, *level 1* the search for occurrences of type 1 (either at steps (2) and (5)), *level 2* the search for occurrences of type 2 (either at steps (3) and (4)), and *level 3* the search for occurrences of type 3 (step (6)). As it can be seen, we try to get fast access to

the pattern occurrences, avoiding as much as we can the trie navigations, which become more expensive if we want to locate just a few occurrences. We shall use this approach also to support efficient `exists` queries (similar to $K = 1$).

6. EXPERIMENTAL RESULTS

We have now a number of practical reduced schemes for LZ-index, each one requiring a different amount of memory space. Hence given an amount of available storage, it is interesting to know which alternative is the best for that space. We hope that larger alternatives are faster in practice, whereas the smaller ones will still be competitive against the best existing indices.

6.1 Experimental Setup

For the experiments of this section we used an Intel(R) Pentium(R) 4 processor at 3 GHz, about 4 gigabytes of main memory and 1024 kilobytes of cache, running version 2.6.13-gentoo of Linux kernel. We compiled our algorithms with `gcc 3.3.6` using full optimization.

6.1.1 *Text Collections.* We test our indices in different practical scenarios, using the texts provided in the *Pizza&Chili* Corpus [Ferragina and Navarro 2005]:

- English Texts: although in many cases natural-language texts are searched for whole words or phrases, there are many other cases where a more powerful full-text search is needed. For the experiments with English text we use the file `http://pizzachili.dcc.uchile.cl/texts/nlang/english.200MB.gz` of 200 megabytes.
- DNA Sequences: nowadays, one of the main applications of full-text indexing is that of computational biology, in particular indexing DNA sequences. We test with the file `http://pizzachili.dcc.uchile.cl/texts/dna/dna.200MB.gz`, of 200 megabytes.
- MIDI Pitch Sequences: a very interesting application that has appeared in recent years is that of processing MIDI pitch sequences. In this case we test with the file of about 53 megabytes downloadable from `http://pizzachili.dcc.uchile.cl/texts/music/pitches.gz`.
- XML Texts: since XML is becoming the standard to represent semi-structured text databases as well as in many other applications, there exists the need of managing a huge amount of texts of this kind. We test with the XML file of 200 megabytes downloadable from `http://pizzachili.dcc.uchile.cl/texts/xml/dblp.xml.200MB.gz`.
- Proteins: another interesting application of text-indexing tools in biology is that of indexing and searching proteins. We use the file `http://pizzachili.dcc.uchile.cl/texts/protein/proteins.200MB.gz`, of 200 megabytes.
- Source Code: to test our indices in applications like software development, we use the source-code file `http://pizzachili.dcc.uchile.cl/texts/code/sources.200MB.gz`, of 200 megabytes.

6.1.2 *Comparison against Other Indices.* We compare our indices against the most efficient indices we are aware of, most of them available in *Pizza&Chili*:

Sadakane's Compressed Suffix Array (CSA). This index [Sadakane 2003] is a representative of the family of *compressed suffix arrays* [Grossi and Vitter 2005; Grossi et al. 2003; Sadakane 2003]. It requires $\epsilon^{-1}uH_0(T) + O(u \log \log \sigma)$ bits of space, a counting time of $O(m \log u)$, a locating time of $O(\log^\epsilon u)$ per occurrence reported, and an extracting time $O(\ell + \log^\epsilon u)$ for any text substring of length ℓ , where $0 < \epsilon \leq 1$ is any constant. We use the code provided at http://pizzachili.dcc.uchile.cl/indexes/Compressed_Suffix_Array/sada_csa.tgz. We have two parameters to set up for this index. The first one is the sample rate of suffix array positions (this information is used to speed up the locating and extracting operations), and the second one is the sample rate of the Ψ function. For `exists` and `count` queries we do not store any suffix array position, but only Ψ values. For `locate` queries, we used values of 4, 8, 16, 32, and 64 for suffix array positions, and the value 128 to sample the Ψ function (as this has shown to be the most efficient alternative [Ferragina et al. 2007]).

Alphabet Friendly FM-index (AF-FMI). This index [Ferragina et al. 2007] is based on the *backward-search* concept [Ferragina and Manzini 2005]. It has a space requirement of $uH_k(T) + o(u \log \sigma)$ bits, a counting time $O(m)$, a locating time $O(\log^{1+\epsilon} u)$ per occurrence reported, and an extracting time $O(\ell + \log^{1+\epsilon} u)$, for $\sigma = O(\text{polylog}(u))$, any constant $\epsilon > 0$, and any $k \leq \alpha \log_\sigma u$, where $0 < \alpha < 1$ is any constant. We use the code provided at http://pizzachili.dcc.uchile.cl/indexes/Alphabet-Friendly_FM-Index/af-index_v2.tgz. We have only one parameter to set up in the code, which is the sample rate of suffix array positions. We have used sample rates of one suffix array position stored out of 4, 8, 16, 32, and 64 text positions. For `exists` and `count` queries we do not store any suffix array position.

Succinct Suffix Array (SSA). This index [Mäkinen and Navarro 2005] is also based on backward search, but uses only one *wavelet tree* [Grossi et al. 2003], achieving $uH_0(T) + o(u \log \sigma)$ bits of space. The time complexities of this index are just as for the AF-FMI. We use the code provided at http://pizzachili.dcc.uchile.cl/indexes/Succinct_Suffix_Array/SSA_v2.tgz. We use the same parameters as for the AF-FMI.

Inverted LZ-index (ILZI). This index [Russo and Oliveira 2008] is based on Lempel-Ziv compression, using a variant of the LZ78 parsing called *maximal parsing*, which has many interesting properties. The ILZI requires at most $(5+\epsilon)uH_k(T) + o(u \log \sigma)$ bits of space, and has a locating complexity of $O((m + occ) \log u)$. We use the implementation by Luis Russo, which does not conform the *Pizza&Chili* API as the other indices. In particular, pattern occurrences are reported in the format `[[t, offset]]`, just like for the original LZ-index. The index does not include the data structure to transform those occurrences into real text positions, as our implementations do. To be fair, we sum the space of the described data structure for text positions to the space of this index. We also sum the average time to transform occurrences into real text positions for `locate` queries. According to our experiments, this is 1.3 microseconds per occurrence. As an additional consequence, this index does not provide `extract` queries, but only `display` queries. So, we are not able to extract arbitrary text substrings.

It is important to note also that the current implementation of the ILZI does not return to the invoking application an array with the pattern occurrences, as required by the *Pizza&Chili* API. This implementation just prints the number of phrases containing the starting position of the occurrences. We get rid of the print operation in the code, and thus there is no any reporting operation (we just find the occurrences). In particular, we are not accounting for the overhead of managing the occurrence array, which grows dynamically as more occurrences are found.

Russo and Oliveira [2008] define a practical variant of LZ78 parsing, the so-called *LZ78 maximal parsing with quorum l* . The idea is that for every phrase $B_i = B_j \cdot c$, B_j is the longest prefix of the rest of the text that appears at least $l + 1$ times in $B_0 \dots B_{i-1}$. Note that by using $l = 0$ we get the original LZ78 parsing. In this way, by using larger quorum values we can reduce the number of phrases in the LZ78 maximal parsing, hence reducing the number of nodes in the trie representing those phrases, and thus reducing the space of the index. We use quorum values $l = 0, 1, 2, 4, 8$, and 16 to get different space/time trade-offs, though this is not actually a trade-off parameter, but an optimization parameter, as we shall see in our experiments. Smaller values of l do not yielded a significant reduction in the space requirement.

6.2 Comparison of Space Requirement and Construction Time

In Table I we show the construction time and final space requirement for the indices we have tested. As it can be seen, we have reduced the space of the original LZ-index, and in the case of Scheme 3 we offer a space/time trade-off. Note that the maximum space requirement of Scheme 3 is about the same as that of the original LZ-index, and that the minimum space requirement we achieve is in all cases around 66% the space of the original LZ-index.

In many cases, such as for XML documents and DNA data, our indices are smaller than the original text. This is important by itself since we are able to provide indexing capabilities with a representation which is smaller than the original text.

We also conclude that our indices are much faster to build than competing schemes. It can be argued that construction time is not so important in indexed text searching, where one constructs the index once and queries it several times, so construction time is amortized over a number of queries. However, as we deal with very large (and sometimes huge) texts (because we would use a classical index otherwise), construction time is not irrelevant. As a comparison between LZ-based schemes, the ILZI (which is in most cases the slowest index to build) is built about 9 times slower than our schemes (e.g., in the case of English text). One reason for this is that our indices are constructed by performing only one pass over the text, while the ILZI needs two passes [Russo and Oliveira 2008]. Recall that the ILZI does not construct any text-position data structure, which would further increase its construction time.

In Table II we show the experimental size of the data structures described in Section 5.2 for reporting text positions in our LZ-indices. This space is already accounted for in Table I.

Table I. Range of space requirement and construction time for the different indices we have tested. The space is shown as a fraction of text size. Indexing time is shown in megabytes per second.

Text	Index	Range of space req. (fract. of text size)	Indexing speed (MB per second)
English	CSA	0.43 – 1.50	0.45 – 0.44
	SSA	0.87 – 2.87	0.93 – 0.90
	AF-FMI	0.65 – 2.65	0.26
	ILZI	1.23	0.15
	Original LZ-index	1.69	1.47
	Scheme 1	1.39	1.30
	Scheme 2	1.38	1.27
	Scheme 3	1.13 – 1.69	0.91 – 1.33
DNA	Scheme 4	1.13 – 1.69	0.71 – 1.31
	CSA	0.46 – 1.53	0.51
	SSA	0.50 – 2.50	1.33 – 1.28
	AF-FMI	0.48 – 2.48	0.43
	ILZI	0.95	0.66
	Original LZ-index	1.24	2.35
	Scheme 1	1.03	2.02
	Scheme 2	1.01	1.99
MIDI Pitches	Scheme 3	0.83 – 1.24	1.36 – 2.12
	Scheme 4	0.83 – 1.24	1.03 – 2.06
	CSA	0.62 – 1.68	0.94 – 0.92
	SSA	1.04 – 3.04	1.80 – 1.71
	AF-FMI	0.93 – 2.94	0.36
	ILZI	1.86	0.24
	Original LZ-index	2.58	1.76
	Scheme 1	2.16	1.49
XML	Scheme 2	2.12	1.47
	Scheme 3	1.76 – 2.58	0.88 – 1.58
	Scheme 4	1.76 – 2.58	0.62 – 1.56
	CSA	0.29 – 1.35	0.68
	SSA	0.98 – 2.98	1.34 – 1.29
	AF-FMI	0.54 – 2.54	0.44 – 0.43
	ILZI	0.61	0.34
	Original LZ-index	0.93	2.38
Proteins	Scheme 1	0.77	2.15
	Scheme 2	0.76	2.15
	Scheme 3	0.63 – 0.93	1.57 – 2.23
	Scheme 4	0.63 – 0.93	1.24 – 2.17
	CSA	0.67 – 1.73	0.57 – 0.56
	SSA	0.82 – 2.82	0.97 – 0.95
	AF-FMI	0.82 – 2.82	0.38 – 0.37
	ILZI	1.73	0.24
Sources	Original LZ-index	2.40	1.55
	Scheme 1	1.99	1.29
	Scheme 2	1.96	1.25
	Scheme 3	1.62 – 2.40	0.79 – 1.35
	Scheme 4	1.62 – 2.40	0.58 – 1.30
	CSA	0.38 – 1.44	0.76 – 0.75
	SSA	1.01 – 3.01	1.37 – 1.32
	AF-FMI	0.73 – 2.73	0.30
	ILZI	1.15	0.17
	Original LZ-index	1.67	1.73
	Scheme 1	1.39	1.54
	Scheme 2	1.37	1.51
	Scheme 3	1.13 – 1.67	1.04 – 1.59
	Scheme 4	1.13 – 1.67	0.79 – 1.54

Table II. Size of the data structure for reporting text positions in our LZ-indices.

Text	Size of text-position data structure (as a fraction of text size)
English	0.14
DNA	0.11
MIDI Pitches	0.27
XML	0.09
Proteins	0.22
Sources	0.16

6.3 Comparison of Search Time

Next we experimentally test whether the trade-offs we provide are competitive for compressed text searching. In our experiments, we call S2 DFUDS the version of Scheme 2 implemented on DFUDS, and S3 DFUDS the DFUDS version of Scheme 3. We do not include Scheme 1 based on DFUDS, since it is outperformed by S2 DFUDS, both of them requiring about the same space. We do not include Scheme 4 in our plots, since Scheme 3 outperforms it in most cases (though they require almost the same space). However, Scheme 4 is interesting by itself since it can be reduced to $(1 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits of space [Arroyuelo et al. 2006; 2008], which cannot be achieved by other schemes.

6.3.1 Extract Queries. Since compressed full-text self-indices replace the text, the fast extraction of arbitrary text substrings is essential in most applications, and thus one of the most relevant problems the indices face. To test the performance of our indices, we extract 10,000 random snippets, each of length 100. We measure the time per symbol extracted, so we average over a total of 1 million extracted symbols.

In Fig. 9 we show the experimental results. As we can see, our indices are very competitive in the range of space they require, in most cases outperforming competing schemes, which in some cases cannot compete even if using more memory than ours. In the particular case of DNA text, our indices are competitive against the SSA, which in the case of small alphabets is a very good alternative. This is because the corresponding *wavelet tree* [Grossi et al. 2003], on which the SSA is based, is shallow, and then each text symbol can be extracted very fast (every symbol is extracted in time proportional to the height of the corresponding wavelet tree).

This shows the superiority of our LZ-indices in this important aspect. Also, we can see that our approach to reduce the space of the LZ-index is effective in this case, since we are able to reduce the space while still maintaining a good extracting performance.

6.3.2 Display Queries. In many applications displaying a context surrounding the occurrences is as important as (and sometimes more important than) the occurrence positions themselves. This is not a problem for classical indices, since the text is available and hence we can get the contexts from the occurrence positions obtained with `locate`. In the case of self-indices, on the other hand, one must ask the index to reproduce the occurrence contexts, which is achieved with `display`

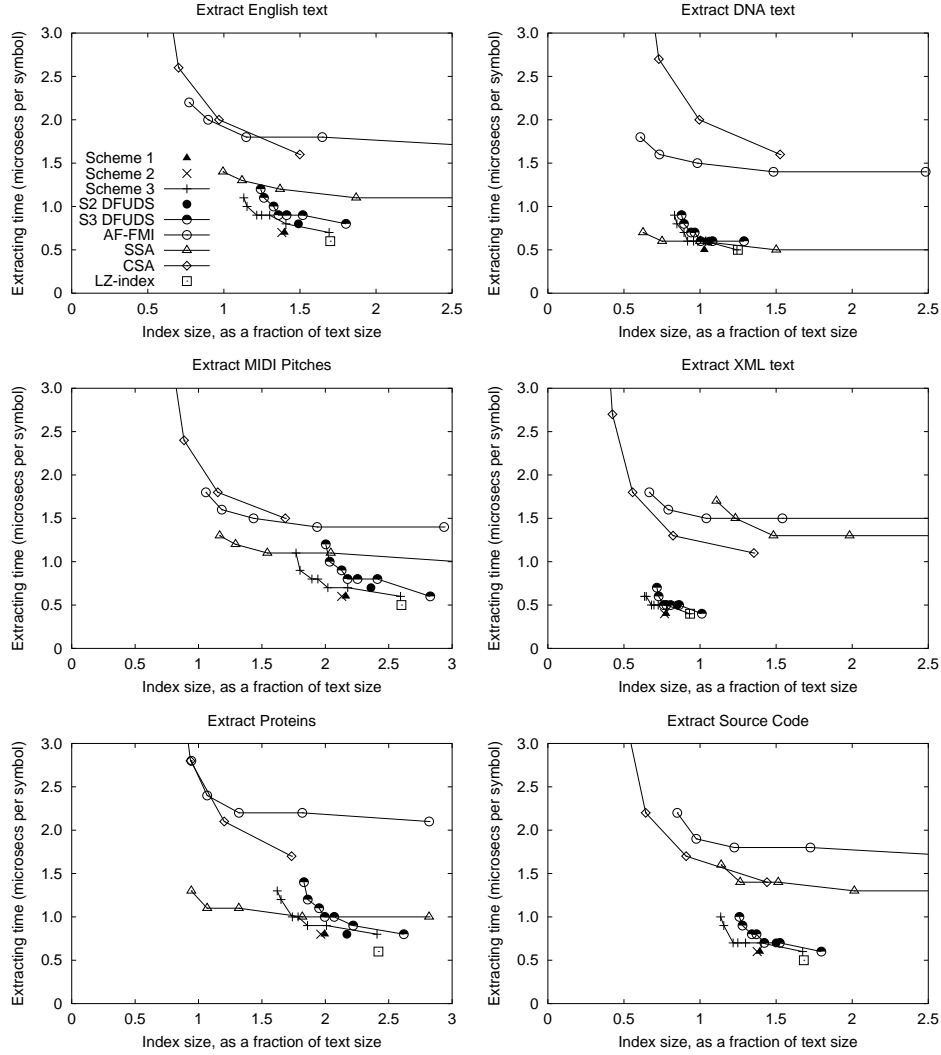


Fig. 9. Experimental extracting time, for random snippets of length $\ell = 100$. Times are measured in microseconds per symbol extracted.

queries since we do not have the text at hand. Thus, fast displaying the occurrence contexts is also important for a self-index.

To test our indices, we search for 5 million pattern occurrences and then show a context of 50 symbols surrounding every occurrence, for patterns of length 10 (in other words, we display 110 symbols per occurrence).

In most indices, $\text{display}(P, \ell)$ queries can be thought of as a $\text{locate}(P)$ query (in order to find the pattern occurrences) followed by an $\text{extract}(i, j)$ query (where i and j are computed by means of the positions obtained with locate and the context length ℓ). In our LZ-indices, however, we originally get the occurrences in the format $\llbracket t, \text{offset} \rrbracket$, to finally transform t and offset into a text position (by

means of the data structure described in Section 5.2). This text position must be transformed by `extract` again into an LZ78 phrase (recall that this involves binary searching the text-position data structures), from where we start the extraction of text in *LZTrie*. To avoid repeating this work, we do not transform the occurrences into text positions when performing `display` queries. We rather display text with a simplified version of `extract` that works on LZ78 phrases rather than on text positions.

In Fig. 10 we show the experimental results. The current implementation of the ILZI shows only a context preceding the pattern occurrences, and not surrounding the occurrences as other schemes do. For our LZ-indices, showing a context surrounding the occurrences (which is usually required) introduces the use of extra operations which are not needed when showing a context preceding an occurrence. As it can be seen, just like for `extract` queries, our indices are among the most competitive schemes for displaying text, in many cases outperforming the competing ILZI.

6.3.3 Locate Queries. For `locate` queries we search for 10,000 patterns extracted at random positions from the text, with length 5, 10, and 15. For short patterns, we limit the total number of occurrences found to 5 million. We measure the time in microseconds per occurrence found.

As we said before, `locate` queries are important in classical full-text indexing, where one has the text at hand in order to access the occurrences and the contexts surrounding them. This is somehow equivalent to `display` queries in the scenario of compressed full-text self-indices. Obtaining just text positions (and nothing else) can be interesting in specific cases, but `display` queries have broader applications.

Partial locate Queries. In many applications we do not need to locate all of the pattern occurrences at once, but just a few of them. This is a very challenging problem for our LZ-indices since, for example, the indices based on suffix arrays are very efficient to find the suffix-array interval containing the occurrences, and hence they are rapidly ready to start locating the occurrences. The ILZI has also a very fast $O(m)$ trie navigation before starting the locating procedure.

We test here our algorithm defined in Section 5.3 (recall that we divide the search process into four levels, level 0 up to 3). In Table III we show the percentage of occurrences that are found in each level of search, for $K = 1$ and for the different text collections. Notice that the search of level 0 (i.e., searching for P as a whole phrase in *LZTrie*) is very effective for patterns of length 5 to 10 (in the case of DNA, for example, almost 100% of the queries can be answered at level 0). Notice also that the percentage found at level 1 is relatively small compared to the corresponding percentage of level 0. Recall that with the search of level 0 we look for a particular case of occurrences of type 1. This means that most of the times a pattern exists as an occurrence of type 1, this can be found at level 0 of the search. For longer patterns, $m = 15$, there is a smaller probability of finding the occurrence contained in a single phrase, and thus the percentage found at level 0 is smaller.

In Figs. 11 up to 14 we show the experimental result for values $K = 1$ and 5, and for $m = 5$ and 10. We do not show the results for $m = 15$ since the results are poorer, as predicted by the results of Table III.

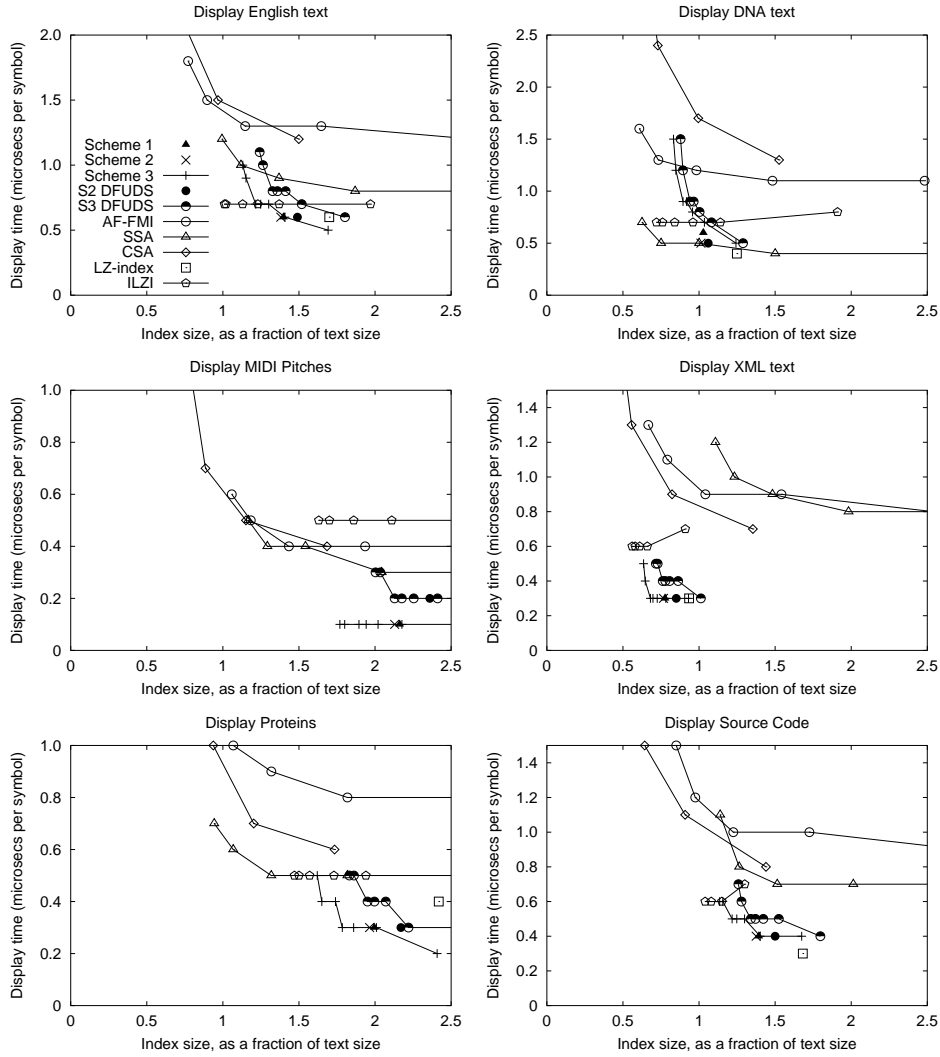


Fig. 10. Experimental display time, for snippets of length $\ell = 110$ around every pattern occurrence. Times are measured in microseconds per symbol extracted.

For $m = 5$ and $K = 1$, the most interesting results are obtained in the cases of English text, DNA, XML, and proteins, though in the latter case our indices do not obtain good compression. For DNA, S3 DFUDS is unbeatable since, as we have seen in Table III, 100% of the occurrences are found at level 0 of the search. Notice that for MIDI pitches our performance is not so competitive (except perhaps for S2 DFUDS), as it was also predicted by Table III: only 73% of the occurrences are found at level 0, and thus we need more trie navigations.

Except for MIDI pitches, in all other cases the best alternatives are the indices based on Lempel-Ziv compression: ours and the ILZI (the best in each case depends on the space usage, as it can be seen). This shows that Lempel-Ziv-based indexing is

Table III. Percentage of patterns found at each level of search, for partial `locate` queries with $K = 1$ and for different pattern lengths. Level 0: P is found as an LZ78 phrase in *LZTrie*; Level 1: P is found as occurrence of type 1, but not as a whole phrase; Level 2: P is found as occurrence of type 2; Level 3: P is found as occurrence of type 3.

Text	Level	Percentage per level		
		$m = 5$	$m = 10$	$m = 15$
English	0	98.64	55.86	7.80
	1	0.54	6.27	1.83
	2	0.82	34.97	63.14
	3	0.0	2.90	27.23
DNA	0	100.0	99.26	11.63
	1	0.0	0.33	1.06
	2	0.0	0.41	78.74
	3	0.0	0.0	8.57
MIDI Pitches	0	72.28	20.01	9.69
	1	3.69	2.10	1.38
	2	23.82	42.66	19.67
	3	0.21	35.23	69.26
XML	0	95.41	65.61	37.56
	1	2.59	11.53	15.53
	2	1.99	19.55	32.13
	3	0.01	3.31	14.78
Proteins	0	98.96	13.29	8.51
	1	0.71	1.41	1.29
	2	0.33	63.87	11.23
	3	0.0	21.43	78.97
Sources	0	94.31	48.40	21.41
	1	2.73	9.34	6.01
	2	2.94	37.15	44.11
	3	0.02	5.11	28.47

a very competitive choice in general, despite that suffix-array-based indices basically compute the suffix-array interval containing the occurrences, and then we ask the index to obtain just one of these occurrences. As we shall see later with `count` queries, these indices are very efficient to find the suffix-array interval; however, asking them to obtain just one occurrence makes them significantly less competitive.

Notice also that our indices based on `DFUDS` outperform (in most cases by far) our indices based on `BP`. This is because the fast navigation on the tries becomes a fundamental aspect, since we are reporting just a few occurrences. For this reason, and in order to make our plots clearer, we do not show the results for Scheme 1, Scheme 2, and the original LZ-index.

As we increase the number of occurrences to locate, in principle the trie navigations are amortized by reporting more occurrences. However, our indices may need to go on more levels of the search, which means more navigations on the tries. Therefore, the total cost depends on the number of occurrences found in each search level. In general, as K grows, it becomes more difficult to compete since occurrences at higher levels are more expensive to obtain. Yet, we still provide some interesting cases for $K = 5$ and $K = 10$ (the latter case is not shown; results

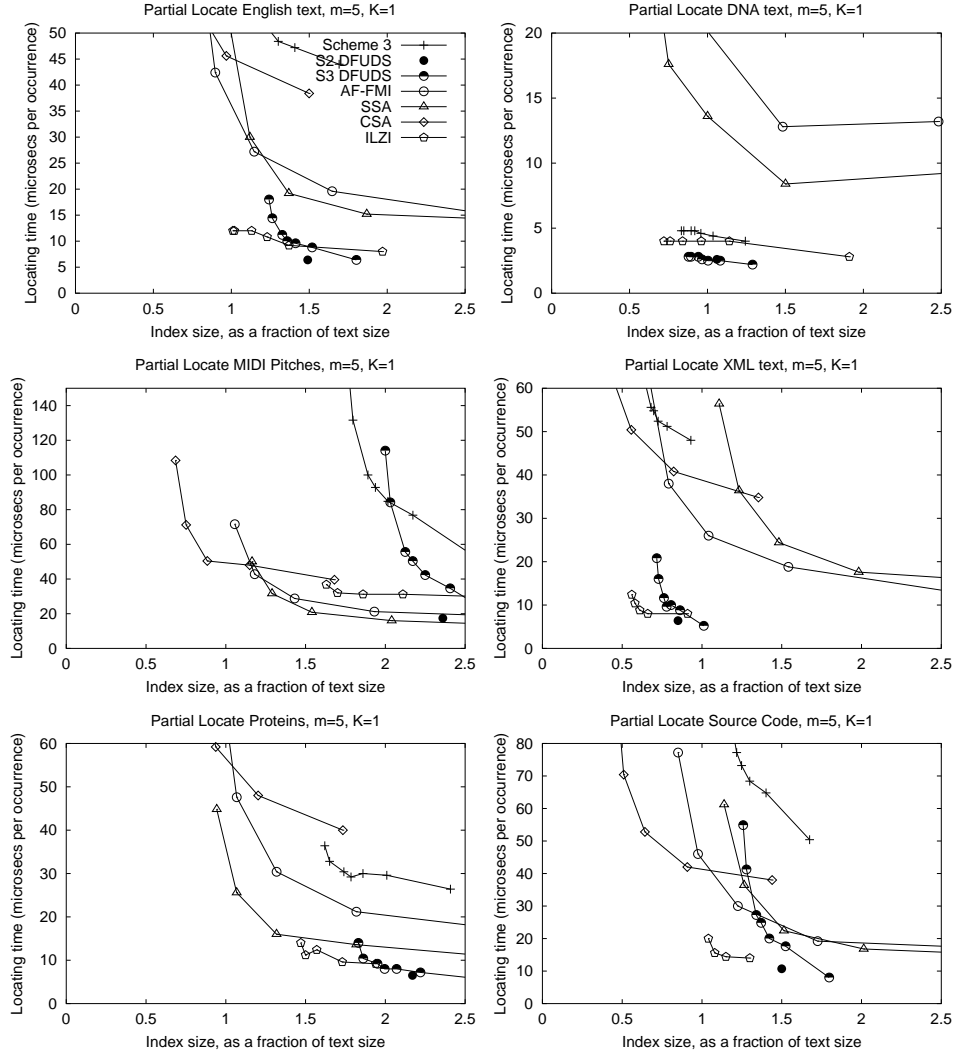


Fig. 11. Experimental time for partial `locate` queries, for patterns of length $m = 5$ and retrieving just $K = 1$ occurrence. We measure the time in microseconds per occurrence reported.

are close to those for $K = 5$).

The trie traversals also raise when we increase the pattern length, as it can be seen for $m = 10$.

Thus, we conclude that our technique for solving partial `locate` queries of Section 5.3 is relevant, and more efficient when the probability of finding the occurrences at level 0 of the search is high, as for example when we search for short patterns, the alphabet is small, or we look for very few occurrences (e.g., $K = 1$).

Full locate Queries. We test here `locate` queries without limiting the number of occurrences. In Fig. 15 we show the experimental results for patterns of length

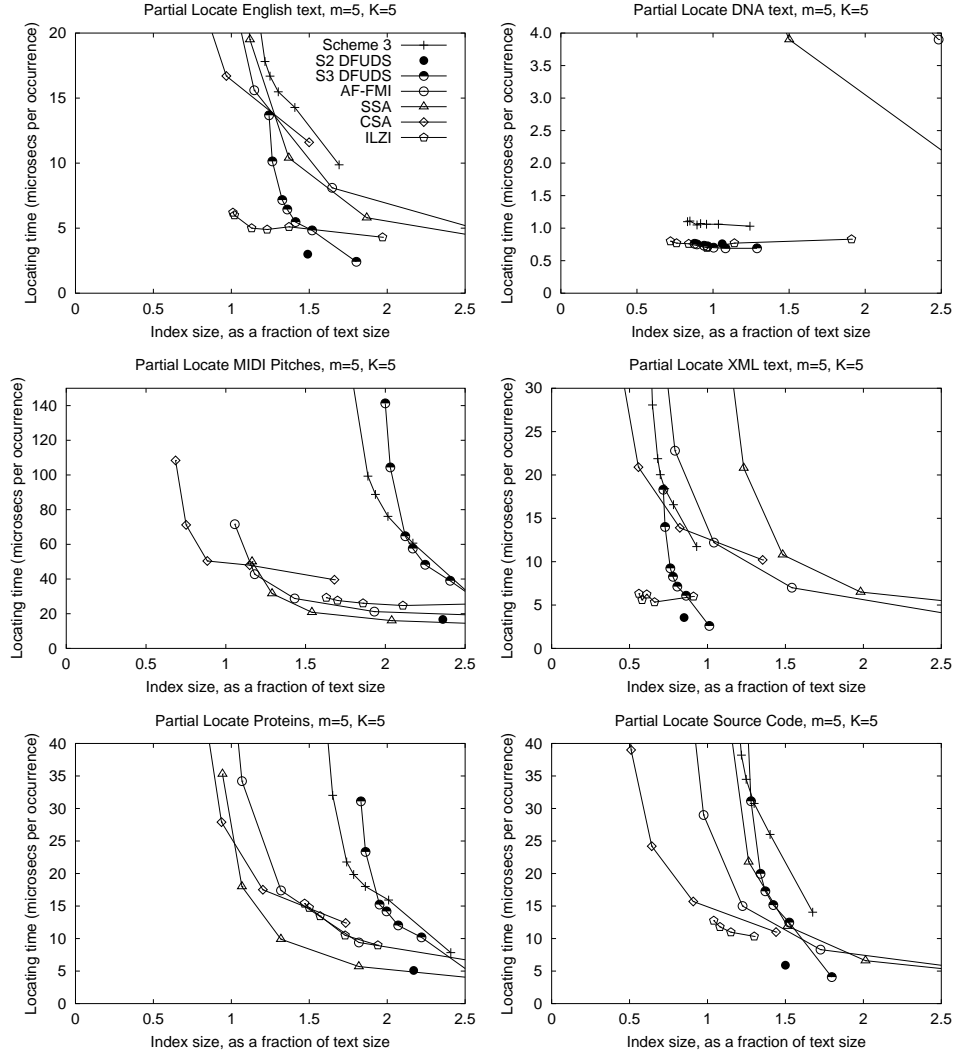


Fig. 12. Experimental time for partial `locate` queries, for patterns of length $m = 5$ and retrieving just $K = 5$ occurrences. We measure the time in microseconds per occurrence reported.

$m = 5$. As we can see, there is no clear winner in all cases, but the performance depends on the available space. However, we can see some clear performance patterns: in most cases our schemes outperform all competing schemes (including the very competitive ILZI) as soon as we have space available to store, at least, Scheme 2. When we reduce the space usage of the index, however, the ILZI outperforms Scheme 3, yet the latter is still competitive (in most cases outperforming the competitive CSA). In general, for `locate` queries with short patterns the superiority of Lempel-Ziv-based indices is clear.

Notice also that, in most cases, Scheme 2 outperforms Scheme 1, both requiring about the same space. This means that the shorter average path length of Scheme ACM Journal Name, Vol. V, No. N, Month 20YY.

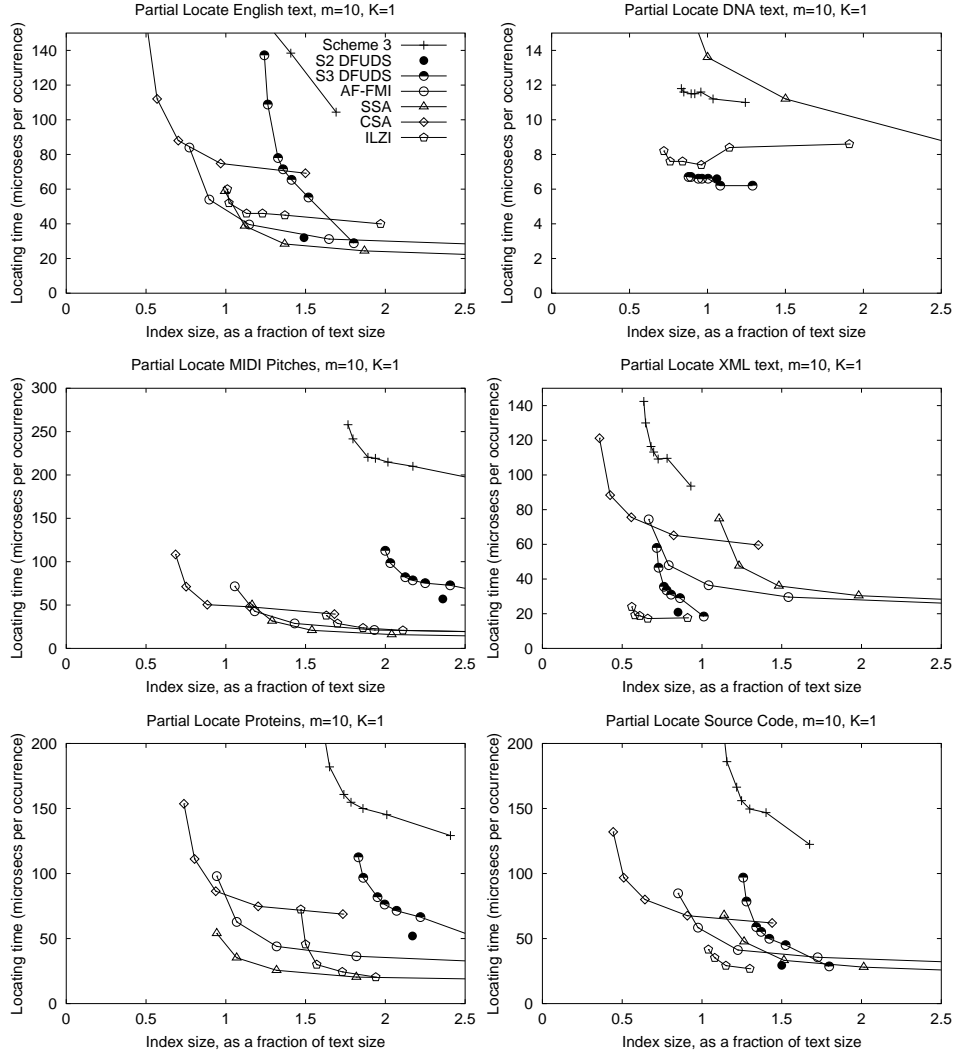


Fig. 13. Experimental time for partial `locate` queries, for patterns of length $m = 10$ and retrieving just $K = 1$ occurrence. We measure the time in microseconds per occurrence reported.

2 is better than having direct access from *RevTrie* to *LZTrie* nodes as in Scheme 1, which is good only for occurrences of type 1. As we said before, occurrences of type 2 are more costly in Scheme 1. It is also important to note that the DFUDS versions of LZ-index have a performance which is similar to the LZ-index based on BP. Notice, however, that DFUDS requires more space than BP, as it was predicted in Section 5.1.2.

In the particular case of XML text, we can see a very important aspect that differentiates LZ-indices from indices based on suffix arrays. The latter need to store extra non-compressible information (the sampled suffix array positions) to efficiently carry out `locate` and `extract` queries. The extra data stored by LZ-

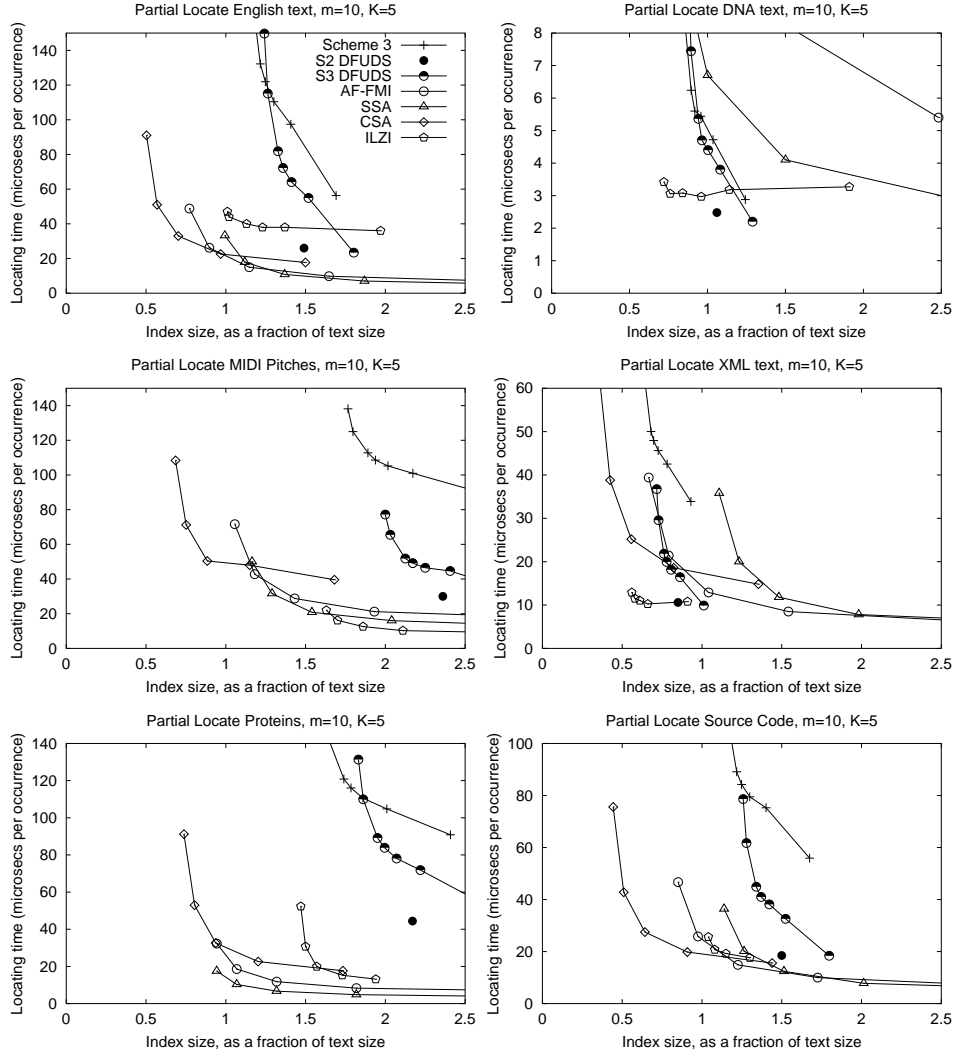


Fig. 14. Experimental time for partial `locate` queries, for patterns of length $m = 10$ and retrieving just $K = 5$ occurrences. We measure the time in microseconds per occurrence reported.

indices, on the other hand, is compressible, for example the size of the arrays for which we sample the inverse-permutation information depends on n , the number of LZ78 phrases of T . Therefore, when the texts are highly compressible, the LZ-indices can be smaller and faster than alternative indices.

For proteins, as an opposite case, our LZ-indices are larger and slower. They are large since the text is not as compressible as others, which can be also noted in the size of competing schemes. Our LZ-indices are in addition slower in this case, since each search retrieves on average only a few patterns, and therefore the work on the tries is not amortized by reporting many occurrences.

We show the results for patterns of length $m = 10$ in Fig. 16. As we can see, ACM Journal Name, Vol. V, No. N, Month 20YY.

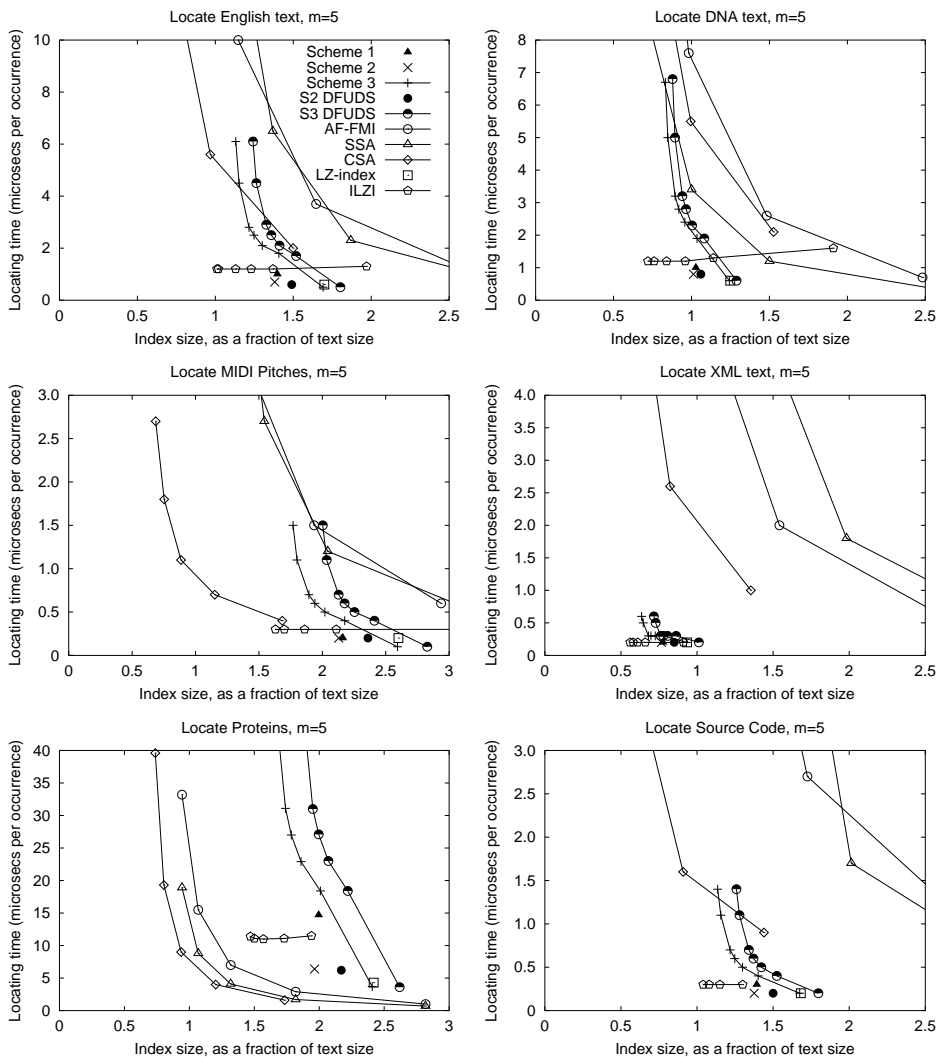


Fig. 15. Experimental locating time, for patterns of length $m = 5$. We measure the time in microseconds per occurrence reported.

our indices are still competitive, yet not as significantly as in the previous case where $m = 5$. The DFUDS implementation of LZ-index outperforms BP in all cases, except for MIDI pitches, XML text, and source code, where these have about the same performance. In particular for proteins and English text, where the number of occurrences per pattern is relatively small, the difference is greater for DFUDS, since the cost of navigating the tries becomes predominant. It is important to note also that Scheme 2 (both for BP and DFUDS implementations) is interesting (in some cases the best) alternative, for the memory space it requires. For patterns of length $m = 15$ (figure omitted), the behavior of the indices is similar to that of length $m = 10$.

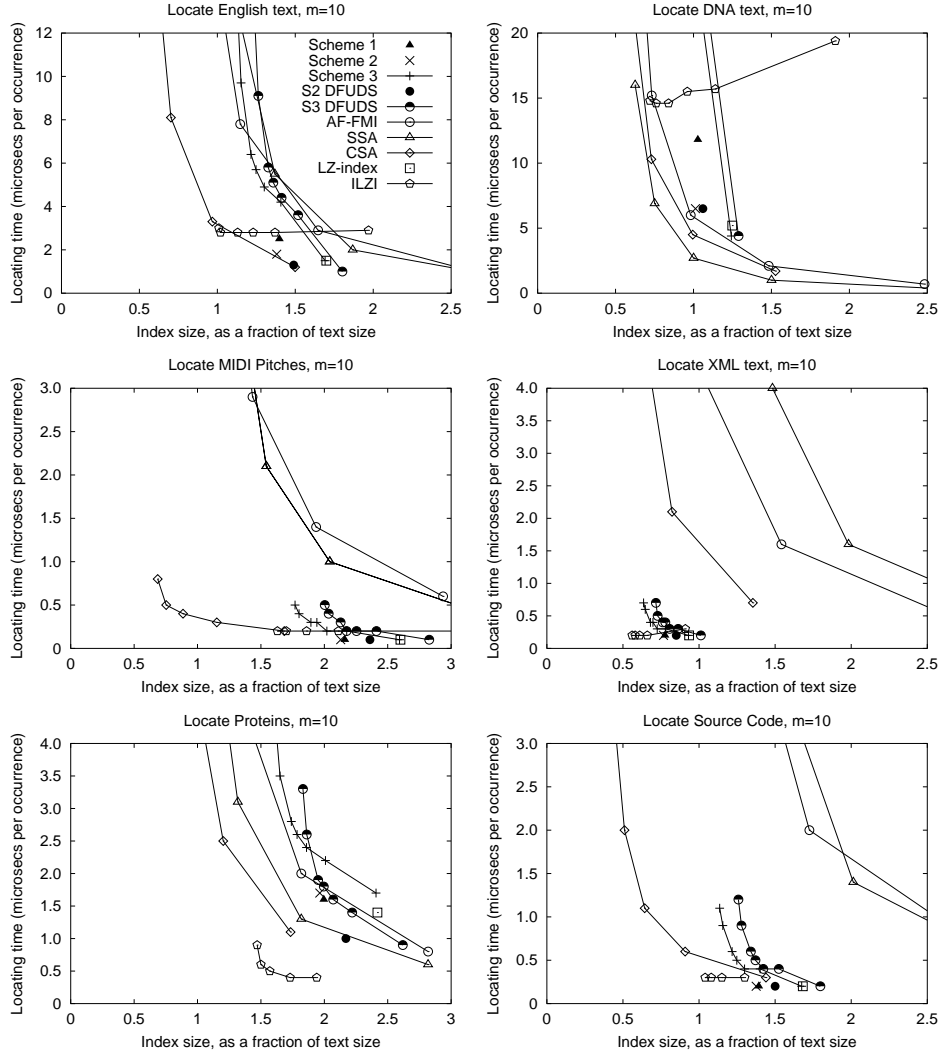


Fig. 16. Experimental locating time, for patterns of length $m = 10$. We measure the time in microseconds per occurrence reported.

Our results indicate that our LZ-indices, despite not being the best in many cases, are an attractive alternative for `locate` queries. In most cases we need at least the space required by Scheme 2 in order to outperform competing schemes. In particular, our LZ-indices are an interesting alternative for full `locate` when the total number of occurrences to report is considerable. For long patterns, the number of pattern occurrences becomes smaller, and therefore the time of searching for the $O(m^2)$ pattern substrings in the tries dominates. This problem can be somehow alleviated by using the DFUDS implementation for the tries. We can conclude that the LZ-indices based on BP and on DFUDS behave very similarly when reporting a great number of occurrences fast is an important issue (e.g., when $m = 5$). Yet

DFUDS requires more space. On the other hand, when the navigation in the tries becomes more important (e.g., for larger m), the DFUDS implementation offers a more attractive alternative, which can replace the BP implementation at the cost of a little extra space requirement.

6.3.4 Count Queries. This kind of queries has applications in much more specific and limited cases, generally conforming the internal machinery of more complex tasks.

In our experiments we search for patterns of length 20 extracted at random positions from the text. We measure the times per symbol of the pattern. As competing schemes based on suffix arrays do not store suffix-array sampling information in order to count (and thus are able to support efficiently only `count` queries), to be fair in this case we do not store the data structure for text positions in our LZ-indices. Yet, note that, within this space, we are still able to support fast `locate` queries (though without reporting text positions), and `display` queries.

The experimental results for `count` queries are shown in Fig. 17. As we can see, our schemes can implement this query, yet they cannot compete against the indices based on suffix arrays, since the counting complexity of these indices is related to the pattern length, and not to the number of pattern occurrences. Our schemes basically need to locate the pattern occurrences in order to count them. We can also see that the DFUDS implementation of LZ-index outperforms in all cases the BP implementation, yet the former requires slightly more space. This is mainly because we are searching for long patterns, and DFUDS provides more efficient descent in the *LZTrie*.

6.3.5 Exists Queries. In some specialized applications we just need to know whether pattern P exists or not in the text. Indices based on suffix arrays basically need to count the number of occurrences, since they first search for the pattern, and then check whether the suffix-array interval they get is empty or not. Despite that our LZ-indices are not competitive for `count` queries, we show here that they are much more efficient for finding the first occurrence of a pattern, which is useful to support `exists` queries (indeed, this has been already shown in the experiment with partial `locate` queries). The key is that we do not necessarily need to count the occurrences, but just to find the first pattern occurrence as fast as we can. We test with patterns of length 5, 10, and 15, and search for 10,000 patterns that exist in the text. We implement existential queries in our indices using the idea of partial `locate` queries, explained in Section 5.3.

In Fig. 18 we show the experimental results for `exists` queries, for patterns of length 5. Excluded plots for patterns of length $m = 10$ produced similar results, while those for patterns of length $m = 15$ gave worse results (this is because, as predicted in Table III, the heuristic of level 0 is not so effective for longer patterns). As it can be seen, our indices are much more competitive than for `count` queries, showing that our approach of first looking for P in *LZTrie* is effective in practice. Our indices achieve the same times (though not the same space) in many cases. As in the case of `count` queries, our indices are larger than competing schemes, yet ours are able to support more than just `count` and `exists` queries within this space.

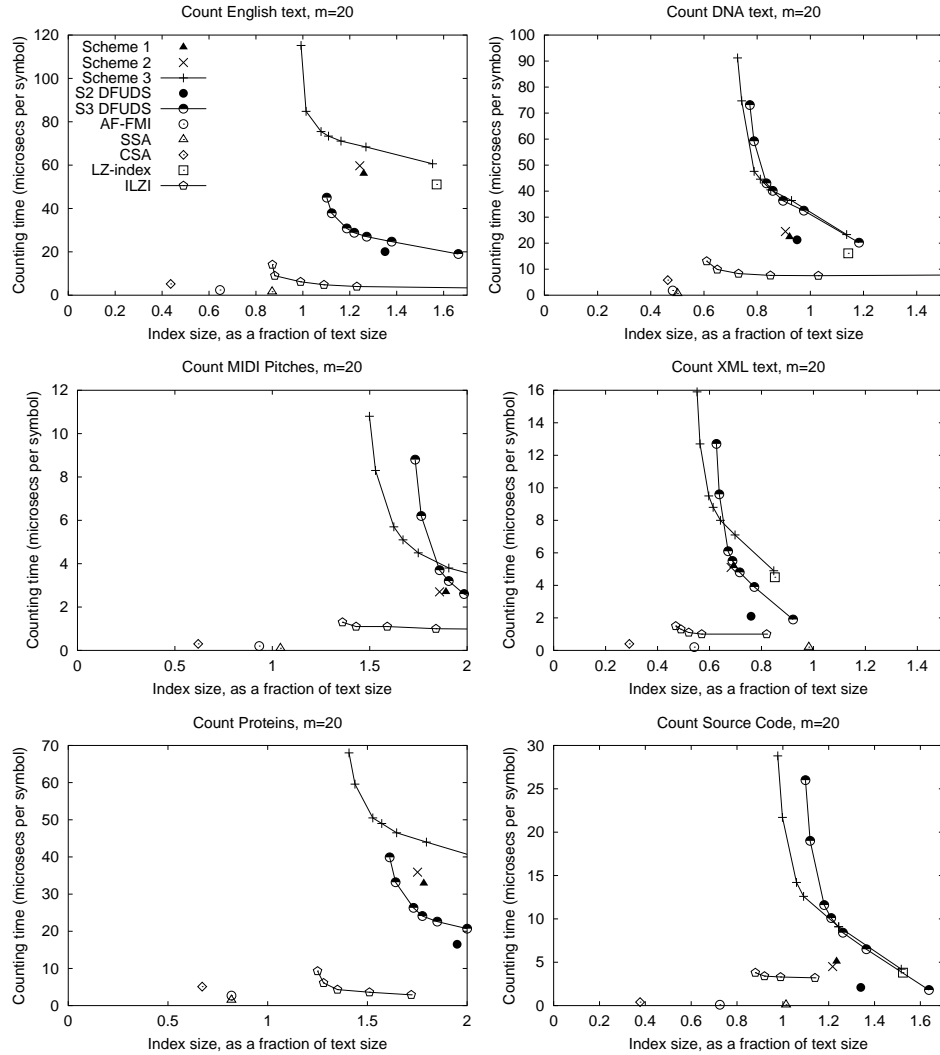


Fig. 17. Experimental counting time, for patterns of length $m = 20$. We measure the time in microseconds per symbol of the pattern.

7. CONCLUSIONS

We have defined a practical approach to reduce the space requirement of LZ-indices, which we call the *navigation-scheme approach*. In particular, we applied this method on Navarro’s LZ-index [Navarro 2004; 2008]. Our approach consists in regarding the original LZ-index as a navigation scheme that allows us moving among index components. From this point of view we can study the redundancy introduced by the original LZ-index components. Then, we define alternative schemes that allow for the same navigation, yet without the redundancy of the original one.

Given a text $T[1..u]$ over an alphabet of size σ , and with k -th order empirical

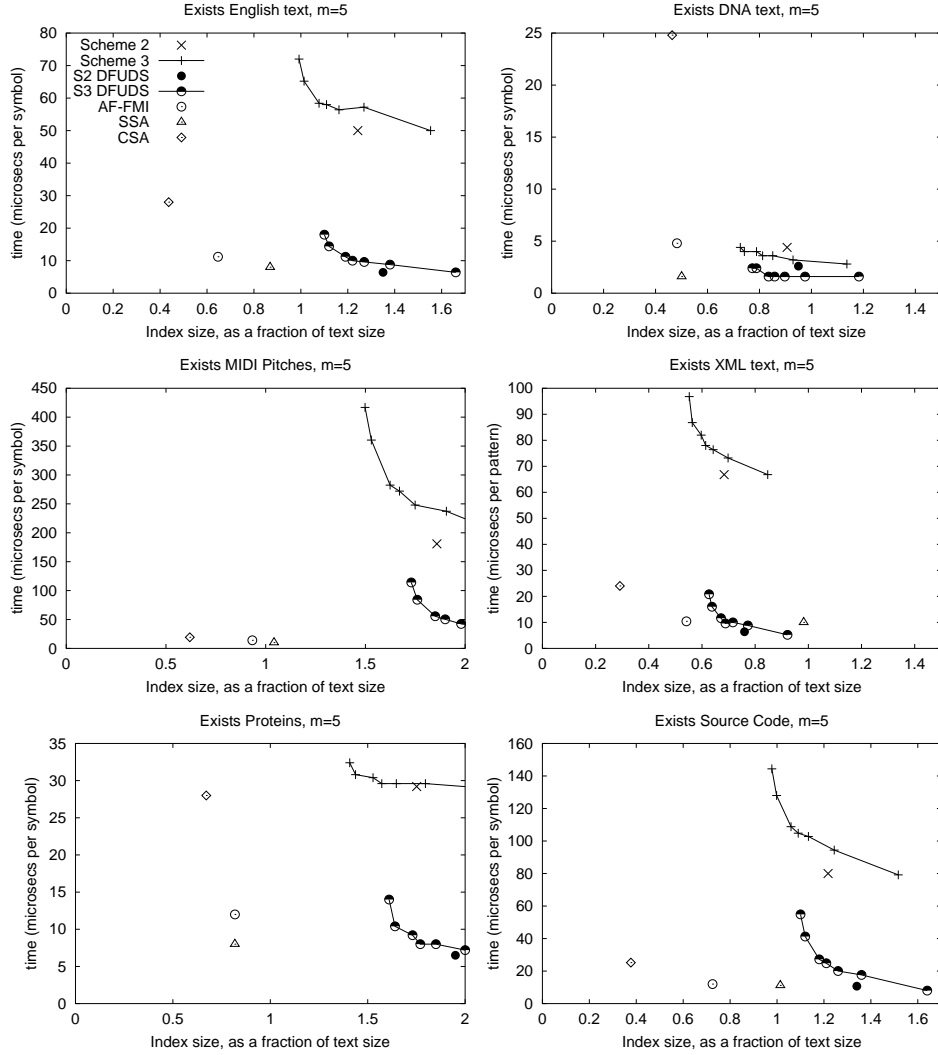


Fig. 18. Experimental time for `exists` queries, for patterns of length $m = 5$. We measure the time in microseconds per pattern.

entropy $H_k(T)$ [Manzini 2001], the original LZ-index of Navarro [2004] requires $4uH_k(T) + o(u \log \sigma)$ bits of space. In this paper we define several new versions of the LZ-index, requiring $3uH_k(T) + o(u \log \sigma)$ and $(2 + \epsilon)uH_k(T) + o(u \log \sigma)$ bits of space. The latter ones allow us to partially overcome one of the drawbacks of the original LZ-index: the lack of space/time tuning parameters. Although our schemes do not provide worst-case guarantees at search time, they ensure $O(\frac{m^2}{\epsilon} + \frac{n}{\epsilon n^{m/2}})$ average-case time for locating the occurrences of pattern $P[1..m]$ in T . This is $O(\frac{m^2}{\epsilon})$ for $m \geq 2 \log_\sigma u$.

We implemented and intensively tested our indices in many practical scenarios,

which cover an interesting range of applications, comparing against the best existing alternatives for compressed full-text self-indexing. From those experiments we conclude that we can effectively reduce the space requirement of the original LZ-index, to about 2/3. We also noted that our indices are the fastest to build, which is important when we deal with very large texts.

When comparing the search performance, we concluded that our indices offer an interesting alternative in practice, for a number of different types of queries: our indices are very competitive (in most cases the best) alternatives for **extract** and **display** queries, which we argue are the most basic queries in the scenario of compressed full-text self-indices, where the text is not available otherwise². For instance, in most cases we are able to extract about 1 to 1.5 million symbols per second, being about twice as fast as the most competitive alternatives. Thus, we can reduce the space, still maintaining the competitiveness of the original LZ-index in this respect.

For **locate** queries we tested two alternatives: partial **locate** queries (where a fixed number of occurrences is located) and full **locate** queries. Our experiments show that our technique for solving partial **locate** queries is efficient in cases of searching for short patterns, of texts with small alphabets, or when we want to locate very few occurrences (e.g., only one occurrence).

For full **locate** queries, we showed that our indices are more effective when the search pattern is not too long, or there are many occurrences to report. In other cases, the $O(m^2)$ -time navigation on the tries is predominant, and thus our performance degrades. For example, for short patterns of length 5, in most scenarios our schemes are the best alternative if we can spend at least 80% the size of the original LZ-index. When less memory space is available, our indices are outperformed by the very competitive Inverted LZ-index (ILZI) [Russo and Oliveira 2008], which is yet another variant of the Lempel-Ziv-based index family. In this case, however, our indices are still competitive with all suffix-array-based schemes (e.g., the competitive Compressed Suffix Array of Sadakane [2003]).

It is well known that the original LZ-index is more adequate for short patterns [Navarro 2008]. We somehow alleviated the problem of long patterns, by using a more efficient representation, so-called DFUDS [Benoit et al. 2005], for the tries that compose the LZ-index. From our experiments we can conclude that our DFUDS implementation allows us for very efficient search, in many cases outperforming the traditional balanced-parentheses representation [Munro and Raman 2001] of the LZ-index. However, more work is needed to compete in this respect against suffix-array-based indices. We hope to get further improvements on this line when working on alphabets larger than those tried in this paper, for example, if defining an LZ-index working on words, for applications of natural-language processing.

We also exhibit an important difference between LZ-indices and those based on suffix arrays: the latter need to store extra non-compressible information (the suffix-array samples) in order to carry out **extract**, **display**, and **locate** queries, whereas the information stored by our LZ-indices is compressible. Thus, when the

²In the literature [Navarro and Mäkinen 2007] the **count** operation is taken as the most basic one, but this is probably biased towards suffix-array-based indices, more than to considering typical applications.

text is highly compressible, we get very small LZ-indices, which are still fast. Indices based on suffix arrays, on the other hand, cannot use a denser suffix-array sampling (because otherwise they become larger), and henceforth their performance is poor. Thus, our indices are also a very good option for highly compressible texts.

We do not consider in this paper the space-efficient construction of our indices, which is a very important issue in practice, since many times a small index requires a large amount of memory space to be build. We plan to adapt the space-efficient algorithm by Arroyuelo and Navarro [2005] for the original LZ-index, so as to construct our LZ-index variants space-efficiently. Currently, our indices are constructed in an uncompressed way, needing about the same space used to build suffix-array-based compressed self-indices. The space-efficient construction of the latter is still at a theoretical stage [Mäkinen and Navarro 2008], or still does not achieve higher-order entropy-bound space [Hon et al. 2007].

We believe that our indices offer an extremely relevant alternative considering their overall performance across the multiple tasks of interest in many real text-search applications.

We made the code of our indices available in the *Pizza&Chili* corpus, throughout the site <http://pizzachili.dcc.uchile.cl/indexes/LZ-index/>.

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. 16th Annual International Symposium on Algorithms and Computation (ISAAC)*. LNCS 3827. 1143–1152.
- ARROYUELO, D., NAVARRO, G., AND SADAKANE, K. 2006. Reducing the space requirement of LZ-index. In *Proc. 17th Annual Symposium on Combinatorial Pattern Matching (CPM)*. LNCS 4009. 319–330.
- ARROYUELO, D., NAVARRO, G., AND SADAKANE, K. 2008. Stronger Lempel-Ziv based compressed text indexing. Tech. Rep. TR/DCC-2008-2, Dept. of Computer Science, University of Chile.
- BARBAY, J., HE, M., MUNRO, J. I., AND RAO, S. S. 2007. Succinct indexes for strings, binary relations and multi-labeled trees. In *Proc. 18th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 680–689.
- BENOIT, D., DEMAINE, E., MUNRO, I., RAMAN, R., RAMAN, V., AND RAO, S. S. 2005. Representing trees of higher degree. *Algorithmica* 43, 4, 275–292.
- FERRAGINA, P., GONZÁLEZ, R., NAVARRO, G., AND VENTURINI, R. 2007. Compressed text indexes: From theory to practice! Submitted. http://pizzachili.dcc.uchile.cl/resources/compressed_indexes_in_practice.pdf.
- FERRAGINA, P. AND MANZINI, G. 2000. Opportunistic data structures with applications. In *Proc. 41st Annual Symposium on Foundations of Computer Science (FOCS)*. 390–398.
- FERRAGINA, P. AND MANZINI, G. 2005. Indexing compressed texts. *Journal of the ACM* 54, 4, 552–581.
- FERRAGINA, P., MANZINI, G., MÄKINEN, V., AND NAVARRO, G. 2007. Compressed representations of sequences and full-text indexes. *ACM Transactions on Algorithms (TALG)* 3, 2, article 20.
- FERRAGINA, P. AND NAVARRO, G. 2005. Pizza&Chili Corpus — Compressed indexes and their testbeds. <http://pizzachili.dcc.uchile.cl>.
- GEARY, R., RAHMAN, N., RAMAN, R., AND RAMAN, V. 2006. A simple optimal representation for balanced parentheses. *Theoretical Computer Science* 368, 3, 231–246.
- GOLYSKI, A., MUNRO, J. I., AND RAO, S. 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.

- GONZÁLEZ, R., GRABOWSKI, S., MÄKINEN, V., AND NAVARRO, G. 2005. Practical implementation of rank and select queries. In *Poster Proc. 4th International Workshop on Efficient and Experimental Algorithms (WEA)*. CTI Press and Ellinika Grammata, 27–38.
- GROSSI, R., GUPTA, A., AND VITTER, J. S. 2003. High-order entropy-compressed text indexes. In *Proc. 14th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 841–850.
- GROSSI, R. AND VITTER, J. S. 2000. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. In *Proc. 32nd Annual ACM Symposium on Theory of Computing (STOC)*. 397–406.
- GROSSI, R. AND VITTER, J. S. 2005. Compressed suffix arrays and suffix trees with applications to text indexing and string matching. *SIAM Journal on Computing* 35, 2, 378–407.
- HON, W. K., LAM, T. W., SADAKANE, K., SUNG, W.-K., AND YIU, M. 2007. A space and time efficient algorithm for constructing compressed suffix arrays. *Algorithmica* 48, 1, 23–36.
- JANSSON, J., SADAKANE, K., AND SUNG, W.-K. 2007. Ultra-succinct representation of ordered trees. In *Proc. 18th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 575–584.
- KÄRKKÄINEN, J. AND UKKONEN, E. 1996. Lempel-Ziv parsing and sublinear-size index structures for string matching. In *Proc. 3rd South American Workshop on String Processing (WSP)*. Carleton University Press, 141–155.
- KIM, D., NA, J., KIM, J., AND PARK, K. 2005. Efficient implementation of rank and select functions for succinct representation. In *Proc. 4th International Workshop on Efficient and Experimental Algorithms (WEA)*. LNCS 3503, 315–327.
- KNESSL, C. AND SZPANKOWSKI, W. 2000. Height in a digital search tree and the longest phrase of the Lempel-Ziv scheme. In *Proc. 11th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 187–196.
- KOSARAJU, R. AND MANZINI, G. 1999. Compression of low-entropy strings with Lempel-Ziv algorithms. *SIAM Journal on Computing* 29, 3, 893–911.
- 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. 2008. Dynamic entropy-compressed sequences and full-text indexes. *ACM Transactions on Algorithms (TALG)*. 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.
- MORRISON, D. R. 1968. Patricia – practical algorithm to retrieve information coded in alphanumeric. *Journal of the ACM* 15, 4, 514–534.
- MUNRO, J. I. 1996. Tables. In *Proc. 16th Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS)*. LNCS 1180. 37–42.
- MUNRO, J. I., RAMAN, R., RAMAN, V., AND RAO, S. S. 2003. Succinct representations of permutations. In *Proc. 30th International Colloquium on Automata, Languages and Programming (ICALP)*. LNCS 2719. 345–356.
- MUNRO, J. I. AND RAMAN, V. 2001. Succinct representation of balanced parentheses and static trees. *SIAM Journal on Computing* 31, 3, 762–776.
- NAVARRO, G. 2004. Indexing text using the Ziv-Lempel trie. *Journal of Discrete Algorithms (JDA)* 2, 1, 87–114.
- NAVARRO, G. 2008. Implementing the LZ-index: Theory versus practice. *ACM Journal of Experimental Algorithmics (JEA)*. To appear.
- NAVARRO, G. AND MÄKINEN, V. 2007. Compressed full-text indexes. *ACM Computing Surveys* 39, 1, article 2.
- OKANOHARA, D. AND SADAKANE, K. 2007. Practical entropy-compressed rank/select dictionary. In *Proc. Workshop on Algorithm Engineering and Experiments (ALENEX)*. 60–70.
- RAMAN, R., RAMAN, V., AND RAO, S. S. 2002. Succinct indexable dictionaries with applications to encoding k -ary trees and multisets. In *Proc. 13th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 233–242.

- RUSO, L. AND OLIVEIRA, A. 2008. A compressed self-index using a Ziv-Lempel dictionary. *Information Retrieval* 11, 4, 359–388.
- SADAKANE, K. 2000. Compressed text databases with efficient query algorithms based on the Compressed Suffix Array. In *Proc. 11th Annual International Symposium on Algorithms and Computation (ISAAC)*. LNCS 1969. 410–421.
- SADAKANE, K. 2003. New Text Indexing Functionalities of the Compressed Suffix Arrays. *Journal of Algorithms* 48, 2, 294–313.
- ZIV, J. AND LEMPEL, A. 1978. Compression of individual sequences via variable-rate coding. *IEEE Trans. Inform. Theory* 24, 5, 530–536.