

# Fast in-memory XPath search using compressed indexes

Diego Arroyuelo<sup>1</sup>, Francisco Claude<sup>2</sup>, Sebastian Maneth<sup>3</sup>, Veli Mäkinen<sup>4</sup>,  
Gonzalo Navarro<sup>5</sup>, Kim Nguyễn<sup>6,\*</sup>, Jouni Sirén<sup>5</sup> and Niko Välimäki<sup>7</sup>

<sup>1</sup>*Departamento de Informática, Universidad Técnica Federico Santa María, Chile*

<sup>2</sup>*Escuela de Informática y Telecomunicaciones, Universidad Diego Portales, Chile*

<sup>3</sup>*School of Informatics, University of Edinburgh, UK*

<sup>4</sup>*HIIT and Department of Computer Science, University of Helsinki, Finland*

<sup>5</sup>*Department of Computer Science, University of Chile, Chile*

<sup>6</sup>*LRI, Université Paris-Sud, France*

<sup>7</sup>*Department of Medical Genetics, Faculty of Medicine, University of Helsinki, Finland*

## SUMMARY

Extensible Markup Language (XML) documents consist of text data plus structured data (markup). XPath allows to query both text and structure. Evaluating such hybrid queries is challenging. We present a system for in-memory evaluation of *XPath search queries*, that is, queries with text and structure predicates, yet without advanced features such as backward axes, arithmetics, and joins. We show that for this query fragment, which contains *Forward Core XPath*, our system, dubbed Succinct XML Self-Index ('SXSI'), outperforms existing systems by 1–3 orders of magnitude. SXSI is based on state-of-the-art indexes for text and structure data. It combines two novelties. On one hand, it represents the XML data in a compact indexed form, which allows it to handle larger collections in main memory while supporting powerful search and navigation operations over the text and the structure. On the other hand, it features an execution engine that uses tree automata and cleverly chooses evaluation orders that leverage the speeds of the respective indexes. SXSI is modular and allows seamless replacement of its indexes. This is demonstrated through experiments with (1) a text index specialized for search of bio sequences, and (2) a word-based text index specialized for natural language search. Copyright © 2013 John Wiley & Sons, Ltd.

Received 12 March 2013; Revised 13 August 2013; Accepted 19 August 2013

KEY WORDS: XML; succinct data structures; XPath; tree automata

## 1. INTRODUCTION

As increasing amounts of data are stored, transmitted, queried, and manipulated in XML, the popularity of XPath and XQuery as query languages for semi-structured data grows. Evaluating such XML queries efficiently is challenging, and has triggered much research. Today there is a wealth of public and commercial XPath/XQuery engines, apart from several theoretical proposals. In this paper, we focus on XPath, which is simpler and forms the basis of XQuery. XPath query engines can be roughly divided into two categories: *sequential* and *indexed*. In the former, which follows a *streaming* approach, no preprocessing of the XML data is performed. Each query sequentially reads the whole document, and the goal is to be as close as possible to making just one pass over the data, while using as little main memory as possible to hold intermediate results and data structures. Instead, the indexed approach preprocesses the XML document to build a data structure on it, so that queries can later be evaluated without traversing the whole document.

\*Correspondence to: Kim Nguyễn, LRI, Université Paris-Sud, France.

†E-mail: kn@lri.fr

A serious shortcoming of the indexed approach is that the index can use much more space than the original data, and thus may have to be manipulated on disk even on moderate-sized collections where the data itself would fit in main memory. Given the way disk performance favors sequential accesses, an index on disk may turn out to be slower than a streaming solution, even if the data are also stored on disk and if the index accesses only a minor part of the data. There are two approaches for dealing with this problem: (1) to load the index only partially (by using clever clustering techniques on disk), or (2) to use less powerful indexes that require less space. Examples of systems using these approaches are Qizx/DB [1], MonetDB/XQuery [2], and Tauro [3].

Most main memory XML query systems (such as Saxon [4], Galax [5], Qizx/Open [1], etc.) use machine pointers to represent XML data. We observe that on various well-established Document Object Model implementations, this representation blows up memory consumption to about 5–10 times the size of the original XML data. As a result, they can only handle in main memory XML collections that are much smaller than what streaming approaches could accommodate (as these require no extra data apart from the plain XML).

In this work, we aim at an index for XML that uses little extra space on top of that of the data (or actually less, as explained soon), yet without giving up on indexing power, but resorting instead to compact data structures. As a result, the index fits in main memory whenever the data does, thereby solving XPath queries without any need of resorting to disk. An in-memory index should outperform streaming approaches by far, even when the latter also operate in main memory. This is confirmed when comparing our indexed approach against two well-known streaming XPath engines (over data coming from a RAM-disk): GCX [6] and SPEX [7] are about 50 and 350 times, respectively, slower than our system.

An XML document can be regarded essentially as a *text collection* (that is, a set of strings) organized into a *tree structure*, so that the strings correspond to the text data and the tree structure corresponds to the nesting of tags. The problem of manipulating text collections and sequences within compressed space is now well understood [8–10], and also much work has been carried out on compact data structures for trees [11–16]. In this paper, we show how those types of compact data structures can be integrated into a compressed index representation for XML data, which is able to efficiently solve XPath queries.

A feature inherited from its components is that the compressed index *replaces* the XML collection, in the sense that the data (or any part of it) can be efficiently reproduced from the index (and thus, the data itself can be discarded). The result is called a *self-index*, as the data is inextricably tied to its index. A self-index might thus require less space than the original data, while representing it and at the same time offering indexed access to it.

Ours is not the first self-index for XML data. The so-called XBW index [17, 18] is a self-index offering some XPath search support, yet this is reduced to a very limited class of queries that are handled particularly well: ‘simple paths’, that is, queries of the form  $//t_1/t_2/ \dots /t_k$ , where each  $t_i$  is a tag name. For such queries, they can count the number of nodes satisfying the query in time  $O(k)$ , and they can report them in time  $O(\log^{1+\epsilon} n)$  per result, for any constant  $\epsilon > 0$ . For those specific queries, the XBW can be between one and two orders of magnitude faster than our system [19]. Likewise, there have been other attempts at using compact representations of the tree and the text. An earlier system in this line is BSBC [20]. They do not implement query evaluation on their compressed format, but their compression results are competitive. They show how some traversal operations used for XPath query evaluation can be sped up by using inverted indexes on the text content. Using inverted indexes, however, limits the applicability of the approach to compress natural language XML collections.

We aim at handling general sequences at the text nodes, and at a much more complete XPath coverage. Our system supports an extension of *Forward Core XPath* [21], that is, all forward navigational axes. The extension includes *text()* and the attribute axis, and the three text predicates = (equality), *contains*, and *starts-with*. We believe this to be a highly relevant and practical subset of XPath, and observe that a large fraction of real-world queries over datasets, such as *Medline* or *DBLP*, fall into this subset. Backward axes, arithmetics, and semi-joins are not yet handled. Our system, dubbed *Succinct XML Self-Index* (SXSI), is the first practical and public tool for compressed indexing of XML data. It takes a little space, solves a significant portion of XPath, and largely

outperforms the best public software supporting XPath we are aware of, namely MonetDB/XQuery [2] and Qizx/DB [1], in many cases by 1–3 orders of magnitude.

The main challenges in achieving our results have been to develop practical implementations of compact data structures (for texts, sequences, trees, and others) that are at a theoretical stage, to develop new compact schemes tailored to this particular problem and to develop query processing strategies tuned for the specific cost model that emerges from the use of these compact data structures. The limitations of our scheme are that it is in-memory, that it is static (i.e., the index must be rebuilt when the XML data changes), and that it does not handle the more sophisticated parts of XPath nor XQuery. The first limitation is a design decision; the last two are subject of future work.

This paper introduces the three main ingredients of SXSI: (i) the text index, (ii) the tree index, and (iii) the query evaluator. Although theoretical descriptions on the first two components can be found elsewhere in the literature, we mention here the main aspects of these components and focus on how they are integrated inside the practical SXSI system. For (iii), we have used tree automata, because they describe queries at a low level, thus allowing one to integrate the calls to our indexes. One important idea in integrating (ii) and (iii) is that of ‘jumping’ to a descendant or following a node (without traversing the intermediate nodes). This provides large speedups and allows an ‘automata-optimal’ evaluation [22]. Another important idea integrating (i) and (iii) is that of ‘true bottom-up runs’: typically, a tree can only be accessed through its root node. In SXSI, we are able to access leaves of the document tree. This allows us to start evaluation at text nodes, a technique that incurs large speedups for queries containing highly selective text predicates. Note that such queries are very common in practice.

The main part of the experimental section is about comparing SXSI against the state-of-the-art XPath engines MonetDB/XQuery and Qizx/DB. Certainly, this comparison must be considered with care: MonetDB and Qizx are full-blown XQuery engines with many features (such as multiuser and transaction support), whereas SXSI is a bare XPath engine. Hence, a comparison is hardly fair. However, our comparison shows the potential of succinct data structures and automata, as alternative core of XML databases. We use two batches of experiments: the ‘tree-oriented’ queries of the XPathMark benchmark [23] (over XMark data [24]) and our own ‘text-oriented’ queries (over Medline documents). Our results show that SXSI outperforms the other systems for virtually all tested queries, in many cases by 1–3 orders of magnitude, and moreover, that the running times of SXSI are more predictable and ‘robust’ than those of other systems. We also demonstrate SXSI’s ability to seamlessly integrate other indexes. We replace the text index by (1) an index tailored toward bio-sequence search and (2) a word-based index tailored toward natural language search.

We extend the work of [25] by (i) detailed explanations and examples, (ii) a description of important general optimizations used to efficiently run automata, (iii) a raw speed comparison of the indexes against naive solutions, (iv) an experimental analysis of the impact of different optimization techniques, and (v) an experimental comparison with alternative text indexes.

## 2. BASIC CONCEPTS AND MODEL

We regard an XML document as (i) an ordered set of strings and (ii) a labeled tree. The latter is the natural XML parse tree defined by the hierarchical tags, where the (normalized) tag name labels the corresponding node. We add an extra root node (labeled ‘&’) on top of the document’s root node; this node is needed for XPath semantics, but could also be used to hold additional information such as the document name. Each text node is represented as a leaf labeled #. Attributes are handled as follows in this model. Each node with attributes obtains an additional single child labeled @ (at the first child position), and for each attribute @attr=value of the node, a child labeled attr is added to its @-node, and a leaf child labeled % to the attr-node. The text content value is then associated to that leaf. Thus, there is exactly one string content associated to each tree leaf labeled # or %. We refer to those strings as *texts*. We do not store empty texts; for instance, the XML document `<a></a>` is stored as a single a-labeled leaf node (which is the unique child of the &-labeled root node).

Let us call  $T$  the concatenation of all the texts, each separated by a symbol ‘\$’ smaller than any other. Let  $n$  the total number of tree nodes,  $\Sigma$  the alphabet of the strings,  $t$  the total number of



In total, there are seven such whitespace texts, which have been omitted in our figure for reasons of readability.

Some notation and measures of compressibility follow, preceding a rough description of our space complexities. The *empirical  $k$ -th order entropy* [26] of a sequence  $S$  over alphabet  $\Sigma$ ,  $H_k(S) \leq \log |\Sigma|$ , is a lower bound to the output size per symbol of any  $k$ -th order compressor applied to  $S$ . The formula of the zero-order entropy is as follows:

$$H_0(S) = \sum_{c \in \Sigma} \frac{s_c}{s} \log \frac{s}{s_c},$$

where  $s_c$  is the number of occurrences of  $c$  in  $S$  and  $s = |S|$ . We assume  $\log = \log_2$  and  $0 \log 0 = 0$  henceforth. Let  $\Sigma^k$  denote the set of words over  $\Sigma$  of length  $k$ . Now let  $S_W$  be the set of characters preceding the occurrences of  $W \in \Sigma^k$  in  $S$ , then for  $k > 0$ ,

$$H_k(S) = \frac{1}{s} \sum_{W \in \Sigma^k} |S_W| H_0(S_W).$$

Note  $0 \leq H_k(S) \leq H_{k-1}(S) \leq \dots \leq H_0(S) \leq \log |\Sigma|$ .

We will build on self-indexes able of handling text collections  $T$  within  $|T|H_k(T) + o(|T| \log |\Sigma|)$  bits [9, 10, 27]. On the other hand, representing an unlabeled tree of  $n$  nodes requires  $2n - O(\log n)$  bits, and several representations using  $2n + o(n)$  bits support many tree query and navigation operations in constant time (e.g., [16]). The labels require in principle other  $n \log t$  bits. Sequences  $S$  of length  $n$  over an alphabet of size  $t$  can be stored within  $n \log t(1 + o(1))$  bits (and even  $nH_0(S) + o(n \log t)$ ), so that any element  $S[i]$  can be accessed, and they can also efficiently answer the following queries [8, 9, 28, 29]:

$rank_c(S, i)$  is the number of  $c$ 's in  $S[1, i]$ ; and  
 $select_c(S, j)$  is the position of the  $j$ -th  $c$  in  $S$ .

These are essential building blocks for more complex functionalities, as seen later.

The final space requirement of our index will include the following:

1.  $|T|H_k(T) + o(|T| \log |\Sigma|)$  bits for representing the text collection  $T$  in self-indexed form. This supports the string searches of XPath and can (slowly) reproduce any text.
2.  $d \log d + o(d \log d)$  bits for the mapping between the self-index and the text identifiers, for example, to determine to which text identifier a self-index position belongs, or restricting self-index searches to some texts.
3.  $2n + o(n)$  bits for representing the tree structure. This supports many navigational operations in constant time.
4.  $4n \log t + 2n + o(n)$  bits to represent the tags in a way that they support very fast XPath searches.
5.  $2n + o(n)$  bits for mapping between tree nodes and text identifiers.
6. Optionally,  $|T| \log |\Sigma|$  or  $|T|H_k(T) + o(|T| \log |\Sigma|)$  bits, plus  $O\left(d \log \frac{|T|}{d}\right)$ , to achieve faster text extraction than in 1).

As a practical yardstick, without the extra storage of texts (item 6), the memory consumption of our system is about the size of the original XML file (and, being a self-index, includes it!), and with the extra text store, the memory consumption is 1–2 times the size of the original XML file.

In Section 3, we describe our representation of the set of strings, including how to obtain text identifiers from text positions. This explains items 1, 2, and 6 previously. Section 4 describes our representation for the tree and the labels, and the way the correspondence between tree nodes and text identifiers works. This explains items 3, 4, and 5. Section 5 describes how we process XPath queries on top of these compact data structures. In Section 6, we give some implementation details and empirically compare our SXSI engine with the most relevant public engines we are aware of. We conclude in Section 7.



3. TEXT REPRESENTATION

Text data in SXSI is represented as a succinct full-text self-index [10], that is, generally known as the *FM-index* [30]. The index supports efficient pattern matching operations that can be easily extended to support different XPath predicates.

3.1. FM-index and backward searching

Given a string  $T$  of total length  $|T|$ , from an alphabet  $\Sigma$ , the *alphabet-friendly FM-index* [9] requires  $|T|H_k(T) + o(|T| \log |\Sigma|)$  bits of space for any  $k \leq \alpha \log_{|\Sigma|} n$  and any constant  $0 < \alpha < 1$ . The index supports counting the number of occurrences of a pattern  $P$  in  $O(|P| \log |\Sigma|)$  time. Locating the occurrences takes extra  $O(\log^{1+\epsilon} |T|)$  time per answer, for any constant  $\epsilon > 0$ .

The FM-index is based on the Burrows–Wheeler transform (BWT) of string  $T$  [31]. Assume  $T$  ends with the special end-marker  $\$$ . Let  $\mathcal{M}$  be a matrix whose rows are all the cyclic rotations of  $T$  in lexicographic order. The first column of  $\mathcal{M}$ , denoted  $F$ , contains all symbols of  $T$  in lexicographic order. The last column  $L$  of  $\mathcal{M}$  forms a permutation of  $T$ , which is the BWT string  $T^{bwt}$ . The matrix is only conceptual; the FM-index uses only on the  $T^{bwt}$  string. Figure 2 illustrates the matrix  $\mathcal{M}$  with its first and last rows ( $F$  and  $T^{bwt}$ ) in bold. Figure 1 (bottom right) shows how this fits in our overall scheme.

The resulting permutation from  $T$  to  $T^{bwt}$  is reversible. There exists a simple last-to-first mapping from symbols in  $T^{bwt}$  to  $F$  [30]: Let  $C[c]$  be the total number of symbols in  $T$  that are lexicographically less than  $c$ . Then, the *LF-mapping* is defined as

$$LF(i) = C[T^{bwt}[i]] + rank_{T^{bwt}[i]}(T^{bwt}, i).$$

Note that  $T^{bwt}[i]$  is the symbol preceding the  $i$ -th lexicographically smallest row of  $\mathcal{M}$ . Thus, if  $T^{bwt}[i] = T[j]$ , then  $T^{bwt}[LF(i)] = T[j - 1]$ . The symbols of  $T$  can therefore be read in reverse order by starting from the location  $i$  such that  $T^{bwt}[i] = \$$  and applying  $LF$  recursively:

$$T[|T|] = \$ = T^{bwt}[1] \quad T[|T| - 1] = T^{bwt}[LF(1)] \quad T[|T| - 2] = T^{bwt}[LF(LF(1))]$$

and so on until, after  $|T|$  steps, we obtain the first symbol  $T[1]$ . The values  $C[c]$  can be stored in a small array of  $|\Sigma| \log |T|$  bits. Function  $rank_c(T^{bwt}, i)$  can be computed in  $O(\log |\Sigma|)$  time with a data structure called *wavelet tree* that, when built on  $T^{bwt}$ , uses only  $|T|H_k(T) + o(|T| \log |\Sigma|)$  bits [9, 28, 32]. In practice, we opt for a Huffman-shaped wavelet tree using uncompressed bitmaps inside [33]. Despite this achieves space  $|T|(H_0(T) + 1)(1 + o(1))$ , it is much faster than the other implementations. In particular, operations cost  $O(H_0(T))$  time on average under some conditions, an improvement that applies to all the  $O(\log |\Sigma|)$  worst-case complexities that follow.

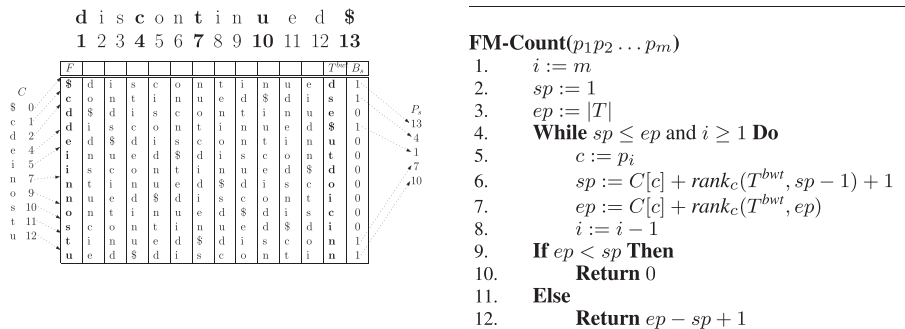


Figure 2. On the left, an example of the FM-index for text  $T = \text{“discontinued”}$  sampled each  $l = 3$  positions. On the right, counting algorithm on the FM-index.

Pattern matching is supported via *backward searching* on the BWT [30]. Given a pattern  $P[1, m]$ , the backward search starts with the range  $[sp, ep] = [1, |T|]$  of rows in  $\mathcal{M}$ . At each step  $i \in \{m, m-1, \dots, 1\}$  of the backward search, the range  $[sp, ep]$  is updated to match all rows of  $\mathcal{M}$  that have  $P[i, m]$  as a prefix. The new range  $[sp', ep']$  is given by  $sp' = C[P[i]] + rank_{P[i]}(T^{bwt}, sp - 1) + 1$  and  $ep' = C[P[i]] + rank_{P[i]}(T^{bwt}, ep)$ . Each step takes  $O(\log |\Sigma|)$  time using the wavelet tree, and finally,  $ep - sp + 1$  gives the number of times  $P$  occurs in  $T$ . Figure 2 gives the pseudocode.

To find out the location of each occurrence, the text is traversed backwards from each  $sp \leq i \leq ep$  (virtually, using  $LF$  on  $T^{bwt}$ ) until a *sampled* position is found. This is a sampling carried out at regular text positions, so that the corresponding positions in  $T^{bwt}$  are marked in a bitmap  $B_s[1, |T|]$ , and the text position corresponding to  $T^{bwt}[i]$ , if  $B_s[i] = 1$ , is stored in a samples array  $P_s[rank_1(B_s, i)]$ . If every  $l$ -th position of  $T$  is sampled, the extra space is  $O((n/l) \log n) + o(n)$  (including the compressed  $B_s$  [34]), and the locating takes  $O(l \log |\Sigma|)$  time per occurrence. Using  $l = \Theta(\log^{1+\epsilon} |T| / \log |\Sigma|)$  for any  $\epsilon > 0$  yields  $o(|T| \log |\Sigma|)$  extra space and locating time  $O(\log^{1+\epsilon} |T|)$ .

Figure 2 illustrates a sampling of  $T$  each  $l = 3$  symbols. Assume we look for  $P = \text{"n"}$ , then backward search finds  $[sp, ep] = [8, 9]$ . Now, to locate the occurrence at 8, we see that  $B_s[8] = 0$ ,  $B_s[LF(8)] = B_s[10] = 0$ , and finally,  $B_s[LF(10)] = B_s[2] = 1$ . This corresponds to position  $P_s[rank_1(B_s, 2)] = P_s[2] = 4$ . Because we applied  $LF$  twice, the answer is  $4 + 2 = 6$ . We have found the occurrence  $T[6..] = \text{"n. . ."}.$

### 3.2. Text collection and queries

The textual content of the XML data is stored as  $\$$ -terminated strings so that each text corresponds to one string. Let  $T$  be the concatenated sequence of the  $d$  texts. Array  $P_s$  is extended to record both the text identifier and the offset inside it. Because there are several  $\$$ 's in  $T$ , we fix a special ordering such that the end marker of the  $i$ -th text appears at  $F[i]$  in  $\mathcal{M}$  (see Figure 1, bottom right). This generates a valid  $T^{bwt}$  of all the texts and makes it easy to extract the  $i$ -th text starting from its  $\$$ -terminator.

Now  $T^{bwt}$  contains all end markers in some permuted order. This permutation is represented with a data structure  $Doc$ , that maps from positions of  $\$$ s in  $T^{bwt}$  to text identifiers. Let  $T^{bwt}[j]$  correspond to the first symbol of the text with identifier  $x$ , thus if  $i = LF(j)$ , it holds  $T^{bwt}[i] = \$$ . Then, we store  $Doc[rank_\$(T^{bwt}, i)] = x$ . Furthermore,  $Doc$  can be stored in a format that allows for range searching (as illustrated in Figure 1 (right)): Given a range  $[sp, ep]$  of  $T^{bwt}$  and a range of text identifiers  $[x, y]$ ,  $Doc$  can be used to output identifiers of all  $\$$ -terminators within the range  $[sp, ep] \times [x, y]$ , in  $O(\log d)$  time per answer [35]. In practice, because we only use the simpler functionality in the current system,  $Doc$  is implemented as a plain array using  $d \log d$  bits.

Note  $Doc$  allows us to never switch from one text to another while looking for the preceding sampled value: If we reach a  $\$$  before finding any  $B_s[i] = 1$ , array  $Doc$  can be used to determine that we are at the first position of some text with identifier  $x$ .

The basic pattern matching feature of the FM-index can be extended to support XPath functions such as *starts-with*, *ends-with*, *contains*, and operators  $=, \leq, <, >, \geq$  for lexicographic ordering. Given a pattern and a range of text identifiers to be searched, these functions return all text identifiers that match the query within the range. In addition, existential (is there a match in the range?) and counting (how many matches in the range?) queries are supported. Time complexities are  $O(|P| \log |\Sigma|)$  for the search phase, plus an extra for reporting. Although we describe the operators in their general form, which needs the range reporting functionality from  $Doc$ , our current prototype implements only the simple case  $[x, y] = [1, d]$ , where  $Doc$  can be an array.

**starts-with( $P, [x, y]$ ):** The goal is to find texts in range  $[x, y]$  prefixed by the given pattern  $P$ . After the normal backward search, the range  $[sp, ep]$  in  $T^{bwt}$  contains the end markers of all the texts prefixed by  $P$ . Now  $[sp, ep] \times [x, y]$  can be mapped to  $Doc$ , and existential and counting queries can be answered in  $O(\log |\Sigma| + \log d)$  time. Matching text identifiers can be reported in  $O(\log d)$  time per identifier. If  $[x, y] = [1, d]$  and  $Doc$  is an array, the counting time after the backward search is  $O(\log |\Sigma|)$ , and each text identifier can be reported in constant time.

**ends-with**( $P, [x, y]$ ): Backward searching is localized to texts in  $[x, y]$  by choosing  $[sp, ep] = [x, y]$  as the starting interval, because we have forced the ordering of  $F[1, d]$ , so that  $F[z] = \$$  is the terminator of text with identifier  $z$ . After the backward search, the resulting range  $[sp, ep]$  contains all possible matches, thus existential and counting queries are answered in constant time after the search. To find out text identifiers for each occurrence, the text must be traversed backwards to find a sampled position (or a  $\$$ ). The cost is  $O(l \log |\Sigma| + \log d)$  per answer, where  $l$  is the sampling step. If  $[x, y] = [1, d]$  and  $Doc$  is an array, the cost is just  $O(l \log |\Sigma|)$ .

**operator** = ( $P, [x, y]$ ): Whole texts that are equal to  $P$ , and with identifiers in the range  $[x, y]$ , can be found as follows. Start with a backward search as in *ends-with*, and then map to the  $\$$ -terminators as in *starts-with*. The time complexities are same as in *starts-with*.

**contains**( $P, [x, y]$ ): To find texts that contain  $P$ , we start with the normal backward search and finish such as in *ends-with*. In this case, there might be several occurrences inside one text, which have to be filtered. Thus, the time complexity is proportional to the total number of occurrences,  $O(l \log |\Sigma|)$  for each. Existential and counting queries are as slow as reporting queries. The basic  $O(|P| \log |\Sigma|)$ -time counting of all the occurrences of  $P$  can still be useful for query optimization.

**operators**  $\leq, <, >, \geq$ : Operator  $\leq$  matches texts that are lexicographically smaller than or equal to the given pattern. It can be solved such as the *starts-with* query, but updating only the  $ep$  of each backward search step, whereas  $sp = 1$  stays constant. Whereas  $[sp, ep]$  delimits the rows of  $\mathcal{M}$  that start with  $P[i, m]$ ,  $[1, ep]$  delimits the rows that start with a prefix lexicographically smaller than or equal to  $P[i, m]$ . If at some point there are no occurrences of  $P[i] = c$  within the prefix  $T^{bwt}[1, ep]$ , this means that  $P[i, m]$  does not appear in  $T$ . To continue the search, we replace  $ep = C[c]$  and continue for  $P[1, i - 1]$ . Other operators can be supported analogously, and costs are as for *starts-with*.

The new XPath extension, *XPath Full Text 1.0* [36], suggests a wider functionality for text searching. Implementation of these extensions requires regular expression and approximate searching functionalities, which can be supported within our index using the general *backtracking framework* [37]: The idea is to alter the backward search to branch recursively to different ranges  $[sp', ep']$  representing the suffixes of the text prefixes (i.e., substrings). This is performed by computing  $sp'_c = C[c] + rank_c(T^{bwt}, sp - 1) + 1$  and  $ep'_c = C[c] + rank_c(T^{bwt}, ep)$  for all  $c \in \Sigma$  at each step and recursing on each  $[sp'_c, ep'_c]$ . Then, the pattern (or regular expression) can be compared with all substrings of the texts, allowing us to search for approximate occurrences [37]. The running time becomes exponential in the number of errors allowed, but different branch-and-bound techniques can be used to obtain practical running times [38, 39]. We omit further details, as these extensions are out of the scope of this paper.

### 3.3. Construction and text extraction

The FM-index can be built by adapting any BWT construction algorithm. Linear time algorithms exist for the task, but their practical bottleneck is the peak memory consumption. Although there exist general time-efficient and space-efficient construction algorithms, it turned out that our special case of text collection admits a tailored incremental BWT construction algorithm [40] (see the references and experimental comparison therein for previous work on BWT construction): The text collection is split into several smaller collections, and a temporary index is built for each of them separately. The temporary indexes are then merged and finally, converted into a static FM-index. The BWT allows extracting the  $i$ -th text by successively applying  $LF$  from  $T^{bwt}[i]$ , at  $O(\log |\Sigma|)$  cost per extracted symbol.

### 3.4. Faster text extraction using more space

To enable faster text extraction, we allow storing the texts in plain format in  $n \log |\Sigma|$  bits, or in an enhanced LZ78-compressed format (derived from the LZ-index [41]) using  $|T|H_k(T) +$



$o(|T| \log |\Sigma|)$  bits. These secondary text representations are coupled with a delta-encoded bit vector storing starting positions of each text in  $T$ . This bitmap requires  $O\left(d \log \frac{|T|}{d}\right)$  more bits.

In fact, keeping next to the FM-index, an additional copy of all texts in plain format has more advantages. As mentioned before, the time complexity of contains-queries is proportional to the total number of occurrences. This implies that for large occurrence numbers, it becomes faster to search over the plain texts than over the FM-index. The precise cut-off point depends on the sampling factor  $l$ , see Section 6.3 for more details. Because a global count over the FM-index is fast ( $O(|P| \log |\Sigma|)$  time), we use it to decide whether to search over the plain text or over the FM-index.

In practice, we opt for a plain text representation, which is much faster for extraction than an LZ-index at the price of not much more space.

## 4. TREE REPRESENTATION

### 4.1. Data representation

The tree structure of an XML collection is represented by the following compact data structures, which provide navigation and indexed access to it. See also the bottom left of Figure 1.

**4.1.1. *Par*.** This is the *balanced parentheses* representation of the tree structure (e.g., [42]). It is obtained by traversing the tree in *depth-first-search* order (or preorder), writing a ‘ ( ’ whenever we arrive at a node and a ‘ ) ’ when we leave it (thus, it follows the sequences of events generated by an XML SAX parser). In this way, every node is represented by a pair of matching opening and closing parentheses. A tree node is identified by the position of its opening parenthesis in *Par* (that is, a node is just an integer index within *Par*). In particular, we use the balanced parentheses implementation of [16], which supports a very complete set of operations, including finding the  $i$ -th child of a node, in constant time; for more information concerning implementation details and performance, see [43]. Overall *Par* uses  $2n + o(n)$  bits. This includes the space needed for constant-time binary *rank* on *Par*, which is very fast in practice.

**4.1.2. *Tag*.** This is the sequence of the tag identifiers of each tree node, including an opening and a closing version of each tag, to mark the beginning and ending point of each node. These tags are numbers in  $[1, 2t]$  and are aligned with *Par* so that the tag of node  $i$  is simply  $Tag[i]$ .

We also need *rank* and *select* queries on *Tag*. They allow us to carry out special operations such as ‘TaggedDesc’, which ‘jumps’ to the first descendant of the given node having a given label (Section 4.2.2). Several sequence representations supporting access, and these operations are known [8, 28, 33]. Given that *Tag* is not too critical in the overall space, but it is in time, we opt for a practical representation that favors speed over space. First, we store the tags in an array using  $\lceil \log 2t \rceil$  bits per field, which gives constant time access to  $Tag[i]$ . The rank and select queries over the sequence of tags are answered by a second structure. Consider the binary matrix  $R[1..2t][1..2n]$  such that  $R[i, j] = 1$  if  $Tag[j] = i$ . We represent each row of the matrix using Okanojara and Sadakane’s structure `sarray` [44]. Its space requirement for each row  $i$  is  $n_i \log \frac{2n}{n_i} + n_i(2 + o(1))$  bits, where  $n_i$  is the number of times symbol  $i$  appears in *Tag*. The total space of both structures adds up to  $2n \log(2t) + 2nH_0(Tag) + n(2 + o(1)) \leq 4n \log t + 2n + o(n)$  bits. Thus, we support access and *select* in  $O(1)$  time, and *rank* in  $O(\log n)$  time.<sup>‡</sup>

### 4.2. Tree navigation

We define the following operations over the tree structure, which are useful to support XPath queries over the tree. Most of these operations are supported in constant time, except when a *rank* over *Tag* is involved. In what follows, we assume that all the operations take an implicit argument *Tree* (we do not write it explicitly to improve the readability). Nodes of *Tree* (that is, positions in *Par*) are

<sup>‡</sup>They report higher complexities, but these are easily improved by using a representation for dense arrays that supports *select* in constant time.

ranged over by  $x$ ,  $y$ , and so on. We assume the existence of a *dummy* node `Nil` guaranteed to be distinct from any node of *Tree* (for instance  $-1$ ).

**4.2.1. Basic tree operations.** These are directly inherited from Sadakane's implementation [16]. We mention only the most important ones for this paper.

- `Close(x)`: The closing parenthesis matching `Par[x]`. If  $x$  is a small subtree, this takes a few local accesses to `Par`, otherwise a few nonlocal table accesses.
- `Preorder(x) = ranki(Par, i)`: Preorder number of  $x$ .
- `SubtreeSize(x) = (Close(x) - x + 1)/2`: Number of nodes in the subtree rooted at  $x$ .
- `IsAncestor(x, y) = x ≤ y ≤ Close(x)`: Whether  $x$  is an ancestor of  $y$ .
- `IsLeaf(x) = (Par[x + 1] = '')`: Whether node  $x$  is a leaf in the tree.
- `FirstChild(x) = x + 1`: First child of  $x$ , if any (i.e., if `Par[x + 1] = '('`), `Nil` otherwise (i.e. if  $x$  denotes a leaf).
- `NextSibling(x) = Close(x) + 1`: Next sibling of  $x$ , if any (i.e., if `Par[Close(x) + 1] = '('`), `Nil` otherwise.
- `Parent(x)`: Parent of  $x$  found as the closest parentheses pair enclosing  $x$ . Somewhat costlier than `Close(x)` in practice, because the answer is less likely to be near  $x$  in `Par`. Return `Nil` for the root node.

**4.2.2. Connecting to tags.** The following operations are essential for our fast XPath evaluation. Let *tag* be a tag identifier.

- `SubtreeTags(x, tag)`: Returns the number of occurrences of *tag* within the subtree rooted at node  $x$ . This is  $rank_{tag}(Tag, Close(x)) - rank_{tag}(Tag, x - 1)$ .
- `Tag(x)`: Gives the tag identifier of node  $x$ . In our representation, this is just `Tag[x]`. Returns `Nil` if there are no such nodes.
- `TaggedDesc(x, tag)`: The first node (in preorder) labeled *tag* strictly within the subtree rooted at  $x$ . It is obtained as  $select_{tag}(Tag, rank_{tag}(Tag, x) + 1)$  if it is  $≤ Close(x)$ ; otherwise, there is no such node, and the function returns `Nil`.
- `TaggedPrec(x, tag)`: The last node labeled *tag* with preorder smaller than that of node  $x$ , and not an ancestor of  $x$ . Let  $r = rank_{tag}(Tag, x - 1)$ . If  $select_{tag}(Tag, r)$  is not an ancestor of node  $x$ , we return it. Otherwise, we set  $r = r - 1$  and iterate. Returns `Nil` when  $r = 0$ .
- `TaggedFoll(x, tag)`: The first node labeled *tag* with preorder larger than that of  $x$ , and not in the subtree of  $x$ . This is  $select_{tag}(Tag, rank_{tag}(Tag, Close(x)) + 1)$ . Return `Nil` if there is no such node.

**4.2.3. Connecting the text and the tree.** Conversion among text numbers, tree nodes, and global identifiers is easily carried out by using `Par` and a bitmap  $B$  of  $2n$  bits that marks the opening parentheses of tree leaves containing text, plus  $o(n)$  extra bits to support rank/select queries. The bitmap  $B$  uses an implementation by Raman *et al.* [34], which is described in [33], and it enables the computation of the following operations:

- `LeafNumber(x)`: Gives the number of leaves up to  $x$  in `Par`. This is  $rank_1(B, x)$ .
- `TextIds(x)`: Gives the range of text identifiers that descend from node  $x$ . This is simply  $[LeafNumber(x - 1) + 1, LeafNumber(Close(x))]$ .
- `XMLIdText(d)`: Gives the global tree preorder identifier for the text with identifier  $d$ . This is  $Preorder(select_1(B, d))$ .
- `XMLIdNode(x)`: Gives the global identifier for a tree node  $x$ . This is just  $Preorder(x)$ .

### 4.3. Displaying contents

Given a node  $x$ , we want to recreate its XML serialization, that is, return (a portion of) the original XML string. We traverse the structure starting from `Par[x]`, retrieving the tag names and the text contents, from the text identifiers. The time is  $O(\log \sigma)$  per text symbol (or  $O(1)$  if we use the redundant text storage described in Section 3) and  $O(1)$  per tag.

- `GetText( $d$ )`: generates the text with identifier  $d$ .
- `GetSubtree( $x$ )`: generates the subtree at node  $x$ .

## 5. XPATH QUERIES

In this section, we define the XPath fragment ‘Core+’, show its translation into automata, and discuss efficient execution of these automata. We do not formally define the semantics of XPath and assume the reader to be familiar with the basics of XPath, see, for example, [21].

### 5.1. The XPath fragment Core+

Our goal is to support a practical subset of XPath, while being able to guarantee efficient evaluation based on the data structures described in the previous sections. As a first shot, we target the forward fragment of ‘Core XPath’ [21]. Here is an Extended Backus–Naur Form for Core XPath.

```

Core      ::= LocationPath | / LocationPath
LocationPath ::= LocationStep (/ LocationStep)*
LocationStep ::= Axis :: NodeTest | Axis :: NodeTest [ Pred ]
Pred       ::= Pred and Pred | Pred or Pred | not ( Pred ) | Core | ( Pred )

```

We focus our presentation on the descendant and child axes, but `self`, `attribute` and `following-sibling` are also supported in our implementation. A node test is either the wild-card (`*`), a tag name, or a node type test, that is, one of ‘`text()`’ or ‘`node()`’.

Our fragment, called *Core+*, supports forward Core XPath and additionally all text predicates of XPath 1.0, that is, the `=` (equality), `contains`, `starts-with`, and `ends-with` predicates. These predicates appear inside filters (square brackets) and are generated by the nonterminal ‘Pred’. Equality tests if a (constant) string is equal to a string selected by a *Core+* expression; `contains` tests if the string is contained in the expression; and `starts-with` tests if the string is a prefix of the expression. Thus, our *Core+* fragment specializes the previous Extended Backus–Naur Form with the following rules:

```

Axis      ::= descendant | child | self | attribute | following-sibling
NodeTest  ::= * | TagName | text() | node()
Pred      ::= Core+ = String | contains(Core+, String) | starts-with(Core+, String)
           | ends-with(Core+, String)

```

An XPath query selects nodes of an XML document. The last axis in a query determines the selected nodes. For instance, the query `/descendant::listitem/child::keyword` selects all `keyword`-children of all `listitem`-nodes in the document. The query `T = /descendant::text()` selects all text nodes of the document, and the query `A = /descendant::*/@attribute::*` selects all attribute nodes of the document. In terms of our example document in Figure 1, the query `T` selects the nodes 7, 9, 11, and 17 of our model. XPath processors return the XML content (subtrees) of the selected nodes. Thus, for `T` the strings ‘Soon discontinued’, ‘blue’, ‘40’, and ‘30’ are returned (plus the whitespace text nodes, compared with the discussion at the beginning of Section 2).

The string held by every text and attribute node is called its *value*. If a subquery only selects text and attribute nodes, then we call it a *value expression*. The ‘Core+’ nonterminals mentioned on the right-hand side of our last ‘Pred’ rule must correspond to value expressions (we have avoided to complicate the grammar to enforce that). As an example, the query

```
/descendant::text() [contains(self::node(), 'and')]
```

is in *Core+*; note that ‘self’ here is a value expression, because `/descendant::text()` is.

### 5.2. From XPath to automata

We use ‘marking tree automata’ as our execution model for *Core+* queries. Before we formally define these automata in the next section, we explain here by means of an example how a *Core+*

1	$q_0, \{\&\}$	$\rightarrow$	$\downarrow_1 q_1$
2	$q_1, \{\text{listitem}\}$	$\rightarrow$	$\downarrow_1 q_1 \wedge \downarrow_1 q_2 \wedge \downarrow_2 q_1$
3	$q_1, \mathcal{L}$	$\rightarrow$	$\downarrow_1 q_1 \wedge \downarrow_2 q_1$
4	$q_2, \{\text{keyword}\}$	$\rightarrow$	$\text{mark} \wedge \downarrow_1 q_2 \wedge \downarrow_1 q_3 \wedge \downarrow_2 q_2$
5	$q_2, \mathcal{L}$	$\rightarrow$	$\downarrow_1 q_2 \wedge \downarrow_2 q_2$
6	$q_3, \{\text{emph}\}$	$\rightarrow$	$\top$
7	$q_3, \mathcal{L}$	$\rightarrow$	$\downarrow_2 q_3$

Figure 3. The tree automaton for the query `/descendant::listitem/descendant::keyword[child::emph]`.

query is translated into an automaton. The translation of an XPath query to an automaton is a simple syntax-directed translation that can be carried out in one pass through the parse tree of the query. Roughly speaking, the resulting automaton is ‘isomorphic’ to the original query. Consider the query

`/descendant::listitem/descendant::keyword[child::emph]` .

This query selects all `keyword`-nodes that are descendants of `listitem`-nodes and that have a child node tagged `emph`. Formally speaking, the query starts with the expression ‘/’, which selects the root node of the document (in our model, the `&`-node), then applies the descendant axis to that, and so on. The automaton has four states,  $q_0$ ,  $q_1$ ,  $q_2$ , and  $q_3$ , which correspond to the four steps in the query (namely, ‘/’, two descendants, and the child one). Its transitions are given in Figure 3. The automaton has start state  $q_0$  and end states  $q_1$  and  $q_2$ . At first, we can consider that this automaton is a classical *nondeterministic* alternating tree automaton (e.g., as the ones documented in [45] and [46]). The transitions have the form: ‘state, set of labels  $\rightarrow$  formula’. In the transitions previously,  $\mathcal{L}$  denotes the whole alphabet of the automaton and in the formulas,  $\top$  denotes the Boolean ‘true’, and  $\downarrow_1 q$  (resp.  $\downarrow_2 q$ ) is true if there exists an accepting run from state  $q$  on the first child (resp. next sibling) of the current node. The notion of (accepting) run for a given input tree is again the usual one:

- the root must be in a start state (here  $q_0$ );
- the leaves must be in an end state (here  $q_1$  or  $q_2$ );
- a node  $x$  of the tree is in state  $q$  if there exists a transition  $q, L \rightarrow \phi$ , if the label of  $x$ , is in  $L$  and if the formula  $\phi$  holds (which possibly requires some conditions on the left or right child of  $x$ ).

The novelty here is the presence of the `mark` predicate (Transition 4) whose intuitive meaning is to remember the nodes in which the transition containing `mark` (here Transition 4) was valid.

Let us now describe informally the correspondence between the XPath formula and the automaton. The latter starts at the root of the input tree in state  $q_0$ . Here, only Transition 1 can be satisfied. That is the case if the first child of the root is in state  $q_1$ . That state has two corresponding transitions (recall that our automaton is *nondeterministic*). Transition 2 requires that (i) the current label is `listitem` and  $q_2$  hold for the first child and  $q_1$  hold for both the first child and next sibling. Transition 3 has no requirement on the label (it can be anything in  $\mathcal{L}$ ) but  $q_1$  must hold for both the first child and next sibling of the current node. This self reference to  $q_1$  in both directions simply encodes the recursion performed by the `descendant` axis in the query. Likewise for  $q_2$  (and in general for any step of the query but the initial ‘/’), there are two transitions. Transition 5 handles the recursion in case the label of the current node is not `keyword` or if the current node is a `keyword`, which has no `emph` in child position. Transition 4 requires that the current node has label `keyword` that  $q_2$  holds for both first child and next sibling (recursion) and that  $q_3$  holds for the first child. If that is the case, the node is *marked*. Lastly  $q_3$  encodes the filter `[child::emph]`. If the current label is `emph`, then the transition is satisfied (there is no need to iterate in that case, because only one `emph`-node is sufficient for the filter to be true). Or the current node is not `emph`, and the automaton looks for an `emph` on the next sibling (recursion on  $\downarrow_2$  only, which encodes a `child` axis).

There are of course several accepting runs for a given input tree (because the automaton is non-deterministic), but if a node is *marked* during a run, then it is a *keyword* node, which has a *listitem*-node above it and an *emph* node amongst its children.

It is well-known that all the runs of a *nondeterministic* automaton can be simulated in a one pass traversal of the input (this holds for word automata as well as tree automata). Essentially, one maintains for each node a *set of states* (all those in which the nondeterministic automaton can be) instead of a single state. We show, after introducing formally our automata model, how we can compute the set of nodes marked during any nondeterministic run for a given automata, first using only FirstChild and NextSibling move, in one traversal of the tree. Then, we show several optimizations techniques that allow us to leverage the speed of the low-level tree and text indices and to compute efficiently the set of marked nodes.

### 5.3. Tree automata representation

Tree automata are a well-known and popular tool for reasoning about XML, see, for example, [47–50]. Only seldom have they been used as a tool for query evaluation. In [51], automata are used to evaluate, on an XML stream, many (very simple) XPath queries in parallel. It is well known that Core XPath can be evaluated using tree automata; see, for example, [52] and [53]. Here, we use alternating tree automata (as in [45] and [46]). Such automata work with Boolean formulas over states, which must become satisfied for a transition to be triggered. This allows a much more compact representation of queries through automata than ordinary tree automata (without formulas). Our tree automata are defined over a binary tree view of the XML tree where the left child is the first child of the XML node, and the right child is the next sibling of the XML node.

#### Definition 5.1

A nondeterministic marking automaton  $\mathcal{A}$  is a tuple  $(\mathcal{L}, \mathcal{Q}, \mathcal{T}, \mathcal{B}, \delta)$ , where

- $\mathcal{L}$  is a countable (possibly infinite) set of tree labels;
- $\mathcal{Q}$  is a finite set of states;
- $\mathcal{T} \subseteq \mathcal{Q}$  is a set of top states (that is, states that must be satisfied at the root node);
- $\mathcal{B} \subseteq \mathcal{Q}$  is a set of bottom states (that is, states that must be satisfied at the leaves);
- $\delta : \mathcal{Q} \times 2_{\mathcal{L}}^{\mathcal{L}} \cup 2_{\text{cof}}^{\mathcal{L}} \rightarrow F$  is a transition function, where  $F$  is the set of Boolean formulas,  $2_{\mathcal{L}}^{\mathcal{L}}$  is the set of finite subsets of  $\mathcal{L}$ , and  $2_{\text{cof}}^{\mathcal{L}}$  is the set of cofinite subsets of  $\mathcal{L}$ . A *Boolean formula*  $\phi$  is produced by the grammar,

$$\begin{aligned} \phi & ::= \top \mid \perp \mid \text{mark} \mid \phi \vee \phi \mid \phi \wedge \phi \mid \neg \phi \mid a \mid p && \text{(formula)} \\ a & ::= \downarrow_1 q \mid \downarrow_2 q && \text{(atom),} \end{aligned}$$

where  $p \in P$  is a built-in *predicate* and  $q$  is a state.

Before explaining in detail the use of formulas, we motivate our use of finite or cofinite sets as guards for transitions. Although traditionally automata transitions are guarded by a state and a single label, this would make the encoding of XPath into automata very tedious and needlessly complicate the algorithms. Indeed, one of the features of XPath is a wildcard element test, namely ‘\*’. One solution could be to suppose that for a given automaton the set of labels of the input document is known in advance and that this set is used as alphabet for the automaton. Unfortunately, this does not accurately reflect the semantics of XPath in which a query can be defined independently of any document and can even be executed on *any* document (it might not yield any result, but its application is valid). Another solution (as in [51]) is to equip automata with a special ‘default’ transition, labeled for instance ‘\_’, which is taken if in the current state no other transition can be evaluated. This has two drawbacks. First, it is only well-defined for deterministic tree automata (our encoding makes heavy use of non-determinism). Second, the evaluation function is polluted by the special cases which handle this default transition. Our solution is more blunt. We guard transitions by finite or co-finite sets of labels, and a transition is taken if the label of the current node is a member of that set. For instance, the ‘\*’ XPath test is encoded



$$\begin{array}{c}
\frac{}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \top = (\top, \emptyset)} \textbf{(true)} \\
\frac{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi_1 = (b_1, R_1) \quad \mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi_2 = (b_2, R_2)}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi_1 \vee \phi_2 = (b_1, R_1) \otimes (b_2, R_2)} \textbf{(or)} \\
\frac{q \in \text{dom}(\mathcal{R}_i)}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \downarrow_i q = (\top, \mathcal{R}(q))} \text{ for } i \in \{1, 2\} \textbf{(left, right)} \\
\frac{\text{EvalPred}(p, t') = b, R}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} p = (b, R)} \textbf{(pred)} \\
\frac{}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \neg \phi = (\bar{b}, \emptyset)} \textbf{(not)} \\
\frac{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi_1 = (b_1, R_1) \quad \mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi_2 = (b_2, R_2)}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi_1 \wedge \phi_2 = (b_1, R_1) \otimes (b_2, R_2)} \textbf{(and)} \\
\frac{}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \text{mark} = (\top, \{t'\})} \textbf{(mark)} \\
\frac{\text{when no other rule applies}}{\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi = (\perp, \emptyset)}
\end{array}$$

where:

$$\begin{array}{l}
\bar{\top} = \perp \quad \text{and} \quad \bar{\perp} = \top \\
(b_1, R_1) \otimes (b_2, R_2) = \begin{cases} \top, R_1 & \text{if } b_1 = \top, b_2 = \perp \\ \top, R_2 & \text{if } b_2 = \top, b_1 = \perp \\ \top, R_1 \cup R_2 & \text{if } b_1 = \top, b_2 = \top \\ \perp, \emptyset & \text{otherwise} \end{cases} \\
(b_1, R_1) \oplus (b_2, R_2) = \begin{cases} \top, R_1 \cup R_2 & \text{if } b_1 = \top, b_2 = \top \\ \perp, \emptyset & \text{otherwise} \end{cases}
\end{array}$$

Figure 4. Inference rules defining the evaluation of a formula.

as a transition guarded by the set  $\mathcal{L} - \{ @, \# \}$ , where ‘@’ and ‘#’ represent labels of subtrees containing attribute nodes and text nodes in our encoding. This allows us to give a very straightforward evaluation function for tree automata, which relies on the evaluation of Boolean formulas, presented next.

#### Definition 5.2 (Evaluation of a formula)

Given an automaton  $\mathcal{A}$  and an input tree  $t$ , the evaluation of a formula is given by the judgment  $\mathcal{R}_1, \mathcal{R}_2, t' \vdash_{\mathcal{A}} \phi = (b, R)$  where  $\mathcal{R}_1$  and  $\mathcal{R}_2$  are mappings from states to sets of nodes of  $t$ ,  $t'$  is a node of  $t$ ,  $\phi$  is a formula,  $b \in \{ \top, \perp \}$ , and  $R$  is a set of nodes of  $t$ . We define the semantics of this judgment by the means of the inference rules given in Figure 4.

These rules are straightforward and combine the rules for a classical alternating automaton, with the rules of a marking automaton. Rules **(or)** and **(and)** implement the Boolean connective of the formula and collect the marking found in their true subformulas. Rules **(left)** and **(right)** (written as a rule scheme for conciseness) evaluate to true if the state  $q$  is in the corresponding set. Intuitively,  $\mathcal{R}_1$  (resp.  $\mathcal{R}_2$ ) is the set of states recognizing the left (resp. right) subtree of the input tree. Rule **(pred)** assumes the existence of an evaluation function for built-in predicates. Among the latter, we assume the existence of a special predicate `mark`, which evaluates to  $\top$  and returns the singleton set containing the current node.

We now give the semantics of an automaton by means of the *run function* **TopDownRun** (Figure 5). This algorithm is based on the textbook algorithm for recursive bottom-up evaluation of tree automata (e.g., [46]). The algorithm performs a recursive first child/next sibling traversal of the tree until a leaf is reached (base case for the recursion). When returning from the recursive evaluation on the left and right subtrees (Lines 6 and 7, Figure 5), the function evaluates the set of transitions for the current node, based on the set of states recognizing the left and right subtree. However, instead of blindly doing a recursive descent from the root to the leaves and evaluating when returning from the recursive calls, the transitions are restricted by the set of states  $Q_{\text{id}}$  (Line 4). This technique is dubbed ‘bottom-up evaluation with top-down preprocessing’ in [46]. We therefore named the run function **TopDownRun** to differentiate it from a real bottom-up run (starting from the leaves of the tree) that we present in Section 5.4.2. The novelty is our use of maps from states to nodes instead of only sets of states. The resulting map associate any state  $q$ , which has an accepting run from  $x$  with the set of nodes that were marked during that run.

**Input** An automaton  $\mathcal{A} = (\mathcal{L}, \mathcal{Q}, \mathcal{T}, \mathcal{B}, \delta)$ , a node  $x$  (of the implicit *Tree*), and a set of states  $Q_{\text{id}} \subseteq \mathcal{Q}$ .

**Output** A mapping  $\mathcal{R}$  from states to sets of nodes, such that  $\mathcal{R}(q)$  is the set of nodes in the subtree rooted at  $x$  that were marked during an accepting run from  $x$  starting in state  $q$ .

**TopDownRun**( $\mathcal{A}, x, Q_{\text{id}}$ )

1. **If**  $x$  is Nil **Then**
2.     **Return**  $\{q \rightarrow \emptyset \mid q \in \mathcal{B} \cap Q_{\text{id}}\}$
3. **Else**
4.      $\text{trans} := \{(q, \ell \rightarrow \phi) \mid q \in Q_{\text{id}}, \text{Tag}(x) \in \ell, (q, \ell \rightarrow \phi) \in \delta\}$
5.      $Q_{\text{id}}^i := \{q \mid \downarrow_i q \text{ occurs in } \phi, (q', \ell \rightarrow \phi) \in \text{trans}\}$  for  $i \in \{1, 2\}$
6.      $\mathcal{R}_1 := \text{TopDownRun}(\mathcal{A}, \text{FirstChild}(x), Q_{\text{id}}^1)$
7.      $\mathcal{R}_2 := \text{TopDownRun}(\mathcal{A}, \text{NextSibling}(x), Q_{\text{id}}^2)$
8.     **Return**  $\{q \mapsto R \mid \mathcal{R}_1, \mathcal{R}_2, x \vdash_{\mathcal{A}} \phi = (\top, R), (q, \ell \rightarrow \phi) \in \text{trans}\}$

Figure 5. Evaluation function for tree automata.

#### 5.4. Leveraging the speed of the low-level interface

We have seen how to evaluate an XPath query by compiling it into a tree automaton and running the latter on the input document. We present now several techniques that make use of the tree and text index presented in Sections 4 and 3. These are the techniques that make our SXSI prototype competitive in speed with state-of-the-art XML databases.

**5.4.1. Jumping to relevant nodes.** Conventionally, the run of a tree automaton visits every node of the input tree. This is, for instance, the behavior of the tree automata presented in [52], which perform two scans of the whole XML document (the latter being stored on disk in a particular format). However, for typical queries, most of the nodes are ‘useless’ in the sense that the automaton only loops through them staying in the same set of states. In other words, the automaton *ignores* most of the nodes. To restrict the run to interesting nodes, we use the notion of *relevant nodes* introduced in [22]. Although the full characterization is out of the scope of this paper, we give a flavor of relevant nodes, using an example. Consider again the query

/descendant::listitem/descendant::keyword[child::emph]

whose corresponding automaton is given in Figure 3. If we consider the starting transition, (Line 1) in Figure 3, we can see that at the root node, (labeled &) the automaton will first go on the first child, in state  $\{q_1\}$ . Then, it will loop, going down on the first child and next sibling of each node until it reaches a `listitem` element, on which it changes state and goes in  $\{q_1, q_2\}$ . Here we see, first, that there are no relevant nodes between the root and the first `listitem`. Indeed, the tree could be of any shape, and labeled with any tag (besides `listitem`), and the result of the query would be the same. Second, we see this `listitem` node is reached through the following sequence of moves (informally):

$$\downarrow_1 q_1 \cdot (\downarrow_1 q_1 \mid \downarrow_2 q_1)^*$$

(one first child move followed by an arbitrary sequence of first child and next sibling moves). Note, that all nodes that can be reached by such a sequence of move are *descendants* of the root node.

Our **TopDownRun** algorithm is therefore specialized as follows:

1. We compute from the formulas of the transitions we consider (Line 4, Figure 5) the set of states  $Q_1$  (resp.  $Q_2$ ) that are reached by a  $\downarrow_1$  move (resp.  $\downarrow_2$  move).
2. For  $Q_1$  (resp.  $Q_2$ ), we compute which set of labels cause a transition from one of the states in  $Q_1$  to go in a state not in  $Q_1$ . These labels make the node that have them relevant.



---

**Input** An automaton  $\mathcal{A} = (\mathcal{L}, \mathcal{Q}, \mathcal{T}, \mathcal{B}, \delta)$  and a sequence  $s$  of nodes in pre-order.

**Output** A mapping  $\mathcal{R}$  from states to sets of nodes, such that  $\mathcal{R}(q)$  is the set of nodes in the tree that were marked during an accepting bottom-up run that maps  $q$  to the root.

**BottomUpRun**( $\mathcal{A}, s$ )

1. **If**  $s$  is empty **Then**
2.     **Return**  $\emptyset$
3. **Else**
4.      $x, s' := \text{Head}(s), \text{Tail}(s)$
5.      $\mathcal{R} := \text{TopDownRun}(\mathcal{A}, x, \mathcal{Q})$
6.      $\mathcal{R}', s'' := \text{MatchAbove}(\mathcal{A}, x, s', \mathcal{R}, \text{Nil})$
7.     **Return**  $\mathcal{R}'$

**Input** An automaton  $\mathcal{A} = (\mathcal{L}, \mathcal{Q}, \mathcal{T}, \mathcal{B}, \delta)$ , a node  $x$ , a sequence  $s$  of nodes (in pre-order) occurring after  $x$  in pre-order, a mapping  $\mathcal{R}_1$  from states to sets of descendant nodes of  $x$ , and a node  $stop$  that is an ancestor of  $x$  (or Nil).

**Output** A mapping  $\mathcal{R}$  from states to sets of nodes and a sequence  $s'$  of nodes such that  $\mathcal{R}(q)$  is the set of nodes that were marked during an accepting bottom-up run that maps  $q$  to the node  $stop$  and a sequence  $s$  of nodes that are not descendants of  $stop$ .

**MatchAbove**( $\mathcal{A}, x, s, \mathcal{R}_1, stop$ )

8.      $p := \text{Parent}(x)$
  9.     **If**  $p = stop$  **Then**
  10.     **Return**  $\mathcal{R}_1, s$
  11.     **Else**
  12.     **If**  $s$  is empty or not(IsAncestor( $p, \text{Head}(s)$ )) **Then**
  13.      $\mathcal{R}_2, s'' := \emptyset, s$
  14.     **Else**
  15.      $x', s' := \text{Head}(s), \text{Tail}(s)$
  16.      $\mathcal{R} := \text{TopDownRun}(\mathcal{A}, x', \mathcal{Q})$
  17.      $\mathcal{R}_2, s'' := \text{MatchAbove}(\mathcal{A}, x', s', \mathcal{R}, p)$
  18.      $trans := \{(q, \ell \rightarrow \phi) \in \delta \mid \exists i \in \{1, 2\}, \exists q' \in \text{dom}(\mathcal{R}_i), \downarrow_i q' \text{ occurs in } \phi, \text{Tag}(p) \in \ell\}$
  19.      $\mathcal{R}' := \{q \mapsto R \mid \mathcal{R}_1, \mathcal{R}_2, p \vdash_{\mathcal{A}} \phi = (\top, R), (q, \ell \rightarrow \phi) \in trans\}$
  20.     **Return**  $\text{MatchAbove}(\mathcal{A}, p, s', \mathcal{R}', stop)$
- 

Figure 6. Bottom-up evaluation function.

root node in a top state  $q \in \mathcal{T}$ . Note that, if naively performed, such a bottom-up run will visit many nodes repeatedly: if a node is the common ancestor of  $m$  potential match nodes, then it would be visited  $m$  times. Instead, we move bottom-up left-to-right, and only move upwards from the left-most potential match until we reach its lowest common ancestor with the next potential match. This technique is similar in spirit to shift-reduce parsing ([56]). Our bottom-up matching algorithm is given in Figure 6.

The behavior of this algorithm is explained in detail on an example in Figure 7 (in this figure, vertical lines denote FirstChild edges, horizontal lines denote NextSibling edges, and dashed lines represent skipped subtrees). Intuitively, it takes as input a sequence of potential matches. For each of them it runs first the **TopDownRun** function to verify the downward context of the query (Lines 5 and 16). It then proceeds to walk upward (using  $\text{Parent}(\_)$ ) from a potential match node but stops (*shift*) when it reaches an ancestor of the next match. The following matches are recursively handled, and the algorithm can restart (*reduce*) when all the descendants of the current node have been treated.

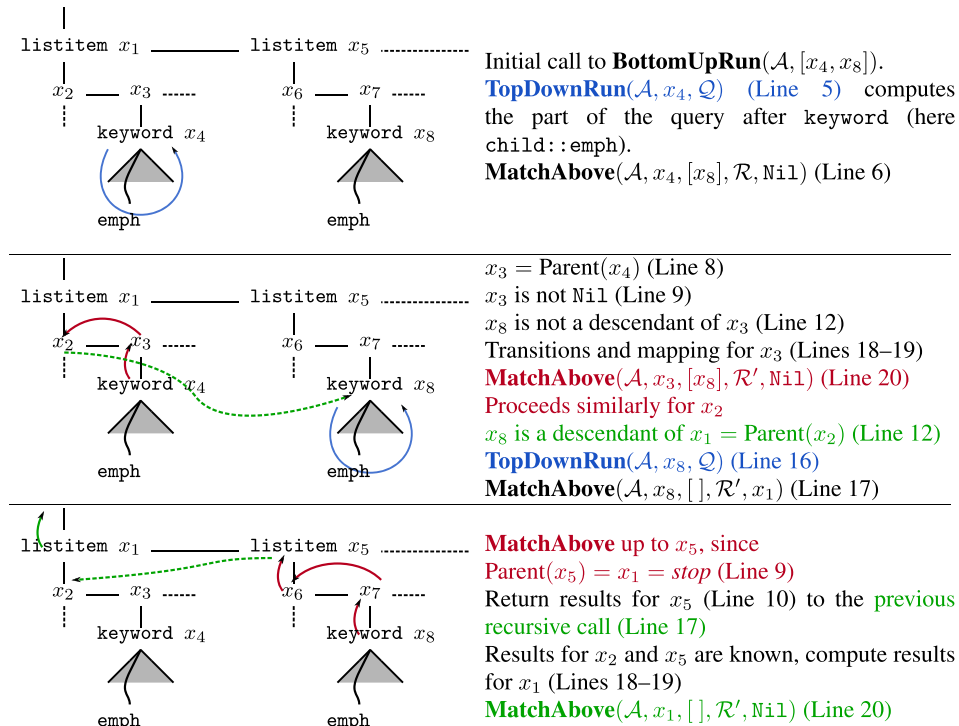


Figure 7. Illustration of the bottom-up run.

For our algorithm to be sound, we need to be sure that the subtrees that are not visited (e.g., left subtree of  $x_2$  and  $x_6$  right subtrees of  $x_3, x_5$ , and  $x_7$  in Figure 7) can safely be ignored. This is the case if the query at issue has the form

$$/axis::step/\dots/axis::step [pred]$$

(as in our previous example). Indeed, if a query has predicates in intermediate steps, for example:

$$/descendant::listitem[descendant::bold]/descendant::keyword[child::emph] ,$$

then we might need to explore the whole subtree of a `listitem` if no `bold` node is found on the path between a keyword node and its `listitem` ancestor. We chose not to explore such extensions to the **BottomUpRun** algorithm, because it would, in our opinion, not yield a significant speedup with respect to the **TopDownRun** algorithm. On the contrary, the speedups provided by the version of **BottomUpRun** we presented in this section were sufficient to justify the use of a dedicated procedure for low-selectivity queries of a particular shape.

### 5.5. General optimizations, on-the-fly determinization

Although the optimizations presented in the previous sections give the most important speedup, we describe hereafter a series of implementation techniques used for the efficient evaluation of automata.

**5.5.1. Hash consing of data structures.** We use hash consing for all critical data structures: sets of states, formulas, sets of transitions, sets of labels, and so on. Hash consed values have the following two properties. First, structurally equal values share the same representation in memory. Therefore, testing for equality of such values (e.g., testing that two sets of transitions are equal) consists in comparing their memory address (which is cheap). Second, to each such value, we can associate a unique integer id (this can be its memory address, e.g., but more interestingly a small



integer assigned at the creation of the value). These two properties—especially the second one—are instrumental to the other optimizations. Indeed, as described in [57], we can memoize (or cache) the results of expensive computations and reuse them when needed instead of recomputing them. We can associate to each function a table, indexed by the argument's id. Although the first computation might be expensive, its result is stored once and for all in the table and can be retrieved with one pointer indirection later on, when the same computation is requested. We explain now how this generic technique comes into play for automata evaluation.

**5.5.2. Just-in-time compilation of automata.** In the **TopDownRun** algorithm (Figure 5), the most expensive operations are in Lines 4, 5, and 8. By expensive, we mean that they take time  $O(m)$  where  $m$  is the number of steps in the original query. At Line 4, we gather all the transitions that can be selected from the current label  $\ell$  and set of states  $Q_{td}$ . From these, we compute, at Line 5, the new set of states  $Q_{td}^1$  and  $Q_{td}^2$  onto which we will launch the recursive call. As explained in Section 5.4, from the set of states  $Q_{td}^1$  (resp.  $Q_{td}^2$ ), we compute the ‘jump’ moves that the automaton will do to reach the next node in the left (resp. right) subtree. If none of the formulas requires the evaluation of a value predicate (such as `contains`, for instance), then we can see that this whole computation of Lines 4 and 5 can be cached in a two-dimensional array, using only  $\ell$  (the current label, identified by a small integer) and  $Q_{td}$  (a hash-consed set of states with a unique small id) as a key. In practice, we store in this table a small sequence of instructions that are computed *at run time* and which represent the behavior of the automaton for the next step (e.g., ‘jump to the next keyword label in state  $\{q_0, q_1\}$ ’). This just-in-time compilation scheme absorbs in practice most of the overhead caused by the automaton machinery and makes running an automaton almost as fast as executing a handwritten, precompiled function. In the same fashion, the computation of the judgment  $\vdash_{\mathcal{A}}$  can be memoized, this time in two parts. First, the sets of states (i.e., the domain of the resulting mapping) is stored once and for all, and second, a sequence of instructions telling how to propagate the results from the left and right subtrees is stored and evaluated for each node.

**5.5.3. Handling of result sets.** Maintaining sets of (result) nodes can be expensive. Our efficient management of sets of nodes relies on the following two observations. First, note that only the states outside of filters actually accumulate nodes. All other states always yield empty bindings. Thus, we can split the set of states into marking and regular states. This reduces the number of  $\heartsuit$  and  $\diamondsuit$  operations on result sets. Note also that given a transition  $q_i, \ell \rightarrow \downarrow_1 q_j \wedge \downarrow_2 q_k$  where  $q_i, q_j$ , and  $q_k$  are marking states, all nodes accumulated in  $q_j$  are in the left subtree of the current node. Likewise, all the nodes accumulated in  $q_k$  are subtrees of the right subtree of the current node. Thus, both sets of nodes are disjoint, and we do not need to keep sorted sets of nodes but only need sequences which support  $O(1)$  concatenation. Computing the union of two result sets  $R_j$  and  $R_k$  can therefore be carried out in constant time, and consequently,  $\heartsuit$  and  $\diamondsuit$  can be implemented in constant time. Furthermore, if we are only interested in obtaining the number of results of a query, these sets can be replaced by integer counters. Marking a node corresponds to incrementing a counter and merging two sets (i.e., performing  $\heartsuit$  or  $\diamondsuit$ ) corresponds to adding two counters. The evaluation of formulas as well as the preorder traversal we perform guarantees that marked nodes are not counted twice (i.e., in our algorithms, the results  $\mathcal{R}$  become mapping from states to integers rather than mapping from states to sets of nodes).

**5.5.4. Lazy result sets.** Another way to leverage the speed and jumping capabilities of our tree index is by making use of a lazy result set. Consider the query `/descendant::listitem/descendant::keyword`. When reaching a `listitem` node, the automaton is in a state that encodes the following behavior: ‘accumulate all `keyword` nodes below this node’. Therefore, instead of having the automaton jump through the subtree to individually put each `keyword` node in the result set, we only store the `listitem` node (i.e., the current node during evaluation) and a flag to remember that during serialization, it is not the `listitem` node that should be printed but rather all its `keyword` descendants. Because our tree index allows us to reach each such node using a constant time jump operation, we delay the process of obtaining all the final result nodes until serialization, therefore speeding up the marking process. This not only saves time but also

memory, because the full set of nodes does not have to be materialized. When evaluating the query in *counting* mode (as described in the previous point), then we replace the lazy result set by a single call to, for example, `SubtreeTags(_, keyword)`, which returns in constant time the number of keyword-labeled nodes. This number is then added to the counter of the corresponding state.

*5.5.5. Early evaluation of formulas.* Another optimization consists in evaluating the Boolean formulas of the automaton as early as possible. First, remark that in the **TopDownRun** algorithm, a node is ‘visited’ three times. Once when the automaton enters the node, during the top-down phase (Line 1). Here, we only know that at most all states in  $Q_{id}$  yield a successful run. Then, when returning from the left subtree (Line 6), we know  $\mathcal{R}_1$ , that is, the states that yield an accepting run for the left subtree. The idea now is to perform a partial evaluation of formulas by using only  $\mathcal{R}_1$ . If this happens to be sufficient to prove or disprove the states in  $Q_{id}$ , then the right subtree can be skipped altogether. This optimization is very important for filters as it insures that, for instance, in a query such as `/descendant::listitem[./descendant::keyword]`, the run function only tests for the presence of the left-most keyword node below a listitem node.

*5.5.6. Relative tag position tables.* As explained earlier, the transitions for the query `.../descendant::keyword/...` would be (just-in-time) compiled into a piece of code performing a subtree traversal using `TaggedDesc(_, keyword)` and `TaggedFoll(_, keyword)` instead of `FirstChild` and `NextSibling`. This is already optimal for documents where keyword nodes may appear arbitrarily. However, it is often the case that labels are not recursive (i.e., nodes with a label  $l$  do not occur below other  $l$ -labeled nodes). To further optimize the compilation of the automaton, we build –while indexing the document– four relative position tables, telling for each label  $l$  in the document the sets of labels that occur respectively in child position, descendant position, following-sibling position and following position. When compiling at run time the automaton and generating a call to `TaggedDescendant` for a label  $l$ , we check that this  $l$  label can indeed appear as descendant of the label of the current node (and similarly for other jumping functions). If the label does not occur, then the `TaggedDescendant` call is replaced by a constant function returning directly the correct sets of states for the left subtree as well as an empty result set, as if the automaton had made a full run on this subtree.

## 6. EXPERIMENTAL RESULTS

This section presents our experimental results and is organized as follows. We first describe our experimental settings, test machine, and benchmark data. We then provide a first round of experiments illustrating the raw performances of the tree and text index: indexing time and resulting index size, direct querying of the text index, and performing full preorder traversal using `FirstChild` and `NextSibling` moves. A third subsection illustrates how the tree index and automata-based engine work together to achieve very fast tree-oriented query evaluation (in particular, using the jumping moves described in Section 4.1.2). We then show how the automaton machinery can leverage the speed of both the text and tree index by evaluating queries containing both text and tree predicates. Lastly, we illustrate the versatility of our approach: our engine is easily extended to support querying of XML document storing bio-genetic data as well as natural language.

We have implemented a prototype XPath evaluator based on the data structures and algorithms presented in the previous sections. Both the tree structure and the FM-index were developed in C++, whereas the XPath engine was written using the OCaml language.

### 6.1. Protocol

To validate our approach, we benchmark our implementation against two well-established XQuery implementations, MonetDB/XQuery and Qizx/DB. We describe our experimental settings hereafter.

*Test machine.* Our test machine features an Intel Core i5 platform featuring eight 3.33 Ghz cores, 16 GB of RAM and a S-ATA hard drive. The OS is a 64-bit version of Ubuntu Linux (11.04). The

kernel version is 3.0 and the file system used to store the various files is ext4, with default settings. All tests are run on a minimal environment where only the tested program and essential services were running. We use the standard compiler and libraries available on this distribution (namely g++ 4.6.1, libxml2 2.7.8 for document parsing, and OCaml 3.12.1).

*Qizx/DB.* We use version 4.1 of Qizx/DB engine (free edition), running on top of the 64-bit version of the JVM (with the `-server` flag set as recommended in the Qizx user manual). The maximal amount of memory of the JVM is set to the maximal amount of physical memory (using the `-Xmx` flag). We also use the flag `-r` of the Qizx/DB command line interface, which allows us to re-run the same query without restarting the whole program (this ensures that the JVM's garbage collector and thread machinery do not impact the performances). We use the timing provided by Qizx debugging flags and report the *serialization time* (which actually includes the materialization of the results in memory and the serialization).

*MonetDB/XQuery.* We use version Oct2010-SP1 of MonetDB, and in particular, version 4.40.3 of MonetDB4 server and version 0.40.3 of the XQuery module (*pathfinder*). We use the timing reported by the `-t` flag of MonetDB client program, `mclient`. We keep the materialization time and the serialization time separate.

*Running times and memory reporting.* Each query is executed 11 times in a row within the same program instance (using `-r 11` for Qizx and by evaluating the same query 11 times for MonetDB without restarting the server). Of these 11 runs, we discard the first and average the remaining 10, which we report as 'running time' in the subsequent experiments. For all engines and all queries, the first run is always much slower (due to cold cache issues, garbage collector adjustments, and so on). The running time does not take into account query parsing or the query optimization phases that take place before the actual query evaluation. We found these times to always be negligible for all queries and engines tested (around 1 ms or less). We monitor the *resident set size* of each process, which corresponds to the amount of process memory actually mapped in physical memory. For the tests in which serialization is involved, we serialize to the `/dev/null` device (i.e., all the results are discarded without causing any output operation). We also ascertained that for all queries, all engines give the same node count and serialize *roughly* the same amount of data (small variations exist, because, e.g., empty elements can be rendered `<a/>` or `<a></a>`).

*Test data.* Our test data is comprised of XMark documents [24] of various sizes (between 116 MB and 1 GB), a 83 MB treebank document<sup>§</sup> and a 122 MB Medline document<sup>¶</sup>. In Section 6.6 that investigates XPath text-oriented queries, we also experiment using a word-based FM-index within SXSI, and do this over a 2.3 GB mediawiki document (part of the English 'wiktionary'). Last, in Section 6.7, we experiment with a 132 MB XML document composed of gene annotations and their DNA sequences.

*Remarks.* We also compared with Tauro [3]. Yet, as it uses a tailored query language, we could not produce comparable results.

## 6.2. Indexing

Our implementation features a versatile index. It is divided into three parts. First, the tree representation composed of the parenthesis structure, as well as the tag structure. Second, the FM-index encoding the text collection. Third, the auxiliary text representation allowing fast extraction of text content.

It is easy to determine from the query which parts of the index are needed in order to solve it, and thus, load only those into main memory. For instance, if a query only involves tree navigation,

<sup>§</sup><http://www.cs.washington.edu/research/xmldatasets>

<sup>¶</sup><http://www.ncbi.nlm.nih.gov/pubmed>

Table II. Search times of FM-index (in milliseconds), sampling factor  $l = 64$ .

Query	Pattern	<i>GlobalCount</i>		<i>ContainsCount</i>		<i>Report-</i>	Mem (MB)
		Number	Time	Number	Time	<i>Contains</i>	
1	Bakst	1	.004	1	0.04	0.012	61
2	ruminants	22	.009	19	2.3	1.6	61
3	morphine	392	.009	144	4.6	4.5	61
4	AUSTRALIA	438	.009	438	29.9	32.7	61
5	molecule	1472	.008	966	128.3	122.0	61
6	brain	2685	.005	1493	218.5	215.2	61
7	human	6897	.005	4690	553.5	548.0	62
8	blood	10,402	.005	8534	401.2	399.7	62
9	from	20,859	.004	12,073	1723	1718	62
10	with	63,332	.004	22,974	5084	5084	63
11	in	238,638	.003	42,586	19,642	19,630	64
12	a	2,932,251	.001	595,716	189,299	188,377	93
13	\n	9,730,750	.001	5,870,474	132,780	132,241	86

Mem, memory.

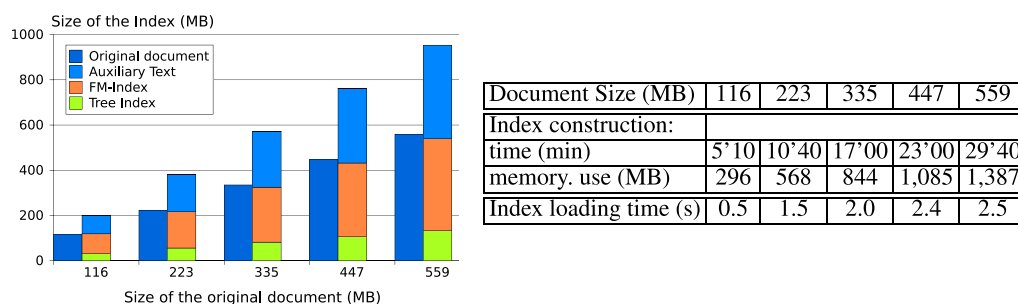


Figure 8. Indexing of XMark documents.

then having the FM-index in memory is unnecessary. On the other hand, if we are interested in very selective text-oriented queries, then only the tree part and FM-index are needed (both for counting and serializing the results). In this case, serialization is a bit slower (due to the cost of text extraction from the FM-index) but remains acceptable because the number of results is low; see Table II.

Figure 8 reports the construction time and memory consumption of the indexing process, the loading time from disk into main memory of a constructed index, and a comparison between the size of the original document and the size of our in-memory structures.

For these indexes, a sampling factor  $l = 64$  (cf. Section 3) was chosen. It should be noted that the size of the tree index plus the size of the FM-index is always less than the size of the original document.

It should further be noted that although loading time is acceptable, it dominates query answering time. This is however not a problem for the use case we have targeted: a main memory query engine where the same large document is queried many times. As mentioned in Section 1, systems such as MonetDB load their indexes only partially; this gives superior performance in a cold-cache scenario when compared with our system.

### 6.3. Raw performance of text index

Here, we give a short overview of the performance of our implementation of the FM-index. We present the search times for different versions of *contains*-queries:

1. *GlobalCount(P)*: returns the global number of occurrences of the pattern  $P$  in all texts.
2. *ContainsCount(P)*: returns the number of texts that contain  $P$ .
3. *ContainsReport(P)*: returns the positions of all occurrences of  $P$  in the texts.

Our experiments are over the text collection obtained from a 122 MB Medline XML document. The size of this text is around 82 MB (if stored in 1-byte per character ASCII format). Our ‘plain’ alternative to the FM-index is a naive (byte-wise) string buffer (using precisely 82 MB of memory). To search over the plain buffer, we use OCaml’s regular string expression library. The naive search time is constant for all our queries at around 2700 ms. For both the naive and the FM-index, the result positions (32 bit integers) for *ContainsReport* queries are materialized in an array. Consider now the performance of our FM-index in comparison. First, at sampling factor  $l = 64$ , shown in Table II. As can be seen, the times for *ContainsCount* and *ContainsReport* for the word ‘from’ are at around 1720 ms. Thus, in this case, it is still faster to search over the FM-index. On the other hand, for the word ‘with’, the search time is over 5000 ms, thus, here the plain search becomes faster. Hence, somewhere between 20,859 and 63,332 occurrences lies the cutoff point from which on searching over the plain text is faster than over the FM-index. Table III shows timings obtained with sampling factor  $l = 4$ . As can be seen the cutoff point is now much later, at a global count somewhere between 411,409 and 748,326. The last columns of Tables II and III show the maximal memory consumption for these queries over the FM-index. As mentioned in the beginning of this section, we measure the maximum memory used by the process, as reported by the operating system (this is a slight over-approximation of the actual memory). The memory overhead for queries with large cardinality, such as the last queries (q13 and q15), is explained by the size of the result array: for both sampling factors, this is around 25 MB. This query has around 6 million results (*ContainsCount*-number); each result is stored as a 4 byte integer. Thus, 23 MB are needed. However, additional memory overhead occurs when results are removed from the *GlobalCount* (because they occur in the same XML text node). For instance, in the second to last query (q12/q14), the ratio of *GlobalCount*-number to *ContainsCount*-number is much larger than for the last query (4.9 versus 1.7). This means that on average, there are around 5 ‘a’-characters per text node, while there are only around 1.7 return-characters per text node. Correspondingly, the maximum memory consumption is much higher too.

#### 6.4. Raw performance of tree index

The performance of some low-level features of our tree index is compared with the corresponding performance of a standard pointer-based implementation of a tree. The latter provides for each tree node two 64-bit pointers to its first child and next sibling nodes (and does not store labels). We first compare construction times. Then, we compare times for a full depth-first left-to-right tree traversal on the different structures. Finally, we test the speed of the taggedDesc and taggedFoll functions.

Table III. Search times of FM-index (in milliseconds), sampling factor  $l = 4$ .

Query		<i>GlobalCount</i>		<i>ContainsCount</i>		<i>Report-</i>	Mem (MB)
		Number	Time	Number	Time	<i>Contains</i>	
1	Bakst	1	.005	1	0.049	0.013	100
2	ruminants	22	.010	19	0.156	0.086	100
3	morphine	392	.009	144	1.7	1.4	100
4	AUSTRALIA	438	.009	438	4.1	3.9	100
5	molecule	1,472	.009	966	6.2	5.9	101
6	brain	2,685	.006	1,493	12.2	11.6	101
7	human	6,897	.005	4,690	25.4	27.3	101
8	blood	10,402	.005	8,534	77.2	73.6	101
9	from	20,859	.003	12,073	84.0	78.7	101
10	with	63,332	.004	22,974	242.8	235.0	102
11	in	238,638	.002	42,586	1,105	1,091	103
12	b	411,409	.001	135,307	1,779	1,762	108
13	g	748,326	.001	320,440	3,412	3,379	119
14	a	2,932,251	.001	595,716	13,183	13,173	133
15	\n	9,730,750	.001	5,870,474	87,771	88,230	126

Mem, memory.



Table IV. Construction times (in milliseconds) for pointer versus SXSI tree store.

File	Parse	Pointers	Parentheses	Tags	Tag-tabs
XMark116M	89,446	373	504	4682	1324
XMark223M	220,143	716	976	9051	2544
XMark559M	620,479	7923	2415	22,857	6283
Treebank83M	67,412	465	615	14,067	18,867
medline122M	67,935	537	760	6933	2036

Table V. Traversal times (in milliseconds), #nodes (in millions).

file	Recursive, all nodes			Element nodes, SXSI		
	#nodes	Pointer	SXSI	#nodes	Rec.	//*
XMark116M	6	33	109	1.7	71	153
XMark223M	12	63	209	3.3	137	296
XMark559M	30	164	535	8.4	345	756
Treebank83M	7	57	184	2.4	136	292
medline122M	9	48	164	2.9	112	244

Table VI. Times (in milliseconds) for tagged traversals over XMark116M.

tag	#nodes	jump(C++)	//(cou)	//(mat)
category	1,040	1.2	1.6	1.7
price	10,141	2.3	2.9	3.1
listitem	63,179	16	22	24
keyword	73,070	11	12	14

We compare different traversals through all nodes with a given label: (i) using a pure C++ function, (ii) using our automata in counting mode, and (iii) using our automata in materialization mode.

*Construction.* As Table IV shows, the construction of the parentheses structure takes roughly 1.5 times the amount of time of allocation a pointer structure for the tree. Constructing the tag sequence is considerably slower, about 10 times as much as building the parentheses structure. This is because, for each opening and for each closing tag, a separate `sarray` is constructed (see bottom left of Figure 1). The last column shows the time for building the four tag-to-tag tables described in Section 5.5.6. We also show the XML parsing time in the first column of the table, which dominates the rest of construction times.

*Full traversals.* The left part of Table V shows that a full tree traversal through all nodes is between 3.2 and 3.4 times slower with SXSI, than with a pointer tree data structure. Note that the pointers are allocated in preorder too, giving optimal performance for preorder traversal. As a comparison, if the pointers are allocated in postorder, then traversal time for the preorder traversal is almost twice as slow as the numbers reported, and if pointers are allocated in in-order, then the times are a bit over twice as slow; see [43] for a discussion of the phenomenon. It should also be noted that for other access patterns, such as random root-to-leaf traversals, the time difference between pointer and succinct trees is much larger, factors of up to 100 are measured in [43].

In the right part of Table V, we see the number of element nodes in these trees, and the time it takes for SXSI to recurse through those nodes: either using a small recursive C-function (column ‘rec.’), or using the automaton for the XPath query `//*`, and executing in counting mode.

*Tagged Traversals.* Here, the speed of the `TaggedDesc` and `TaggedFoll` functions is investigated. Using these two functions, three different traversal through all nodes with a given label are considered: first, by a small C++ function, and second and third by our automata through a `//label` query in counting and materializing modes, respectively. For instance, Table VI shows that iterating through

```

T01 //NP
T02 //S[//VP and //NP]/VP/PP[IN]/NP/VBN
T03 //NP[//JJ or //CC]
T04 //CC[ not(//JJ) ]
T05 //NN[//VBZ or //IN]*/[//NN or //_QUOTE_]
X01 /site/regions
X02 /site/regions/*/item
X03 /site/closed_auctions/closed_auction/annotation/description/text/keyword
X04 //listitem//keyword
X05 /site/closed_auctions/closed_auction[ annotation/description/text/keyword ]/date
X06 /site/closed_auctions/closed_auction[ //keyword]/date
X07 /site/people/person[ profile/gender and profile/age]/name
X08 /site/people/person[ phone or homepage]/name
X09 /site/people/person[ address and (phone or homepage) and (creditcard or profile)]/name
X10 //listitem[not(//keyword/emph)]//parlist
X11 //listitem[ (//keyword or //emph) and (//emph or //bold)]//parlist
X12 //people[ //person[not(address)] and //person[not(watches)]]//person[watches]
X13 /*[ // * ]
X14 /*
X15 /*/*
X16 /*/*/*
X17 /*/*/*/*

```

Figure 9. Tree-oriented queries based on Treebank (T01–T05) and XMark (X01–X17).

all keyword-nodes of the 116 MB XMark document takes essentially the same time for all three methods (11–14 ms). This is in contrast to some other labels: for `listitem`, for instance, the count-automaton traversal is 1.5 times slower than the C++ traversal. This can be explained by the fact that `listitem` is a recursive tag: there are in fact 23,298 `listitem` nodes that appear as descendants of `listitem` nodes. Hence, at each `listitem` node the automaton issues a tagged-Descendant to search for further nodes. The other labels such as `keyword` and `category` do not appear recursively. Because this information is part of our tree index (cf. Section 5.5.6), the automaton run function avoids all these taggedDesc calls, which brings the speed almost up to the one of the C++ function.

### 6.5. XPath tree queries

We benchmark tree queries using the queries given in Figure 9. Queries X01 to X12 are taken from the XPathMark benchmark [58], derived from the XMark XQuery benchmark suite. X13 to X17 are ‘crash tests’ that are either simple (X13 selects only the root, because it always has at least one descendant in our files) or generate roughly the same amount of results but with various intermediate result sizes. Queries T01 to T05 work on the Treebank file.

*Query answering time.* For this experiment, we use the Treebank document (83 MB) and XMark documents (116 MB and 1 GB). In the cases of MonetDB and Qizx, the files were indexed using the default settings. Let us first describe in detail Figure 10, which summarizes the running time for XMark queries. Each of the six graphs should be read as follows. For each query (X01 to X17), the graph reports as vertical bars the relative running time of the three engines with respect to SXSI’s running time (therefore, SXSI’s score is always 100%). In these graphs, a higher bar means that the engine was slower. We also give at the top of each bar the average running time for the query in milliseconds (or seconds, if the number is suffixed with an ‘s’). For instance, in the first graph—labelled ‘116 MB (counting)’—we can see that for query X01, SXSI evaluates the query in 1.8 ms, MonetDB 7.4 ms (or roughly 400% of SXSI’s time) and QizX 3.6 ms (or roughly 200% of SXSI’s time). For count queries, the timing for all three engines are given side by side (SXSI, MonetDB, and QizX in that order). For full reporting queries however, we want to gauge precisely the amount of time spent during materialization and during serialization. The definition of materialization seems to fit the evaluation model of both MonetDB and SXSI: create a data structure in memory, which holds the resulting nodes *in order* and *without duplicates* such that access of the first result in preorder can be carried out in constant time, and accessing the next resulting node from the current one in preorder can also be carried out in constant time. The timing for both SXSI and MonetDB are

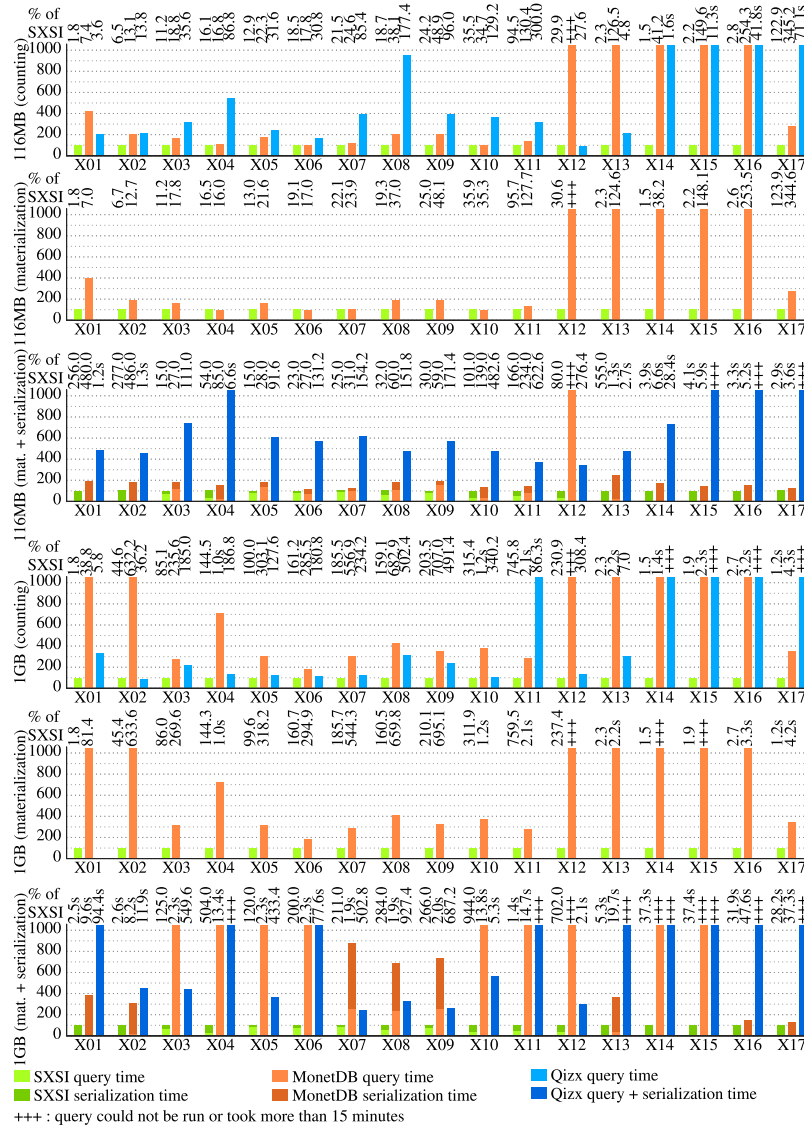


Figure 10. Running times for XMark queries, in milliseconds or seconds and as percent of SXSI’s speed. Lower bars are better.

given in the graphs labelled ‘(materialization)’. Because Qizx interleaves evaluation of the query and serialization, we only compared it with SXSI and MonetDB in the ‘(materialization+serialization)’ series. We also checked that all three engines generated in the end the same amount of data while in serialization mode and that they generated valid XML documents (in particular, special characters were replaced by their corresponding XML entities—for instance ‘&’ is rendered as ‘&amp;’).

From the results of Figure 10, we see how the different components of SXSI contribute to the efficient evaluation model. Fully qualified paths, such as queries X01–03 and X5 illustrate the sheer speed of the tree structure and in particular, the efficiency of its basic operations (such as FirstChild and NextSibling, which are used for the child axis), as well as the efficient execution scheme provided by the automaton. The descendant axes (used, e.g., in X04, X06, X10–12) show the impact of the jumping primitives and the computation of relevant nodes. Complex filters (X06–12) show how the alternating automata can efficiently evaluate complex Boolean formulas corresponding to structural conditions over subtrees of a given node, including negations of paths and nested predicates. Finally, X12 to X17 illustrate the robustness of our automaton model. Indeed while such queries might seem unrealistic, the good performances that we obtain are the combination of (i)

using an automata-based evaluator (which factors in the states of the automaton all the necessary computation and thus, do not materialize unneeded intermediate results) and (ii) our implementation of lazy result sets, which shifts the burden of walking through the document as much as possible to the serialization process.

As for Treebank queries, the results are reported in Figure 11. As one can see, SXSI behaves well for simple queries (such as T01) and for more complex ones (T02, T04 and T05). An important point is also the following one: the running times for all queries and for *all engines* are much worse than the ones obtained for larger XMark documents. It seems that, despite using three different approaches (tree automata+jump, staircase join or Qizx’s proprietary one), a high number of distinct paths and labels impact query results much more than the final node count.

*Jumping and memoization of computations.* As the experiments of the previous section show, SXSI often outperforms MonetDB and Qizx. To better highlight how the optimizations described in Sections 5.4 and 5.5 impact the running time of Algorithm 5, we selectively disabled some of them and executed queries X01 to X17 on a 116 MB XMark document. The results are shown in Figure 12. The first bar of each cluster gives the running time of a naive version of the **TopDownRun** algorithm, where each node of the document is traversed, and the computation of the next transitions to evaluate is recomputed for each node. In other words, in this run, the  $|D|$  (size of the document) and  $|Q|$  (size of the automaton/query) factors are paid in full. Here, we can see that query answering time depends on both the size of the query (which increases from X01 to X12) but also the number of selected nodes (e.g., X14 is a short query but selects all nodes, causing memory allocations and result sets book-keeping—which are expensive operations—for each node in the document).

The second bar in each cluster reports query answering time when computation and jumping to relevant nodes are performed (but these computations are not memoized). This includes not only

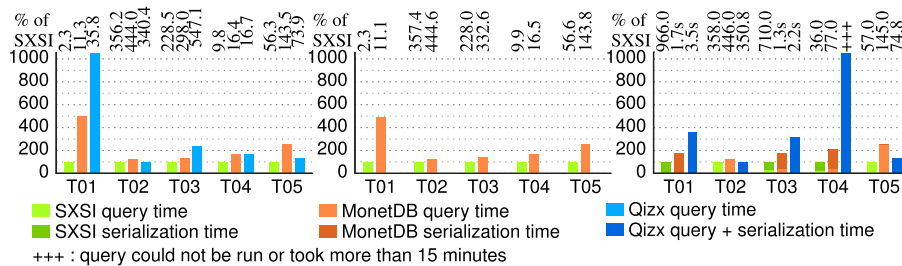


Figure 11. Running times for Treebank queries, in milliseconds or seconds and as percent of SXSI’s speed. From left to right, counting, materialization, and serialization times. Lower bars are better.



Figure 12. Impact of jumping and hash consing on the running time of **TopDownRun** (Figure 5). Logarithmic scale, lower bars are better.

jumping to relevant nodes but also discarding whole subtrees and using the constant-time subtree counting function of the tree interface (Section 5.5.4). It comes to no surprise that skipping nodes aggressively yields important speedups in query answering time (some queries being translated into a single low-level function call). However, queries such as X17 show the limit of the approach. For these queries ( $//*//*$  and  $//*//*$  resp.) the engine reaches all the element nodes at depth 3 (resp. 4) and select in constant time the set of elements in their subtrees. Because of the shape of the document, there are much more nodes at depth 4 than at depth 3, and therefore, the automaton has to traverse more nodes to touch all those at depth 4. This explains why the impact of jumping is much lower for query X17 than it is for X16.

The third bar in each cluster reports query answering time of a version of the **TopDownRun** that traverses the whole document, but memoizes the computations of Lines 4, 5, and 8 in Figure 5 in look-up tables. This also decreases the query answering time greatly. Furthermore, because this technique does not rely on a particular index, it shows that it would also benefit more traditional ‘Document Object Model based’ query engines, which lack the low-level jumping primitives.

Lastly, the fourth bar reports query answering time with both optimizations enabled. Here, we see that the jumping optimization is further improved by the memoization: the somewhat expensive process that computes the next relevant node and the corresponding jump function is stored in a look-up table and can be accessed in constant time when the automaton reaches again the same configuration.

*Memory use and precision.* Although it is straightforward to predict the memory consumption of our engine with respect to the index part (the full index is mapped in memory excluding the auxiliary text, see Figure 8), the behavior of the automaton evaluation function is unclear. Indeed, to speedup the computation, we create memoization tables, we handle partial result sets, and we perform recursive procedures that might be as deep as the *binary encoding* of the XML document (because we recurse on FirstChild and NextSibling move), thus increasing the size of the call stack.

We report in Figure 13 (left) the memory consumption for the automata evaluation of materialization queries. This includes the size of the recursive call stack, the size of OCaml’s heap (which is grown dynamically by OCaml’s garbage collector to accommodate the memory need). On the heap are allocated the memoization tables, intermediary structures and final result sets. As one can see, the memory use is very modest, peaking at 32 MB for query Q11. Although we do not compare directly with MonetDB or QizX for memory consumption (because these engines try to maximize the memory used to achieve better speed), we see that we can reach comparable (if not greater) speed while being very conservative memory wise.

To gauge the precision of our automata-based approach, we report in Figure 13 (right) for each query:

- the number of visited nodes (i.e., the number of nodes onto which the **TopDownRun** function is called);
- the number of marked nodes (i.e., the number of nodes that are marked as potential results during the evaluation);
- and finally, the number of result nodes for the query.

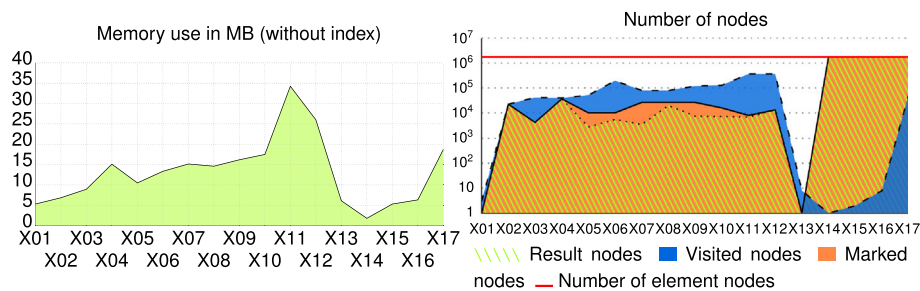


Figure 13. On the left, memory use in megabytes for XMark 116 MB documents (excluding the index). On the right, comparison of visited, marked and result nodes for each query (logarithmic scale).



A first observation is that the number of marked nodes is almost always the same as the number of result nodes. In particular, for queries of the form:

$$/axis::step/\dots/axis::step[pred] ,$$

we only mark nodes that are part of the final result (the first part of the path is resolved during the top-down phase, whereas the filter is validated during the bottom-up phase. Thus, once reaching a potential result the automaton has recognized its path to the root as well as its subtree and can decide whether to select it or not). For general queries, the experimental results show that early evaluation of Boolean formulas helps to decide early during query evaluation whether a node is indeed a result or not (even though some nodes are still selected and discarded later on). Another point of interest is that for several queries (X02, X04, X14–X17), we *only* visit result nodes. Although this might be expected for queries X14–X17, for which virtually every node is a result, queries such as X02 or X04 are very selective. However, these queries provide enough information for the runtime analysis of relevant nodes to be *exact* and therefore, only touch result nodes. In general of course, the number of traversed node is larger than the number of resulting nodes but always far less than the whole document. Queries X14 to X17 show the impact of lazy result sets where we mark several nodes (whole subtrees actually) in one function call and therefore, manage to return more nodes than we have actually visited.

Lastly, we can see that the shapes of the ‘Visited Nodes’ curve in Figure 13 (right) and the memory use (Figure 13 left) are quite similar (the former being flattened by the logarithmic scale). This (quite expectedly) shows that the number of visited nodes (and not the number of result nodes or intermediary results) impacts directly the memory consumption of our query engine.

## 6.6. XPath text queries

In this section, we illustrate how our approach can leverage the speed of the underlying text engine. Evaluation of text queries is performed as follows:

1. We determine during parsing whether the query can be run using the **BottomUp** algorithm (Section 5.4.2).
2. We determine whether the text predicates (`starts-with`, `ends-with`, `contains`, ...) are applied to a single text node, for example, if the context node in the XPath expression is reached using `axis::text()`, or if the content of a selected element is known to be PCDATA (this information is kept in the index). If that is the case, we use the FM-index and run the query bottom-up. Otherwise, we revert to using the naive text representation to ensure that the semantics of the XPath text function is preserved, and run the query using the **TopDownRun** algorithm.

Here, the check in Step 2 is necessary to implement the semantics of XPath’s text predicates over mixed content. Indeed, for such elements (containing both text and other elements), the semantics is to first create a text node resulting in the concatenation of all text elements and only then to perform the text predicate. For instance, on the document

```
<a>01<b>23</b>45</a>
```

the following query returns the root node as a result (i.e., the predicate evaluates to true):

```
/child::a[ contains( . , "1234" ) ] .
```

Because the text value of ‘.’ (i.e., the root node) is the concatenated string ‘012345’, which indeed contains the substring ‘1234’. Although this ‘feature’ of text predicates is rarely used in practice, its presence in the XPath specification can hinder the efficiency of ‘normal’ text queries.

**6.6.1. FM-index.** For this experiment, we compared the generic version of XPath text searching function (such as `contains`, `ends-with`, and `start-with`) for all three engines on a 659 MB Medline document (used to store bibliographic data about medical publications). The queries we

used for this experiment are given in Figure 14. In queries M01, M03, and M04, SXSI can use the FM-index but cannot run bottom-up, because they feature complex filters. On the other hand, queries M02 and M05–M09 can be run purely bottom-up. Lastly, queries M10 and M11 need to use the ‘naive’ text collection because the string that is searched for could overlap several text elements of the TextCollection (for M10, the `MedlineCitation` element has mixed content, whereas for M11, it is not known at query compile time which element will be searched for the text).

The timings for these queries for all three engines are given in Figure 15.

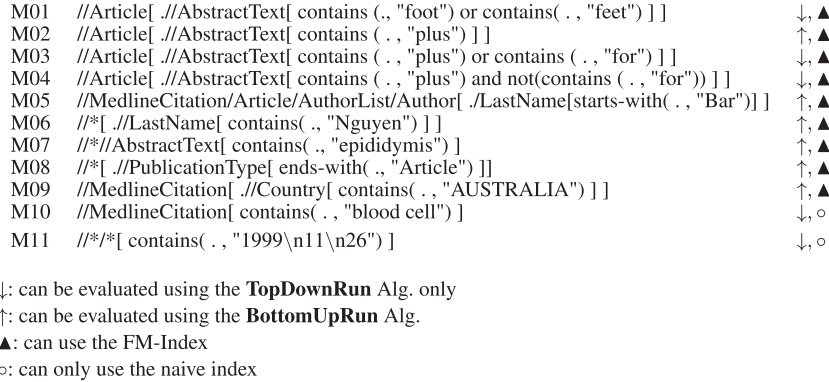


Figure 14. Text-oriented XPath queries over Medline and their evaluation strategies.

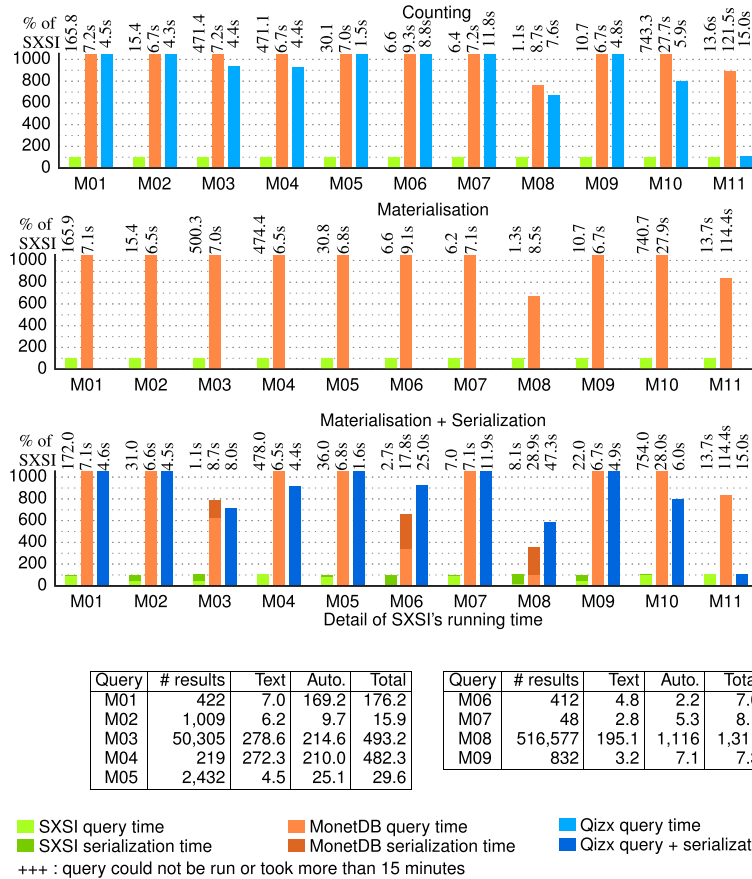


Figure 15. Running times for Medline queries, in milliseconds or seconds and as percent of SXSI’s speed. Lower bars are better.

```

W01 //Article[ //AbstractText[ contains (., "blood sample") ] ]
W02 //Article[ //AbstractText[ contains (., "is such that") ] ]
W03 //Article[ //AbstractText[ contains(., "various types of") and contains(., "immune cells") ] ]
W04 //Article[ //AbstractText[ contains(., "of the bone marrow") ] ]
W05 //Article[ //AbstractText[ contains(., "cell") and not(contains(., "blood")) ] ]

W06 //text[ contains (., "dark horse")]
W07 //text[ contains (., "horse") and contains(., "princess") ]
W08 //page/child::title[ contains (., "crude oil") ]
W09 //page[//text[ contains(., "played on a board")]]/title
W10 //page[//text[ contains(., "whether accidentally or purposefully")]]/title

```

Figure 16. Word-based queries on Medline (W01–W05) and Wikimedia (W06–W10) documents.

Table VII. Experimental results for word-based text queries.

Query	W01	W02	W03	W04	W05	W06	W07	W08	W09	W10
SXSI (ms)	5.7	4.9	136.8	5.7	143.5	10.5	1600	10.6	7.8	10.2
Qizx (ms)	86.3	78.3	137.3	158.0	150.7	–	–	–	–	–

The interpretation of the histograms is the same as for the tree-oriented queries. For queries M01–M09, for which the FM-index can be used, we also give a table with the details of the amount of time spent querying the whole text collection for matches and the amount of time spent traversing the tree with the automaton. We also give the number of matches of the query.

As expected, in the case where queries can be evaluated in a bottom-up fashion the improvement in running time is of several orders of magnitude. Even when this is not the case, the improvement of using the FM-index over the naive substring search completely justifies the somewhat longer indexing time and memory use. Lastly, for queries where the naive text representation must be used, the fact that our tree representation allows to return in constant time, the set of texts that occurs as descendant of a given node allows us to avoid costly materialization of large temporary text elements (as we can see in query M11, Qizx seems to also manage to avoid allocations, yielding a query answering time close to ours, whereas MonetDB seems to allocate large amounts of temporary character buffers, which impacts query answering time).

**6.6.2. Word-based text index.** The loose coupling of SXSI's components (automata-based query engine, succinct tree index, and text index) allows us to readily plug other sorts of tree structures or text indices. The text collection presented in Section 3 achieves exact symbol-based text pattern-matching at the cost of indexing time and query time for large results. We show in this section how one can choose another tradeoff, by plugging in the simple word-based text index by Fariña *et al.* [59]. This index achieves high indexing and querying speed, with little memory by limiting queries at a word boundary. In this index, distinct words are treated as distinct symbols, and the text collection is therefore viewed as a suffix array over a very large alphabet (the alphabet size is the number of distinct words in the original text). We compare SXSI equipped with this index to the full-text extension of Qizx. This extension implements the XQuery Full-Text facility [36], which allows to perform token-based (essentially word-based) queries. We used the queries in Figure 16 to test both indexes. Although in this figure we used the 'contains' function for conciseness, Qizx queries are implemented using the special 'ftcontains' operator (which performs word based queries).

We tested queries for both Qizx and SXSI on the 122 MB Medline document and for SXSI only<sup>†</sup> on a 2.3 GB mediawiki document (a snapshot of the english 'wiktionary' [60]). The results of the experiment are given in Table VII. As we can see, even when both SXSI and Qizx are allowed to make use of an efficient index, SXSI's bottom-up evaluation strategy clearly improves query answering time (all queries, but W03, W05 and W06 can be answered in 10 ms or less). For queries that involve several text predicates (i.e., that need to be evaluated in a top-down fashion), SXSI and Qizx perform similarly (and efficiently, even for large documents).

<sup>†</sup>Unfortunately only the standard version of Qizx hits a hard limit of around 2 GB for a single document (while the commercial 'XL' version of the engine supports documents of 1 TB).

```

<!ELEMENT chromosome (name, gene*) >
<!ELEMENT name #PCDATA >
<!ELEMENT gene (name, strand, biotype, status,
description?, promoter, sequence, transcript*) >
<!ELEMENT strand #PCDATA >
<!ELEMENT biotype #PCDATA >
<!ELEMENT status #PCDATA >
<!ELEMENT description #PCDATA >
<!ELEMENT promoter #PCDATA >
<!ELEMENT sequence #PCDATA >
<!ELEMENT transcript (name, start, end, exon*,
sequence, protein?) >
<!ELEMENT start #PCDATA >
<!ELEMENT end #PCDATA >
<!ELEMENT exon (name, start, end, sequence) >

```

Figure 17. DTD for bio-genetic data.

### 6.7. Biological sequence queries

As a last experiment, we demonstrate the versatility of SXSI by showing that it can be used as a very efficient biological database manager, answering queries that make use of both the tree structure and a tailored text index. More precisely, we create XML files that combine gene annotations with their DNA sequences. A sample DTD for these files is given in Figure 17. In this DTD, the elements `promoter` and `sequence` are of particular interest: they store the DNA represented as long sequences of A, T, C, and G characters. The other `#PCDATA` elements store the gene annotation data such as positions, names, and so on.

Our experiment data are composed from human chromosome five\*\*, which contains 2719 genes having in total 8330 different transcripts. For each gene, we include 1000 base pairs of its upstream promoter sequence, the gene sequence itself (all exons and introns included), and annotation information such as gene's biotype and description. Additionally, we include all known transcripts of each gene, that is, sequences of the exons they contain as well as the concatenation of these exons. The resulting textual content is highly repetitive, because each one of the exon sequences can appear in many transcripts. Highly repetitive data have been shown to compress well using certain run-length encoded text indexes [61], thus, here, the text index implementation is switched to use a run-length compressed suffix array [61] instead of the FM-index. In this example, the final XML file†† is 132 MB, whereas the text index requires only 63 MB of memory plus 59 MB for the samples array. The full index, including tree and text, is around 135 MB, that is, only as big as the original document. The resulting XML document contains 323,318 elements, of which 65,286 are either `promoter` or `sequence` nodes containing genetic data.

To do biologically relevant XML queries, we extend our engine to support position specific scoring matrix queries (PSSM), which allow us to search for transcription factor binding sites from genes' promoter regions. The input for this query is a position frequency matrix (PFM) and a minimum threshold for a valid match. The matrices can be found from the Jaspas database [62].

In a nutshell, PFMs have one row for each symbol of the alphabet (in our case 4 rows A, T, C, and G) and one column for each position in the pattern to search. For instance, the PFM:

$$\begin{bmatrix} A & 0 & 20 & 10 & 1 \\ T & 30 & 10 & 0 & 0 \\ C & 0 & 0 & 10 & 20 \\ G & 18 & 6 & 6 & 6 \end{bmatrix}$$

denotes patterns of length four, and the substring AGCT would obtain the score  $0+6+10+10 = 26$ . To form the PSSM query, the PFM matrix is first converted into log-odds form to take into account

\*\*Ensemble Human genome release 59, August 2010.

††<http://www.cs.helsinki.fi/group/suds/sxsi/data/>

query	# results	Text	Auto	Total
//promoter[ PSSM(., M1)]	134	85.1	7.40	92.5
//promoter[ PSSM(., M2)]	4	4.4	1.15	5.5
//promoter[ PSSM(., M3)]	1	6.5	0.38	7.0
//exon[ //sequence[ PSSM(., M1) ] ]	434	85.5	7.50	92.6
//exon[ //sequence[ PSSM(., M2) ] ]	25	4.3	1.28	5.6
//exon[ //sequence[ PSSM(., M3) ] ]	9	6.4	0.62	7.0
//*[ PSSM(., M1) ]	1,875	85.0	7.60	92.6
//*[ PSSM(., M2) ]	184	4.3	1.19	5.5
//*[ PSSM(., M3) ]	51	6.4	0.58	6.9

M1 : Jaspar ID =  
MA0031.1, length =  
8, threshold = 5,000

M2 : Jaspar ID =  
MA0050.1, length =  
12, threshold = 100,000

M3 : Jaspar ID =  
MA0017.1, length =  
14, threshold = 300,000

Figure 18. Running times for position specific scoring matrix (PSSM) queries (in milliseconds).

the uneven background distribution of nucleotide frequencies. Then, the PSSM query takes such a matrix as well as a threshold and returns all text elements whose content scores more than the given threshold. Figure 18 gives the running times for XPath queries using the PSSM predicates and run-length compressed suffix array, with block size 128 and sample rate 16, as the text index. The table summarizes also the number of results, the length of the search pattern and the value of the threshold.

It is interesting to remark that since the document has a very flat and shallow structure, the automaton/tree part of the query evaluates always very quickly (7 ms or under). The PSSM scheme also allows us to write biologically meaningful queries that would otherwise be impossible or very hard to write with regular expressions or a regular full-text extension. Yet, we did not have to modify our core engine, only the text index was modified in isolation to add PSSM capabilities; the automata and tree machinery remained unchanged.

## 7. CONCLUSIONS AND FUTURE WORK

We have presented SXSI, a system for compact in-memory representation of an XML collection and for fast-indexed XPath queries over the representation. Even in its current prototype stage, SXSI is already competitive with well-known efficient systems such as MonetDB and Qizx. A number of avenues for future work are open. We mention the broadest ones.

Handling updates to the collections is possible in principle, as there are dynamic data structures for sequences, trees, and text collections [16, 27, 63–65]. However, their practicality has not yet been established nor how they relate to classical schemes that maintain a log of changes and re-index periodically.

The compact data structures used in SXSI support several fancy operations beyond those actually used by our XPath evaluator. A matter of future work is to explore other evaluation strategies that take advantage of those nonstandard capabilities. As an example, the current XPath evaluator does not use the range search capabilities of the structure *Doc* of Section 3. This could be useful in the case of top-down evaluation of queries that contain nonselective text searches; after a top-down phase the search on *Doc* could be restricted to the range of a particular subtree.

An interesting challenge is to support XPath string-value semantics, where strings spanning more than one text node can be searched for. This, at least at a rough level, is not hard to achieve with our FM-index, by removing the \$-terminators and marking them on a separate bitmap instead.

We would like to extend our implementation to full XPath 1.0 and to add core functionalities of XQuery. The first step here is to add backward axes to our XPath fragment. Ideally, we would like to extend the lazy on-the-fly determinization procedure of our automata to a lazy remove-backward procedure, which removes backward axes during the run of the automaton. We do however need to execute automata with backward axes too, because not every query can be rewritten into a forward one (see Section 3.2 of [66]). Even with backward axes and with the semi-joins of XPath, evaluation can be performed in linear time (cf. [67]) and should be fast and predictable in SXSI. This efficiency will have to be given up in general, when we want to support more complex features such as numerical operations and loops of XQuery.



## ACKNOWLEDGEMENTS

We would like to thank Schloss Dagstuhl for the very pleasant and stimulating research environment it provides; the work of this paper was initiated during the Dagstuhl seminar ‘Structure-Based Compression of Complex Massive Data’ (Number 08261). In particular, the idea of sorting according to end markers came alive during the meeting, as briefly sketched in the report [68], and was then independently developed in different directions [25,64]. We are grateful to Kunihiro Sadakane for making available to us his implementation of parentheses structure for succinct trees, and to Juha Karjalainen for composing the BioXML data. Diego Arroyuelo and Francisco Claude were partially funded by NICTA, Australia. Francisco Claude was partially funded by NSERC of Canada and the Go-Bell Scholarships Program. Diego Arroyuelo and Gonzalo Navarro were partially funded by Fondecyt Grant 1-110066, Chile. Gonzalo Navarro was partially funded by Millennium Institute for Cell Dynamics and Biotechnology (ICDB), Grant ICM P05-001-F, Mideplan, Chile. Veli Mäkinen and Jouni Sirén were partially funded by the Academy of Finland under grant number 1140727. Niko Välimäki was partially funded by the Helsinki Graduate School in Computer Science and Engineering.

## REFERENCES

1. XML Mind products. Qizx XML query engine, 2007. Available from: <http://www.axyana.com/qizx> [last accessed 30 August 2013].
2. Boncz PA, Grust T, van Keulen M, Manegold S, Rittinger J, Teubner J. MonetDB/XQuery: a fast XQuery processor powered by a relational engine. *SIGMOD*, Chicago, Illinois, USA, 2006; 479–490.
3. Signum. Tauro, 2008. Available from: <http://tauro.signum.sns.it/> [last accessed 30 August 2013].
4. Kay M. Ten reasons why Saxon XQuery is fast. *IEEE Data Engineering Bulletin* 2008; **31**(4):65–74.
5. Fernández MF, Siméon J, Choi B, Marian A, Sur G. Implementing XQuery 1.0: the Galax experience. *VLDB*, Berlin, Germany, 2003; 1077–1080.
6. Koch C, Scherzinger S, Schmidt M. The GCX system: dynamic buffer minimization in streaming xquery evaluation. *VLDB*, Vienna, Austria, 2007; 1378–1381.
7. Olteanu D. SPEX: streamed and progressive evaluation of XPath. *IEEE Transactions on Knowledge and Data Engineering* 2007; **19**(7):934–949.
8. Golynski A, Munro I, Rao S. Rank/select operations on large alphabets: a tool for text indexing. *SODA*, Miami, Florida, USA, 2006; 368–373.
9. Ferragina P, Manzini G, Mäkinen V, Navarro G. Compressed representations of sequences and full-text indexes. *ACM Transactions on Algorithms* 2007; **3**(2): article 20.
10. Navarro G, Mäkinen V. Compressed full-text indexes. *ACM Computing Surveys* 2007; **39**(1): article 2.
11. Jacobson G. Space-efficient static trees and graphs. *FOCS*, Research Triangle Park, North Carolina, USA, 1989; 549–554.
12. Munro J, Raman V. Succinct representation of balanced parentheses and static trees. *SIAM Journal on Computing* 2001; **31**:762–776.
13. Geary R, Rahman N, Raman R, Raman V. A simple optimal representation for balanced parentheses. *CPM*, Istanbul, Turkey, 2004; 159–172.
14. Benoit D, Demaine E, Munro JI, Raman R, Raman V, Rao SS. Representing trees of higher degree. *Algorithmica* 2005; **43**(4):275–292.
15. Arroyuelo D. An improved succinct representation for dynamic  $k$ -ary trees. *CPM*, Pisa, Italy, 2008; 277–289.
16. Sadakane K, Navarro G. Fully-functional static and dynamic succinct trees. *SODA*, Austin, Texas, USA, 2010; 134–149.
17. Ferragina P, Luccio F, Manzini G, Muthukrishnan S. Structuring labeled trees for optimal succinctness, and beyond. *FOCS*, Pittsburgh, PA, USA, 2005; 184–196.
18. Ferragina P, Luccio F, Manzini G, Muthukrishnan S. Compressing and searching XML data via two zips. *WWW*, Edinburgh, Scotland, 2006; 751–760.
19. Maneth S, Sebastian T. Fast and tiny structural self-indexes for XML. *CoRR* 2010. [abs/1012.5696](http://arxiv.org/abs/1012.5696).
20. Böttcher S, Hartel R, Heinzemann C. BSBC: towards a succinct data format for XML streams. *WEBIST*, Vol. 1, Funchal, Madeira, Portugal, 2008; 13–21.
21. Gottlob G, Koch C, Pichler R. Efficient algorithms for processing XPath queries. *ACM Transactions on Database Systems* 2005; **30**(2):444–491.
22. Maneth S, Nguyen K. XPath whole query optimization. *PVLDB* 2010; **3**(1):882–893.
23. Franceschet M. XPathMark: an XPath benchmark for the XMark generated data. *XSym*, Trondheim, Norway, 2005; 129–143.
24. Schmidt A, Waas F, Kersten ML, Carey MJ, Manolescu I, Busse R. XMark: a benchmark for XML data management. *VLDB*, Hong Kong, China, 2002; 974–985.
25. Arroyuelo D, Claude F, Maneth S, Mäkinen V, Navarro G, Nguyen K, Sirén J, Välimäki N. Fast in-memory XPath search using compressed indexes. *ICDE*, Long Beach, California, USA, 2010; 417–428.
26. Manzini G. An analysis of the Burrows–Wheeler transform. *Journal of the ACM* 2001; **48**(3):407–430.

27. Mäkinen V, Navarro G. Dynamic entropy-compressed sequences and full-text indexes. *ACM Transactions on Algorithms* 2008; **4**(3): article 32.
28. Grossi R, Gupta A, Vitter JS. High-order entropy-compressed text indexes. *SODA*, Baltimore, Maryland, USA, 2003; 841–850.
29. Barbay J, Gagie T, Navarro G, Nekrich Y. Alphabet partitioning for compressed rank/select and applications. *ISAAC* (2), Jeju Island, Korea, 2010; 315–326.
30. Ferragina P, Manzini G. Indexing compressed text. *Journal of the ACM* 2005; **54**(4):552–581.
31. Burrows M, Wheeler DJ. A block-sorting lossless data compression algorithm. *Technical Report 124*, Digital Equipment Corporation, 1994.
32. Mäkinen V, Navarro G. Implicit compression boosting with applications to self-indexing. *SPIRE*, Santiago, Chile, 2007; 229–241.
33. Claude F, Navarro G. Practical rank/select queries over arbitrary sequences. *SPIRE*, Melbourne, Australia, 2008; 176–187.
34. Raman R, Raman V, Rao SS. Succinct indexable dictionaries with applications to encoding  $k$ -ary trees and multisets. *SODA*, San Francisco, CA, USA, 2002; 233–242.
35. Mäkinen V, Navarro G. Rank and select revisited and extended. *Theory of Computing Systems* 2007; **387**(3):332–347.
36. XQuery and XPath Full Text 1.0. Available from: <http://www.w3.org/TR/xpath-full-text-10> [last accessed 30 August 2013].
37. Lam TW, Sung WK, Tam SL, Wong CK, Yiu SM. Compressed indexing and local alignment of DNA. *Bioinformatics* 2008; **24**(6):791–797.
38. Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology* 2009; **10**(3): R25.
39. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* 2009; **25**(14):1754–1760.
40. Sirén J. Compressed suffix arrays for massive data. *SPIRE*, Saariselkä, Finland, 2009; 63–74.
41. Arroyuelo D, Navarro G, Sadakane K. Reducing the space requirement of LZ-index. *CPM*, Barcelona, Spain, 2006; 319–330.
42. Munro I, Raman V. Succinct representation of balanced parentheses, static trees and planar graphs. *FOCS*, Miami Beach, Florida, USA, 1997; 118–126.
43. Arroyuelo D, Cánovas R, Navarro G, Sadakane K. Succinct trees in practice. *ALLENEX*, Austin, Texas, USA, 2010; 84–97.
44. Okanohara D, Sadakane K. Practical entropy-compressed rank/select dictionary. *ALLENEX*, New Orleans, Louisiana, USA, 2007; 60–70.
45. Comon H, Dauchet M, Gilleron R, Jacquemard F, Löding C, Lugiez D, Tison S, Tommasi M. Tree automata techniques and applications, 2007. Available from: <http://www.grappa.univ-lille3.fr/tata> [last accessed 30 August 2013].
46. Hosoya H. *Foundations of XML Processing: The Tree Automata Approach*. Cambridge University Press: Cambridge, UK, 2010.
47. Neven F. Automata theory for XML researchers. *SIGMOD Record* 2002; **31**:39–46.
48. Schwentick T. Automata for XML - a survey. *Journal of Computer and System Sciences* 2007; **73**:289–315.
49. Genevès P, Layaïda N. XML reasoning made practical. *ICDE*, Long Beach, California, USA, 2010; 1169–1172.
50. Libkin L, Sirangelo C. Reasoning about XML with temporal logics and automata. *Journal of Applied Logic* 2010; **8**:210–232.
51. Green TJ, Gupta A, Miklau G, Onizuka M, Suciu D. Processing XML streams with deterministic automata and stream indexes. *ACM Transactions on Database Systems* 2004; **29**:752–788.
52. Koch C. Efficient processing of expressive node-selecting queries on XML data in secondary storage: a tree automata-based approach. *VLDB*, Berlin, Germany, 2003; 249–260.
53. Björklund H, Gelade W, Marquardt M, Martens W. Incremental XPath evaluation. *ICDT*, Saint-Petersburg, Russia, 2009; 162–173.
54. Grust T, van Keulen M, Teubner J. Staircase join: teach a relational DBMS to watch its (axis) steps. *VLDB*, Berlin, Germany, 2003; 524–525.
55. McHugh J, Widom J. Query optimization for xml. *VLDB*, Edinburgh, Scotland, 1999; 315–326.
56. Aho AV, Sethi R, Ullman JD. *Compilers: Principles, Techniques and Tools*. Addison Wesley: Boston, USA, 1986.
57. Conchon S, Filliâtre J-C. Type-safe modular hash-consing. *Proceedings of the ACM SIGPLAN Workshop on ML*, Portland, Oregon, USA, 2006. Available from: <http://www.lri.fr/~filliatr/ftp/publis/hash-consing2.ps> [last accessed 30 August 2013].
58. Franceschet M. XPathMark: functional and performance tests for XPath. *Xquery Implementation Paradigms*, 2007. Available from: <http://drops.dagstuhl.de/opus/volltexte/2007/892> [last accessed 30 August 2013].
59. Fariña A, Brisaboa N, Navarro G, Claude F, Places A, Rodríguez E. Word-based self-indexes for natural language text. *ACM Transactions on Information Systems* 2012; **30**(1): article 1.
60. English Wiktionary. Available from: <http://en.wiktionary.org/> [last accessed 30 August 2013].
61. Mäkinen V, Navarro G, Siren J, Välimäki N. Storage and retrieval of highly repetitive sequence collections. *Journal of Computational Biology* 2010; **17**(3):281–308.
62. JASPAR database. Available from: <http://jaspar.genereg.net> [last accessed 30 August 2013].

63. Chan H-L, Hon W-K, Lam T-W, Sadakane K. Compressed indexes for dynamic text collections. *ACM Transactions on Algorithms* 2007; **3**(2): article 21.
64. Böttcher S, Bültmann A, Hartel R. Search and modification in compressed texts. *DCC*, Snowbird, Utah, US, 2011; 403–412.
65. Navarro G, Nekrich Y. Optimal dynamic sequence representations. *SODA*, New Orleans, Louisiana, USA, 2013; 865–876.
66. Benedikt M, Koch C. XPath leashed. *ACM Computing Surveys* 2008; **41**(1): article 3.
67. Bojanczyk M, Parys P. XPath evaluation in linear time. *Journal of the ACM* 2011; **58**(4):17.
68. Bonifati A, Leighton G, Makinen V, Maneth S, Navarro G, Pugliese A. An in-memory XQuery/XPath engine over a compressed structured text representation. *Structure-Based Compression of Complex Massive Data. Dagstuhl Seminar Proceedings 08261*, Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany, 2008.