

# Efficient Algorithms for Processing XPath Queries\*

Georg Gottlob, Christoph Koch, and Reinhard Pichler

Database and Artificial Intelligence Group  
Technische Universität Wien, A-1040 Vienna, Austria  
{gottlob, koch}@dbai.tuwien.ac.at, reini@logic.at

## Abstract

Our experimental analysis of several popular XPath processors reveals a striking fact: Query evaluation in each of the systems requires time exponential in the size of queries in the worst case. We show that XPath can be processed much more efficiently, and propose main-memory algorithms for this problem with polynomial-time combined query evaluation complexity. Moreover, we present two fragments of XPath for which linear-time query processing algorithms exist.

## 1 Introduction

XPath has been proposed by the W3C [17] as a practical language for selecting nodes from XML document trees. The importance of XPath stems from (1) its potential application as an XML query language per se and it being at the core of several other XML-related technologies, such as XSLT, XPointer, and XQuery and (2) the great and well-deserved interest such technologies receive [1]. Since XPath and related technologies will be tested in ever-growing deployment scenarios, its implementations need to scale well *both* with respect to the size of the XML data and the growing size and intricacy of the queries (usually referred to as *combined complexity*).

Recently, there has been some work on related problems such as query containment for XPath [6, 11, 16], XPath axis rewriting to deal with streaming XML data [4], the expressiveness and complexity of various fragments of XSLT [2, 12], and contributions towards a

\* This work was supported by the Austrian Science Fund (FWF) under project No. Z29-INF. All methods and algorithms presented in this paper are covered by a pending patent. Further resources, updates, and possible corrections will be made available at <http://www.xmltaskforce.com>.

*Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.*

formal semantics definition of XPath [7, 14]. However, to the best of our knowledge, no research results on good or even reasonable methods for processing XPath have been published which may serve as yardsticks for new algorithms.

## Contributions

In this paper, we show that it is possible to noticeably improve the efficiency of existing and future XPath engines. We claim that current implementations of XPath processors do not live up to their potential. The way XPath is defined in [17] motivates an implementation approach that leads to highly inefficient (exponential-time) XPath processing, and many implementations seem to have naively followed this intuition. Likewise, the semantics of a fragment of XPath defined in [14], which uses a fully functional formalism, motivates an exponential-time algorithm.

To get a better understanding of the state-of-the-art of XPath implementations, we experiment with three existing XPath processors, namely XALAN, XT, and Microsoft Internet Explorer 6 (IE6). XALAN [19] is a framework for processing XPath and XSLT which is freely available from the Apache foundation. XT [5] is a freely available XSLT<sup>1</sup> processor written by James Clark. IE6 is a commercial Web browser which supports the formatting of XML documents using XSL. Our experiments show that the time consumption of all three systems grows exponentially in the size of XPath queries in general. This exponentiality is a very practical problem. Of course, queries tend to be short, but we will argue that meaningful practical queries are *not short enough* to allow the existing systems to handle them.

The main contributions of this paper, apart from our experiments, are the following:

- We define a formal bottom-up semantics of XPath (i.e., for the full language as proposed in [17]), which leads to a bottom-up main-memory XPath processing algorithm that runs in low-degree polynomial time in terms of the data and of the query size in the worst case. By a *bottom-up algorithm* we mean

<sup>1</sup>Of course, XSLT allows to embed and execute arbitrary XPath queries.

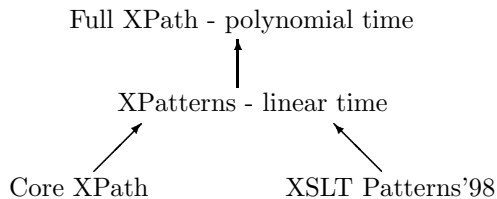


Figure 1: XPath fragments considered in this paper.

a method of processing XPath while traversing the *parse tree* of the query from its leaves up to its root.

- We discuss a general mechanism for translating our bottom-up algorithm into a top-down one. (“Top-down” again relates to the parse tree of the query.) Both have the same worst-case bound on running times but the latter may compute fewer useless intermediate results than the bottom-up algorithm.
- We present a linear-time algorithm (in both data and query size) for a practically useful fragment of XPath, which we will call *Core XPath* in the sequel.

In the experiments presented in this paper, we show that evaluating such queries in XALAN and XT already takes exponential time in the size of the queries in the worst case. The processing time of IE6 for this fragment grows polynomially in the size of queries, but requires quadratic time in the size of the XML *data* (when the query is fixed).

- We discuss the now superseded language of *XSLT Patterns* of the XSLT draft of December 16th, 1998 [18]. Since then, full XPath has been adopted as the XSLT Pattern language. This language remains interesting, as it shares many features with XPath and is a useful practical query language. We extend this language with all of the XPath axes and call it *XPatterns* to keep it short. Surprisingly, XPatterns queries can be evaluated very efficiently, in linear time in the size of the data and the query.

The rationale for presenting these fragments is their relevance to the efficiency of engines for full XPath on common queries. An overview of the various query language fragments considered in this paper and data complexity bounds of the associated algorithms is given in Figure 1. By  $\mathcal{L}_1 \leftarrow \mathcal{L}_2$ , we denote that language  $\mathcal{L}_1$  subsumes language  $\mathcal{L}_2$ : XPatterns fully subsumes the Core XPath language, and subsumes XSLT Patterns’98 (except for a minor detail). XPatterns is a fragment of XPath.

## Structure

The structure of this paper is as follows. In Section 2, we provide experimental results for existing XPath processors. Section 3 presents basic notions, including the data model and auxiliary functions. Section 4 introduces XPath axes. Section 5 defines the semantics of XPath in a concise way. Section 6 houses the bottom-up semantics definition and algorithm for full XPath, and Section 7 comes up with the modifications to obtain a top-down algorithm. Section 8

presents linear-time fragments of XPath (Core XPath and XPatterns). We conclude with Section 9.

## 2 State-of-the-Art of XPath Systems

In this section, we evaluate the efficiency of three XPath engines, namely Apache XALAN (the Lotus/IBM XPath implementation which has been donated to the Apache foundation) and James Clark’s XT, which are, as we believe, the two most popular freely available XPath engines, and Microsoft Internet Explorer 6 (IE6), a commercial product. The reader is assumed familiar with XPath and standard notions such as *axes* and *location steps* (cf. [17]).

The version of XALAN used for the experiments was Xalan-j\_2\_2\_D11 (i.e., a Java release). We used the current version of XT (another Java implementation) with release tag 19991105, as available on James Clark’s home page, in combination with his XP parser through the SAX driver. We ran both XALAN and XT on a 360 MHz (dual processor) Ultra Sparc 60 with 512 MB of RAM running Solaris. IE6 was evaluated on a Windows 2000 machine with a 1.2 GHz AMD K7 processor and 1.5 GB of RAM.

XT and IE6 are not literally XPath engines, but are able to process XPath embedded in XSLT transformations. We used the `xsl:foreach` performative to obtain the set of all nodes an XPath query would evaluate to.

We show by experiments that all three implementations require time exponential in the size of the queries in the worst case. Furthermore, we show that even the simplest queries, with which IE6 can deal efficiently in the size of the queries, take quadratic time in the size of the data. Since we used two different platforms for running the benchmarks, our goal of course was not to compare the systems against each other, but to test the scalabilities of their XPath processing algorithms. The reason we used two different platforms was that Solaris allows for accurate timing, while IE6 is only available on Windows. (The IE6 timings reported on here have the precision of  $\pm 1$  second).

The XML documents we used were of very simple structure. *The* document of size  $n$  was of the form

$$\langle a \rangle \underbrace{\langle b \rangle \dots \langle b \rangle}_{n \text{ times}} \langle /a \rangle$$

and its tree thus contained  $n + 1$  nodes.

### Experiment 1: Exponential-time Query Complexity of XALAN and XT

In this experiment, we used the document of size 2 (i.e.,  $\langle a \rangle \langle b \rangle \langle b \rangle \langle /a \rangle$ ). Queries were constructed using a simple pattern. The first query was `‘//a/b’` and the  $i + 1$ -th query was obtained by taking the  $i$ -th query and appending `‘/parent::a/b’`. For instance, the third query was `‘//a/b/parent::a/b/parent::a/b’`.

It is easy to see that the time measurements reported in Figure 2, which uses a log scale Y axis, grow exponentially with the size of the query. The sharp

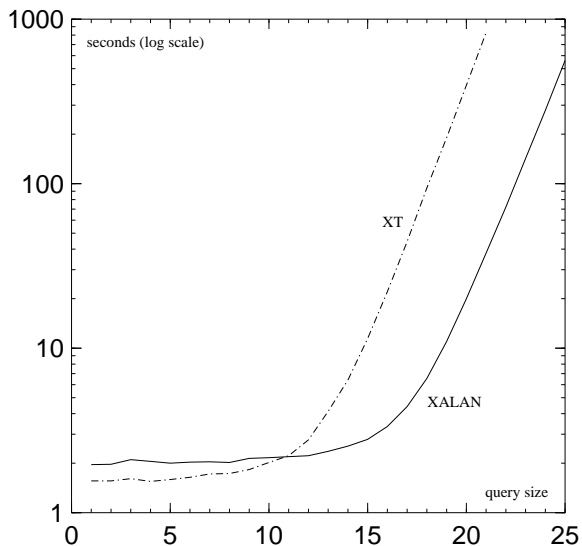


Figure 2: Exponential-time query complexity of XT and XALAN (Experiment 1).

bend in the curves is due to the near-constant runtime overhead of the Java VM and of parsing the XML document.

### Discussion

This behavior can be explained with the following pseudocode fragment, which seems to appropriately describe the basic query evaluation strategy of XALAN and XT.

```

procedure process-location-step( $n_0, Q$ )
/*  $n_0$  is the context node;
   query  $Q$  is a list of location steps */
begin
  node set  $S := \text{apply } Q.\text{first}$  to node  $n_0$ ;
  if ( $Q.\text{tail}$  is not empty) then
    for each node  $n \in S$  do
      process-location-step( $n, Q.\text{tail}$ );
end

```

It is clear that each application of a location step to a context node may result in a set of nodes of size linear in the size of the document (e.g., each node may have a linear number of descendants or nodes appearing after it in the document). If we now proceed by recursively applying the location steps of an XPath query to individual nodes as shown in the pseudocode procedure above, we end up consuming time exponential in the size of the query in the worst case, even for very simple path queries. As a (simplified) recurrence, we have

$$\text{Time}(|Q|) := \begin{cases} |D| * \text{Time}(|Q| - 1) & \dots & |Q| > 0 \\ 1 & \dots & |Q| = 0 \end{cases}$$

where  $|Q|$  is the length of the query and  $|D|$  is the size of the document, or  $\text{Time}(|Q|) = |D|^{|Q|}$ .

The class of queries used puts an emphasis on simplicity and reproducibility (using the *very* simple

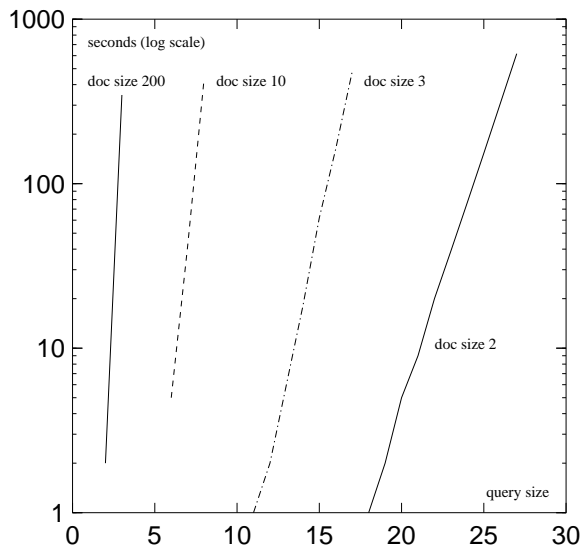


Figure 3: Exponential-time query complexity of IE6, for document sizes 2, 3, 10, and 200 (Experiment 2).

document  $\langle a \rangle \langle b \rangle \langle b \rangle \langle /a \rangle$ ). Interestingly, each ‘parent::a/b’ sequence quite exactly doubles the times both systems take to evaluate a query, as we first jump (back) to the tree root labeled “a” and then experience the “branching factor” of two due the two child nodes labeled “b”.

Our class of queries may seem contrived; however, it is clear that we make a practical point. First, more realistic document sizes allow for very short queries only<sup>2</sup>. At the same time, XPath query engines need to be able to deal with increasingly sophisticated queries, along the current trend to delegate larger and larger parts of data management problems to query engines, where they can profit from their efficiency and can be made subject to optimization. The intuition that XPath can be used to match a large class of *tree patterns* [13, 10, 3] in XML documents also implies to a certain degree that queries may be extensive.

Moreover, similar queries using antagonist axes such as “following” and “preceding” instead of “child” and “parent” do have practical applications, such as when we want to put restrictions on the relative positions of nodes in a document. Finally, if we make the realistic assumption that the documents are always much larger than the queries ( $|Q| \ll |D|$ ), it is not even necessary to jump back and forth with antagonist axes. We can use queries such as `//following::*//following::*/.../following::*` to observe exponential behavior.

### Experiment 2: Exponential-time Query Complexity of Internet Explorer 6

In our second experiment, we executed queries that nest two important features of XPath, namely paths

<sup>2</sup>We will show this in the second experiment for IE6 (see Figure 3), and have verified it for XALAN and XT as well.

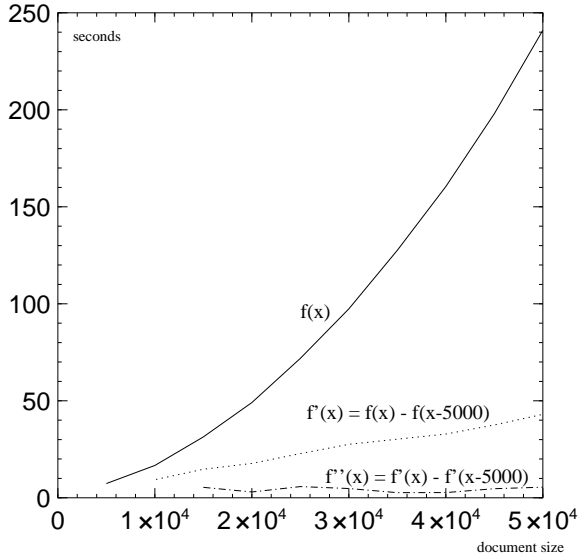


Figure 4: Quadratic-time data complexity of IE6.  $f'$  and  $f''$  are the first and second derivatives, respectively, of our graph of timings  $f$  (Experiment 3).

and arithmetics, using IE6. The first three queries were

```
//a/b[count(parent::a/b) > 1]
//a/b[count(parent::a/b[
  count(parent::a/b) > 1]) > 1]
//a/b[count(parent::a/b[
  count(parent::a/b[
    count(parent::a/b) > 1]) > 1]) > 1]
```

and it is clear how to continue this sequence.

The experiment was carried out for four document sizes (2, 3, 10, and 200). Figure 3 shows clearly that IE6 requires time exponential in the size of the query.

### Experiment 3: Quadratic-time Data Complexity for Simple Path Queries (IE6)

For our third experiment, we took a fixed query and benchmarked the time taken by IE6 for various document sizes. The query was  $'//a' + q(20) + '//b'$  with

$$q(i) := \begin{cases} '//b[ancestor::a' + q(i-1) & \dots & i > 0 \\ '+ '//b]/ancestor::a' & \dots & i = 0 \\ ', & \dots & i = 0 \end{cases}$$

(Note: The size of queries  $q(i)$  is of course  $O(i)$ .)

**Example 2.1** For instance, the query of size two according to this scheme, i.e.  $'//a' + q(2) + '//b'$ , is

```
//a//b[ancestor::a//b[ancestor::a//b
  ]/ancestor::a//b
]/ancestor::a//b □
```

The granularity of measurements (in terms of document size) was 5000 nodes. Figure 4 shows that IE6

takes quadratic time w.r.t. the size of the data already for this simple class of path queries.

The query complexity of IE6 w.r.t. such queries is polynomial as well. Due to space limitations, we do not provide a graph for this experiment.

By virtue of our experiments, the following question naturally arises: Is there an algorithm for processing XPath with guaranteed polynomial-time behavior (combined complexity), or even one that requires only linear time for simple queries? In the remainder of this paper, we are able to provide a positive answer to this.

## 3 Basic Notions

In this paper, we use an XML document model simplified as follows. All of the artifacts of this section are defined in the context of a given XML document. In our data model, an XML document is an unranked, ordered, and labeled tree. Let  $\text{dom}$  be the set of all nodes in this tree, and let us use the two functions

$$\text{firstchild}, \text{nextsibling} : \text{dom} \rightarrow \text{dom},$$

to represent the tree<sup>3</sup>. “firstchild” returns the first child of a node (if there are any children, i.e., the node is not a leaf), and otherwise “null”. Let  $n_1, \dots, n_k$  be the children of some node in-order. Then,  $\text{nextsibling}(n_i) = n_{i+1}$ , i.e., “nextsibling” returns the neighboring node to the right, if it exists, and “null” otherwise (if  $i = k$ ). We define the functions  $\text{firstchild}^{-1}$  and  $\text{nextsibling}^{-1}$  as the inverses of the former two functions, where “null” is returned if no inverse exists for a given node. Where appropriate, we will use binary relations of the same name instead of the functions. ( $\{(x, f(x)) \mid x \in \text{dom}, f(x) \neq \text{null}\}$  is the binary relation for function  $f$ .)

Let  $\Sigma$  be a finite labeling alphabet. We define a function  $T : (\Sigma \cup \{*\}) \rightarrow 2^{\text{dom}}$  (“node test”)<sup>4</sup> which assigns to each label (XML tag) the set of nodes labeled with it; moreover,  $T(*) := \text{dom}$ . Let  $<_{\text{doc}}$  be the document order, where  $x <_{\text{doc}} y$  (for two nodes  $x, y \in \text{dom}$ ) iff the opening tag of  $x$  precedes the opening tag of  $y$  in the (well-formed) document. The function  $\text{first}_{<_{\text{doc}}}$  returns the first node in a set w.r.t. document order.

All nodes are of the same type; thus, we do not distinguish among element, attribute, and processing instruction nodes, among others. For the same reason, we do not discuss the “namespace” and “attribute” axes<sup>5</sup> as well as the “local-name”, “namespace-uri”, and “name” core library functions [17] either. The reason for this is lack of space; however, extending our semantics and algorithms (without a penalty w.r.t. efficiency bounds) is an easy exercise.

<sup>3</sup>Actually, “firstchild” and “nextsibling” are part of the XML Document Object Model (DOM).

<sup>4</sup>Our  $*$  node test coincides with  $\text{node}()$  of [17], while  $*$  of [17] is a node-typed version of  $\text{node}()$ .

<sup>5</sup>To cover these two kinds of axes in our simple node-labeled tree data model, we could for instance assume that attributes are modeled in the document tree as child nodes with special labels.

child := firstchild.nextsibling*
parent := (nextsibling <sup>-1</sup> )*.firstchild <sup>-1</sup>
descendant := firstchild.(firstchild ∪ nextsibling)*
ancestor := (firstchild <sup>-1</sup> ∪ nextsibling <sup>-1</sup> )*.firstchild <sup>-1</sup>
descendant-or-self := descendant ∪ self
ancestor-or-self := ancestor ∪ self
following := ancestor-or-self.nextsibling. nextsibling*.descendant-or-self
preceding := ancestor-or-self.nextsibling <sup>-1</sup> . (nextsibling <sup>-1</sup> )*.descendant-or-self
following-sibling := nextsibling.nextsibling*
preceding-sibling := (nextsibling <sup>-1</sup> )*.nextsibling <sup>-1</sup>

Table 1: Axis definitions in terms of “primitive” tree relations “firstchild”, “nextsibling”, and their inverses.

Each node in an XML document may be identified by a unique id. The function `deref_ids : string → 2dom` interprets its input string as a whitespace-separated list of keys and returns the set of nodes whose id’s are contained in that list. The function `strval : dom → string` returns the string value of a node, which is the concatenation of non-tag strings between the node’s start and end tags in the document. The functions `to_string` and `to_number` convert a number to a string resp. a string to a number according to the rules specified in [17].

## 4 XPath Axes

The XPath axes *self*, *child*, *parent*, *descendant*, *ancestor*, *descendant-or-self*, *ancestor-or-self*, *following*, *preceding*, *following-sibling*, and *preceding-sibling* are binary relations  $\chi \subseteq \text{dom} \times \text{dom}$ . Let  $\text{self} := \{\langle x, x \rangle \mid x \in \text{dom}\}$ . The other axes are defined in terms of our “primitive” relations “firstchild” and “nextsibling” as shown in Table 1 (cf. [17]).  $R_1.R_2$ ,  $R_1 \cup R_2$ , and  $R_1^*$  denote the concatenation, union, and reflexive and transitive closure, respectively, of binary relations  $R_1$  and  $R_2$ . Let  $E(\chi)$  denote the regular expression defining  $\chi$  in Table 1. It is important to observe that some axes are defined in terms of other axes, but that these definitions are acyclic / non-recursive.

**Definition 4.1** Let  $\chi$  denote an XPath axis relation. We define the function  $\chi : 2^{\text{dom}} \rightarrow 2^{\text{dom}}$  as  $\chi(X_0) = \{x \mid \exists x_0 \in X_0 : x_0\chi x\}$  (and thus overload the relation name  $\chi$ ), where  $X_0 \subseteq \text{dom}$  is a set of nodes.  $\square$

**Algorithm 4.2** (Axis Evaluation)

**Input:** A set of nodes  $S$  and an axis  $\chi$

**Output:**  $\chi(S)$

**Method:**  $\text{eval}_\chi(S)$

```

function  $\text{eval}_{(R_1 \cup \dots \cup R_n)^*}(S)$  begin
   $S' := S$ ; /*  $S'$  is represented as a list */
  while there is a next element  $x$  in  $S'$  do
    append  $\{R_i(x) \mid 1 \leq i \leq n, R_i(x) \neq \text{null},$ 
       $R_i(x) \notin S'\}$  to  $S'$ ;
  return  $S'$ ;
end;

```

**function**  $\text{eval}_\chi(S) := \text{eval}_{E(\chi)}(S)$ .

**function**  $\text{eval}_{\text{self}}(S) := S$ .

**function**  $\text{eval}_{e_1.e_2}(S) := \text{eval}_{e_2}(\text{eval}_{e_1}(S))$ .

**function**  $\text{eval}_R(S) := \{R(x) \mid x \in S\}$ .

**function**  $\text{eval}_{\chi_1 \cup \chi_2}(S) := \text{eval}_{\chi_1}(S) \cup \text{eval}_{\chi_2}(S)$ .

where  $S \subseteq \text{dom}$  is a set of nodes of an XML document,  $e_1$  and  $e_2$  are regular expressions,  $R, R_1, \dots, R_n$  are primitive relations,  $\chi_1$  and  $\chi_2$  are axes, and  $\chi$  is an axis other than “self”.  $\square$

Clearly, some axes could have been defined in a simpler way in Table 1 (e.g., ancestor equals `parent.parent*`). However, the definitions, which use a limited form of regular expressions only, allow to compute  $\chi(S)$  in a very simple way, as evidenced by Algorithm 4.2.

The function  $\text{eval}_{(R_1 \cup \dots \cup R_n)^*}$  essentially computes graph reachability (not transitive closure). It can be implemented to run in linear time in terms of the data in a straightforward manner; (non)membership in  $S'$  is checked in constant time using a direct-access version of  $S'$  maintained in parallel to its list representation (naively, this could be an array of bits, one for each member of  $\text{dom}$ , telling which nodes are in  $S'$ ).

**Lemma 4.3** Let  $S \subseteq \text{dom}$  be a set of nodes of an XML document and  $\chi$  be an axis. Then, (1)  $\chi(S) = \text{eval}_\chi(S)$  and (2) Algorithm 4.2 runs in time  $O(|\text{dom}|)$ .

**Proof Sketch** ( $O(|\text{dom}|)$  running time) The time bound is due to the fact that each of the eval functions can be implemented so as to visit each node at most once and the number of calls to eval functions and relations joined by union is constant (see Table 1).  $\square$

**Definition 4.4** (Document order relative to an axis)

We define the relation  $<_{\text{doc}, \chi}$  relative to the axis  $\chi$  as follows. For  $\chi \in \{\text{self}, \text{child}, \text{descendant}, \text{descendant-or-self}, \text{following-sibling}, \text{following}\}$ ,  $<_{\text{doc}, \chi}$  is the standard document order  $<_{\text{doc}}$ . For the remaining axes, it is the reverse document order.

Moreover, given a node  $x$  and a set of nodes  $S$  with  $x \in S$ , let  $\text{idx}_\chi(x, S)$  be the index of  $x$  in  $S$  w.r.t.  $<_{\text{doc}, \chi}$  (where 1 is the smallest index).  $\square$

## 5 Semantics of XPath

In this section, we present a concise definition of the semantics of XPath 1 [17]. We assume the syntax of this language known, and cohere with its *unabbreviated* form [17]. We use a normal form syntax of XPath, which is obtained by the following rewrite rules, applied initially:

1. Location steps  $\chi::t[e]$ , where  $e$  is an expression that produces a number (see below), are replaced by the equivalent expression  $\chi::t[e = \text{position}()]$ .
2. All type conversions are made explicit (using the conversion functions `string`, `number`, and `boolean`, which we will define below).

$P[\chi::t[e_1] \cdots [e_m]](x) :=$   
**begin**  
 $S := \{y \mid x\chi y, y \in T(t)\};$   
**for**  $1 \leq i \leq m$  (in ascending order) **do**  
 $S := \{y \in S \mid \llbracket e_i \rrbracket(y, \text{id}_{X_\chi}(y, S), |S|) = \text{true}\};$   
**return**  $S;$   
**end;**  
 $P[\pi_1/\pi_2](x) := P[\pi_1](x) \cup P[\pi_2](x)$   
 $P[\pi](x) := P[\pi](\text{root})$   
 $P[\pi_1/\pi_2](x) := \bigcup_{y \in P[\pi_1](x)} P[\pi_2](y)$

Figure 5: Standard semantics of location paths.

- Each variable is replaced by the (constant) value of the input variable binding.

The main structural feature of XPath are *expressions*, which are of one of four types, namely *node set*, *number*, *string*, or *boolean*. Each expression evaluates relative to a context  $\vec{c} = \langle x, k, n \rangle$  consisting of a *context node*  $x$ , a *context position*  $k$ , and a *context size*  $n$  [17].  $\mathbf{C} = \text{dom} \times \{\langle k, n \rangle \mid 1 \leq k \leq n \leq |\text{dom}|\}$  is the domain of contexts. Let  $ArithOp \in \{+, -, *, \text{div}, \text{mod}\}$ ,  $RelOp \in \{=, \neq, \leq, <, \geq, >\}$ ,  $EqOp \in \{=, \neq\}$ , and  $GtOp \in \{\leq, <, \geq, >\}$ . By slight abuse of notation, we identify these arithmetic and relational operations with their symbols in the remainder of this paper. However, it should be clear whether we refer to the operation or its symbol at any point. By  $\pi, \pi_1, \pi_2, \dots$  we denote location paths.

**Definition 5.1** (*Semantics of XPath*) Each XPath expression returns a value of one of the following four types: number, node set, string, and boolean (abbreviated num, nset, str, and bool, respectively). Let  $\mathcal{T}$  be an expression type and the semantics  $\llbracket e \rrbracket : \mathbf{C} \rightarrow \mathcal{T}$  of XPath expression  $e$  be defined as follows.

$$\begin{aligned}
\llbracket \pi \rrbracket(\langle x, k, n \rangle) &:= P[\pi](x) \\
\llbracket \text{position}() \rrbracket(\langle x, k, n \rangle) &:= k \\
\llbracket \text{last}() \rrbracket(\langle x, k, n \rangle) &:= n
\end{aligned}$$

For all other kinds of expressions  $e = Op(e_1, \dots, e_m)$  mapping a context  $\vec{c}$  to a value of type  $\mathcal{T}$ ,  $\llbracket Op(e_1, \dots, e_m) \rrbracket(\vec{c}) := \mathcal{F}[Op](\llbracket e_1 \rrbracket(\vec{c}), \dots, \llbracket e_m \rrbracket(\vec{c}))$ , where  $\mathcal{F}[Op] : \mathcal{T}_1 \times \dots \times \mathcal{T}_m \rightarrow \mathcal{T}$  is called the *effective semantics function* of  $Op$ . The function  $P$  is defined in Figure 5 and the effective semantics function  $\mathcal{F}$  is defined in Table 2.  $\square$

To save space, we at times re-use function definitions in Table 2 to define others. However, our definitions are not circular and the indirections can be eliminated by a constant number of unfolding steps. Moreover, for lack of space, we define neither the number operations floor, ceiling, and round nor the string operations concat, starts-with, contains, substring-before, substring-after, substring (two versions), string-length, normalize-space, translate, and lang in Table 2, but it is very easy to obtain these definitions from the XPath 1 Recommendation [17].

Expr. $E$ : Operator Signature	Semantics $\mathcal{F}[E]$
$\mathcal{F}[\text{constant number } v : \rightarrow \text{num}]()$	$v$
$\mathcal{F}[ArithOp : \text{num} \times \text{num} \rightarrow \text{num}](v_1, v_2)$	$v_1 \text{ ArithOp } v_2$
$\mathcal{F}[\text{count} : \text{nset} \rightarrow \text{num}](S)$	$ S $
$\mathcal{F}[\text{sum} : \text{nset} \rightarrow \text{num}](S)$	$\sum_{n \in S} \text{to\_number}(\text{strval}(n))$
$\mathcal{F}[\text{id} : \text{nset} \rightarrow \text{nset}](S)$	$\bigcup_{n \in S} \text{deref\_ids}(\text{strval}(n))$
$\mathcal{F}[\text{id} : \text{str} \rightarrow \text{nset}](s)$	$\text{deref\_ids}(s)$
$\mathcal{F}[\text{constant string } s : \rightarrow \text{str}]()$	$s$
$\mathcal{F}[\text{and} : \text{bool} \times \text{bool} \rightarrow \text{bool}](b_1, b_2)$	$b_1 \wedge b_2$
$\mathcal{F}[\text{or} : \text{bool} \times \text{bool} \rightarrow \text{bool}](b_1, b_2)$	$b_1 \vee b_2$
$\mathcal{F}[\text{not} : \text{bool} \rightarrow \text{bool}](b)$	$\neg b$
$\mathcal{F}[\text{true}() : \rightarrow \text{bool}]()$	true
$\mathcal{F}[\text{false}() : \rightarrow \text{bool}]()$	false
$\mathcal{F}[RelOp : \text{nset} \times \text{nset} \rightarrow \text{bool}](S_1, S_2)$	$\exists n_1 \in S_1, n_2 \in S_2 : \text{strval}(n_1) \text{ RelOp } \text{strval}(n_2)$
$\mathcal{F}[RelOp : \text{nset} \times \text{num} \rightarrow \text{bool}](S, v)$	$\exists n \in S : \text{to\_number}(\text{strval}(n)) \text{ RelOp } v$
$\mathcal{F}[RelOp : \text{nset} \times \text{str} \rightarrow \text{bool}](S, s)$	$\exists n \in S : \text{strval}(n) \text{ RelOp } s$
$\mathcal{F}[RelOp : \text{nset} \times \text{bool} \rightarrow \text{bool}](S, b)$	$\mathcal{F}[\text{boolean}](S) \text{ RelOp } b$
$\mathcal{F}[EqOp : \text{bool} \times (\text{str} \cup \text{num} \cup \text{bool}) \rightarrow \text{bool}](b, x)$	$b \text{ EqOp } \mathcal{F}[\text{boolean}](x)$
$\mathcal{F}[EqOp : \text{num} \times (\text{str} \cup \text{num}) \rightarrow \text{bool}](v, x)$	$v \text{ EqOp } \mathcal{F}[\text{number}](x)$
$\mathcal{F}[EqOp : \text{str} \times \text{str} \rightarrow \text{bool}](s_1, s_2)$	$s_1 \text{ EqOp } s_2$
$\mathcal{F}[GtOp : (\text{str} \cup \text{num} \cup \text{bool}) \times (\text{str} \cup \text{num} \cup \text{bool}) \rightarrow \text{bool}](x_1, x_2)$	$\mathcal{F}[\text{number}](x_1) \text{ GtOp } \mathcal{F}[\text{number}](x_2)$
$\mathcal{F}[\text{string} : \text{num} \rightarrow \text{str}](v)$	$\text{to\_string}(v)$
$\mathcal{F}[\text{string} : \text{nset} \rightarrow \text{str}](S)$	<b>if</b> $S = \emptyset$ <b>then</b> "" <b>else</b> $\text{strval}(\text{first}_{<doc}(S))$
$\mathcal{F}[\text{string} : \text{bool} \rightarrow \text{str}](b)$	<b>if</b> $b = \text{true}$ <b>then</b> "true" <b>else</b> "false"
$\mathcal{F}[\text{boolean} : \text{str} \rightarrow \text{bool}](s)$	<b>if</b> $s \neq \text{" "}$ <b>then</b> true <b>else</b> false
$\mathcal{F}[\text{boolean} : \text{num} \rightarrow \text{bool}](v)$	<b>if</b> $v \neq \pm 0$ <b>and</b> $v \neq \text{NaN}$ <b>then</b> true <b>else</b> false
$\mathcal{F}[\text{boolean} : \text{nset} \rightarrow \text{bool}](S)$	<b>if</b> $S \neq \emptyset$ <b>then</b> true <b>else</b> false
$\mathcal{F}[\text{number} : \text{str} \rightarrow \text{num}](s)$	$\text{to\_number}(s)$
$\mathcal{F}[\text{number} : \text{bool} \rightarrow \text{num}](b)$	<b>if</b> $b = \text{true}$ <b>then</b> 1 <b>else</b> 0
$\mathcal{F}[\text{number} : \text{nset} \rightarrow \text{num}](S)$	$\mathcal{F}[\text{number}](\mathcal{F}[\text{string}](S))$

Table 2: XPath effective semantics functions.

Expression Type	Associated Relation $R$
num	$R \subseteq \mathbf{C} \times \mathbb{R}$
bool	$R \subseteq \mathbf{C} \times \{\text{true}, \text{false}\}$
nset	$R \subseteq \mathbf{C} \times 2^{\text{dom}}$
str	$R \subseteq \mathbf{C} \times \text{char}^*$

Table 3: Expression types and associated relations.

The compatibility of our semantics definition (modulo the assumptions made in this paper to simplify the data model) with [17] can easily be verified by inspection of the latter document.

It is instructive to compare the definition of  $P[\pi_1/\pi_2]$  in Figure 5 with the procedure process-location-step of Section 2 and the claim regarding exponential-time query evaluation made there. In fact, if the semantics definition of [17] (or of this section, for that matter) is followed rigorously to obtain an analogous functional implementation, query evaluation using this implementation requires time exponential in the size of the queries.

## 6 Bottom-up Evaluation of XPath

In this section, we present a bottom-up semantics and algorithm for evaluating XPath queries in polynomial time. We discuss the intuitions which lead to polynomial time evaluation (which we call the “context-value table principle”), and establish the correctness and complexity results.

**Definition 6.1** (*Semantics*) We represent the four XPath expression types nset, num, str, and bool using relations as shown in Table 3. The bottom-up semantics of expressions is defined via a semantics function

$$\mathcal{E}_\uparrow : \text{Expression} \rightarrow \text{nset} \cup \text{num} \cup \text{str} \cup \text{bool},$$

given in Table 4 and as

$$\mathcal{E}_\uparrow[\text{Op}(e_1, \dots, e_m)] := \{ \langle \vec{c}, \mathcal{F}[\text{Op}](v_1, \dots, v_m) \rangle \mid \vec{c} \in \mathbf{C}, \langle \vec{c}, v_1 \rangle \in \mathcal{E}_\uparrow[e_1], \dots, \langle \vec{c}, v_m \rangle \in \mathcal{E}_\uparrow[e_m] \}$$

for the remaining kinds of XPath expressions.  $\square$

Now, for each expression  $e$  and each  $\langle x, k, n \rangle \in \mathbf{C}$ , there is exactly one  $v$  s.t.  $\langle x, k, n, v \rangle \in \mathcal{E}_\uparrow[e]$ .

**Theorem 6.2** *Let  $e$  be an arbitrary XPath expression. Then, for context node  $x$ , position  $k$ , and size  $n$ , the value of  $e$  is  $v$ , where  $v$  is the unique value such that  $\langle x, k, n, v \rangle \in \mathcal{E}_\uparrow[e]$ .*

The main principle that we propose at this point to obtain an XPath evaluation algorithm with polynomial-time complexity is the notion of a *context-value table* (i.e., a relation for each expression, as discussed above).

Expr. $E$ : Operator Signature Semantics $\mathcal{E}_\uparrow[E]$
location step $\chi::t \rightarrow \text{nset}$ $\{ \langle x_0, k_0, n_0, \{x \mid x_0 \chi x, x \in T(t)\} \rangle \mid \langle x_0, k_0, n_0 \rangle \in \mathbf{C} \}$
location step $E[e]$ over axis $\chi$ : $\text{nset} \times \text{bool} \rightarrow \text{nset}$ $\{ \langle x_0, k_0, n_0, \{x \in S \mid \langle x, \text{id}_{x_\chi}(x, S),  S , \text{true} \rangle \in \mathcal{E}_\uparrow[e] \} \rangle \mid \langle x_0, k_0, n_0, S \rangle \in \mathcal{E}_\uparrow[E] \}$
location path $/\pi$ : $\text{nset} \rightarrow \text{nset}$ $\mathbf{C} \times \{ S \mid \exists k, n : \langle \text{root}, k, n, S \rangle \in \mathcal{E}_\uparrow[\pi] \}$
location path $\pi_1/\pi_2$ : $\text{nset} \times \text{nset} \rightarrow \text{nset}$ $\{ \langle x, k, n, z \rangle \mid 1 \leq k \leq n \leq  \text{dom} , \langle x, k_1, n_1, Y \rangle \in \mathcal{E}_\uparrow[\pi_1], \bigcup_{y \in Y} \langle y, k_2, n_2, z \rangle \in \mathcal{E}_\uparrow[\pi_2] \}$
location path $\pi_1 \mid \pi_2$ : $\text{nset} \times \text{nset} \rightarrow \text{nset}$ $\mathcal{E}_\uparrow[\pi_1] \cup \mathcal{E}_\uparrow[\pi_2]$
position(): $\rightarrow \text{num}$ $\{ \langle x, k, n, k \rangle \mid \langle x, k, n \rangle \in \mathbf{C} \}$
last(): $\rightarrow \text{num}$ $\{ \langle x, k, n, n \rangle \mid \langle x, k, n \rangle \in \mathbf{C} \}$

Table 4: Expression relations for location paths, position(), and last().

**Context-value Table Principle.** Given an expression  $e$  that occurs in the input query, the context-value table of  $e$  specifies all valid combinations of contexts  $\vec{c}$  and values  $v$ , such that  $e$  evaluates to  $v$  in context  $\vec{c}$ . Such a table for expression  $e$  is obtained by first computing the context-value tables of the direct subexpressions of  $e$  and subsequently combining them into the context-value table for  $e$ . Given that the size of each of the context-value tables has a polynomial bound and each of the combination steps can be effected in polynomial time (all of which we can assure in the following), query evaluation in total under our principle also has a polynomial time bound<sup>6</sup>.  $\square$

**Query Evaluation.** The idea of Algorithm 6.3 below is so closely based on our semantics definition that its correctness follows directly from the correctness result of Theorem 6.2.

**Algorithm 6.3** (Bottom-up algorithm for XPath)

**Input:** An XPath query  $Q$ ;

**Output:**  $\mathcal{E}_\uparrow[Q]$ .

**Method:**

let  $\text{Tree}(Q)$  be the parse tree of query  $Q$ ;

$\mathbf{R} := \emptyset$ ;

**for each** atomic expression  $l \in \text{leaves}(\text{Tree}(Q))$  **do**

    compute table  $\mathcal{E}_\uparrow[l]$  and add it to  $\mathbf{R}$ ;

**while**  $\mathcal{E}_\uparrow[\text{root}(\text{Tree}(Q))] \notin \mathbf{R}$  **do**

**begin**

        take an  $\text{Op}(l_1, \dots, l_n) \in \text{nodes}(\text{Tree}(Q))$

        s.t.  $\mathcal{E}_\uparrow[l_1], \dots, \mathcal{E}_\uparrow[l_n] \in \mathbf{R}$ ;

        compute  $\mathcal{E}_\uparrow[\text{Op}(l_1, \dots, l_n)]$  using  $\mathcal{E}_\uparrow[l_1], \dots, \mathcal{E}_\uparrow[l_n]$ ;

        add  $\mathcal{E}_\uparrow[\text{Op}(l_1, \dots, l_n)]$  to  $\mathbf{R}$ ;

**end**;

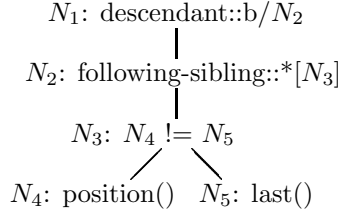
**return**  $\mathcal{E}_\uparrow[\text{root}(\text{Tree}(Q))]$ .  $\square$

<sup>6</sup>The number of expressions to be considered is fixed with the parse tree of a given query.

**Example 6.4** Let  $D$  denote the simple XML document of size 4 from Section 2, i.e., it has one node labeled  $a$  and four child nodes labeled  $b$ . Hence, we write  $\text{dom} = \{a, b_1, \dots, b_4\}$ , where  $a$  denotes the unique node labeled  $a$  and  $b_i$  denotes (in document order) the  $i$ -th node labeled  $b$ . Now suppose that we want to evaluate the XPath query  $Q$ , which reads as

`descendant::b/following-sibling::*[position() != last()]`

over the input context  $\langle a, 1, 1 \rangle$ . We illustrate how this evaluation can be done by the context-value table principle: First of all, we have to set up the parse tree



of  $Q$  with its 5 nodes  $N_1, \dots, N_5$ , corresponding to the 5 subexpressions of  $Q$  ( $N_1$ , the root node of the parse tree, corresponds to  $Q$ ). Then we compute the context-value tables of the leaf nodes  $N_4$  and  $N_5$  in the parse tree. (All context-value tables are shown in Figure 6.) Note that it suffices to compute the result value “ $val$ ” for those combinations  $\langle x, k, n \rangle$  of context-node/size/position which can possibly be generated by the “following-sibling” axis. From the context-value tables at the nodes  $N_4$  and  $N_5$ , it is easy to compute the table at  $N_3$ , and subsequently the tables at  $N_2$  and  $N_1$ . From the context-value table at the root  $N_1$  of the parse tree, we can read out the final result  $\{b_2, b_3\}$ .

In Figure 6, we have made a number of simplifying assumptions to keep the tables short. The table of  $N_3$  contains only those contexts for which the value is “true”; thus, the values are not displayed. Similarly, in the context-value tables at  $N_2$  and  $N_1$ , we have only considered the dependence of the result (node sets have been unnested to obtain a flat table) on the context node  $x$ . In contrast, the context position and size have been omitted since they have no influence on the overall result (they would only blow up the tables).  $\square$

**Theorem 6.5** *XPath can be evaluated bottom-up in polynomial time (combined complexity).*

**Proof Sketch** During the bottom-up computation of a query  $Q$ ,  $O(|Q|)$  relations are created. The size of bool relations is bounded by  $O(|D|^3)$  and the size of nset relations by  $O(|D|^4)$ . All relations have a functional dependency from the expression and the context (columns one to three) to the value (column four). Thus, num and string relations are of size  $O(|D|^3)$  times the maximum size of such values. It remains to be shown that numbers and strings computable in XPath are of size  $O(|D| * |Q|)$ .

Indeed, strings and numbers that may be obtained from the document are guaranteed to be of linear size.

$N_4$			
$x$	$k$	$n$	$val$
$b_2$	1	3	1
$b_3$	2	3	2
$b_4$	3	3	3
$b_3$	1	2	1
$b_4$	2	2	2
$b_4$	1	1	1

$N_5$			
$x$	$k$	$n$	$val$
$b_2$	1	3	3
$b_3$	2	3	3
$b_4$	3	3	3
$b_3$	1	2	2
$b_4$	2	2	2
$b_4$	1	1	1

$N_3$		
$x$	$k$	$n$
$b_1$	1	3
$b_1$	2	3
$b_2$	2	3

$N_2$	
$x$	$val$
$b_1$	$b_2$
$b_1$	$b_3$
$b_2$	$b_3$

$N_1$	
$x$	$val$
$a$	$b_2$
$a$	$b_3$

Figure 6: Context-value tables of Example 6.4.

(Note: For the conversion from a node set to a string or number, only the first node in the set is chosen.)

Of the string functions, only “concat” may produce a string longer than the input strings. (Note that the “translate” function of [17] does not allow for arbitrary but just single-character replacement, e.g. for case-conversion purposes.) The size of such strings is bounded by  $O(|D| * |Q|)$  (In each operation, the can only grow by a constant factor). Likewise, the numbers obtained through arithmetic operations can be represented within the same space bound.

The overall space bound of  $O(|D|^4 * |Q|^2)$  follows. Note that no significant additional amount of space is required for intermediate computations.

Regarding time, the algorithm for evaluating a query  $Q$  bottom-up is as follows. First, we prepare a number of data structures which will allow to evaluate each of the effective semantics functions of Figure 2 in time  $O(I^2)$ , where  $I$  is the size of their arguments. These preprocessing steps can be carried out with a time bound better than the overall time bound to be shown. By considering the definition of  $\mathcal{E}_\uparrow$ , it becomes clear that each of the expression relations can be computed in time  $O(|D|^5 * |Q|)$  at worst when the expression semantics tables of the direct subexpressions are given. (The  $|Q|$  factor is due the size bound on strings and numbers generated during the computation.) Moreover,  $O(|Q|)$  such computations are needed in total to evaluate  $Q$ .

As a final remark, note that contexts can also be represented in terms of pairs of a current and a “previous” context node (rather than triples of a node, a position, and a size), which are defined relative to an axis and a node test (which, however, are fixed with the query). For instance, the corresponding ternary context for  $\vec{c} = \langle x_0, x \rangle$  w.r.t. axis  $\chi$  and node test  $t$  is  $\langle x, \text{id}_{x_\chi}(x, Y), |Y| \rangle$ , where  $Y = \{y \mid x_0 \chi y, y \in T(t)\}$ . Thus, position and size values can be recovered on demand. Using this more sophisticated representation, it is possible to obtain an improved time bound of  $O(|D|^3 * |Q|^2)$  for XPath query evaluation.  $\square$



$\mathcal{S}_\downarrow[\chi::t[e_1] \cdots [e_m]](X_1, \dots, X_k) :=$   
**begin**  
 $S := \{\langle x, y \rangle \mid x \in \bigcup_{i=1}^k X_i, x \chi y, \text{ and } y \in T(t)\};$   
**for each**  $1 \leq i \leq m$  (in ascending order) **do**  
**begin**  
fix some order  $\vec{S} = \langle \langle x_1, y_1 \rangle, \dots, \langle x_l, y_l \rangle \rangle$  for  $S$ ;  
 $\langle r_1, \dots, r_l \rangle := \mathcal{E}_\downarrow[e_i](t_1, \dots, t_l)$   
where  $t_j = \langle y_j, \text{idx}_\chi(y_j, S_j), |S_j| \rangle$   
and  $S_j := \{z \mid \langle x_j, z \rangle \in S\};$   
 $S := \{\langle x_i, y_i \rangle \mid r_i \text{ is true}\};$   
**end;**  
**for each**  $1 \leq i \leq k$  **do**  
 $R_i := \{y \mid \langle x, y \rangle \in S, x \in X_i\};$   
**return**  $\langle R_1, \dots, R_k \rangle;$   
**end;**

$\mathcal{S}_\downarrow[\pi](X_1, \dots, X_k) := \mathcal{S}_\downarrow[\pi](\overbrace{\{\text{root}\}, \dots, \{\text{root}\}}^{k \text{ times}})$   
 $\mathcal{S}_\downarrow[\pi_1/\pi_2](X_1, \dots, X_k) :=$   
 $\mathcal{S}_\downarrow[\pi_2](\mathcal{S}_\downarrow[\pi_1](X_1, \dots, X_k))$   
 $\mathcal{S}_\downarrow[\pi_1 \mid \pi_2](X_1, \dots, X_k) :=$   
 $\mathcal{S}_\downarrow[\pi_1](X_1, \dots, X_k) \cup^\diamond \mathcal{S}_\downarrow[\pi_2](X_1, \dots, X_k)$

Figure 7: Top-down evaluation of location paths.

## 7 Top-down Evaluation of XPath

In the previous section, we obtained a bottom-up semantics definition which led to a polynomial-time query evaluation algorithm for XPath. Despite this favorable complexity bound, this algorithm is still not practical, as usually many irrelevant intermediate results are computed to fill the context-value tables which are not used later on. Next, building on the context-value table principle of Section 6, we develop a top-down algorithm based on vector computation for which the favorable (worst-case) complexity bound carries over but in which the computation of a large number of irrelevant results is avoided.

We introduce a family of tuple operators to simplify the following discussion. Given a unary operation  $Op : D \rightarrow D$ , the operation  $Op^\diamond : D^k \rightarrow D^k$  is defined as  $Op^\diamond(x_1, \dots, x_k) := \langle Op(x_1), \dots, Op(x_k) \rangle$ . Analogously, given a binary operation  $\circ : D \times D \rightarrow D$ , the operation  $\circ^\diamond : D^k \times D^k \rightarrow D^k$  is defined as  $\langle x_1, \dots, x_k \rangle \circ^\diamond \langle y_1, \dots, y_k \rangle := \langle x_1 \circ y_1, \dots, x_k \circ y_k \rangle$ . For instance,  $\langle X_1, \dots, X_k \rangle \cup^\diamond \langle Y_1, \dots, Y_k \rangle := \langle X_1 \cup Y_1, \dots, X_k \cup Y_k \rangle$ .  $\text{count}^\diamond(X_1, \dots, X_k) := \langle |X_1|, \dots, |X_k| \rangle$  computes the cardinalities of a tuple of sets element-wise,  $\text{ss}^\diamond(x_1, \dots, x_k) := \langle \{x_1\}, \dots, \{x_k\} \rangle$  lifts a tuple of elements to a tuple of singleton sets, and  $\text{proj}_i^\diamond(\vec{c}_1, \dots, \vec{c}_l)$  selects the  $i$ -th elements from the tuples  $\vec{c}_1, \dots, \vec{c}_l$ .

For practical reasons, before we arrive at the point of defining a top-down semantics function  $\mathcal{E}_\downarrow$  for XPath, we introduce an auxiliary semantics definition for location paths,  $\mathcal{S}_\downarrow$ :

$\mathcal{S}_\downarrow : \text{LocationPath} \rightarrow \text{List}(2^{\text{dom}}) \rightarrow \text{List}(2^{\text{dom}})$

That is, given a location path  $\pi$  and a list  $\langle X_1, \dots, X_k \rangle$

of node sets,  $\mathcal{S}_\downarrow$  determines a list  $\langle Y_1, \dots, Y_k \rangle$  of node sets, s.t. for every  $i \in \{1, \dots, k\}$ , the nodes reachable from the context nodes in  $X_i$  via the location path  $\pi$  are precisely the nodes in  $Y_i$ .  $\mathcal{S}_\downarrow[\pi]$  can be obtained from the relations  $\mathcal{E}_\uparrow[\pi]$  as follows. Let

$$\mathcal{S}_\downarrow[\pi](X_1, \dots, X_k) = \langle Y_1, \dots, Y_k \rangle$$

Then, a node  $y$  is in  $Y_i$  iff there is an  $x \in X_i$  and some  $p, s$  such that  $\langle x, p, s, y \rangle \in \mathcal{E}_\uparrow[\pi]$ . For the actual computation of  $\mathcal{S}_\downarrow[\pi]$ , we basically distinguish the same cases (related to location paths) as for the bottom-up semantics  $\mathcal{E}_\uparrow[\pi]$  (see Figure 7).

**Definition 7.1** The semantics function  $\mathcal{E}_\downarrow$  for arbitrary XPath expressions is of the following type:

$$\begin{aligned} \mathcal{E}_\downarrow : \text{XPathExpression} &\rightarrow \text{List}(\mathbf{C}) \\ &\rightarrow \text{List}(\text{XPathType}) \end{aligned}$$

Given an XPath expression  $e$  and a list  $\langle \vec{c}_1, \dots, \vec{c}_l \rangle$  of contexts,  $\mathcal{E}_\downarrow$  determines a list  $\langle r_1, \dots, r_l \rangle$  of results of one of the XPath types number, string, boolean, or node set.  $\mathcal{E}_\downarrow$  is defined as

$$\mathcal{E}_\downarrow[\pi](\vec{c}_1, \dots, \vec{c}_l) := \mathcal{S}_\downarrow[\pi](\text{ss}^\diamond(\text{proj}_1^\diamond(\vec{c}_1, \dots, \vec{c}_l)))$$

$$\mathcal{E}_\downarrow[\text{position}()](\langle x_1, k_1, n_1 \rangle, \dots, \langle x_l, k_l, n_l \rangle) := \langle k_1, \dots, k_l \rangle$$

$$\mathcal{E}_\downarrow[\text{last}()](\langle x_1, k_1, n_1 \rangle, \dots, \langle x_l, k_l, n_l \rangle) := \langle n_1, \dots, n_l \rangle$$

and

$$\begin{aligned} \mathcal{E}_\downarrow[Op(e_1, \dots, e_m)](\vec{c}_1, \dots, \vec{c}_l) := \\ \mathcal{F}[Op]^\diamond(\mathcal{E}_\downarrow[e_1](\vec{c}_1, \dots, \vec{c}_l), \dots, \mathcal{E}_\downarrow[e_m](\vec{c}_1, \dots, \vec{c}_l)) \end{aligned}$$

for the remaining kinds of expressions.  $\square$

For example,

$$\begin{aligned} \mathcal{E}_\downarrow[e_1 \text{ ArithOp } e_2](\vec{c}_1, \dots, \vec{c}_l) := \\ \mathcal{E}_\downarrow[e_1](\vec{c}_1, \dots, \vec{c}_l) \text{ ArithOp}^\diamond \mathcal{E}_\downarrow[e_2](\vec{c}_1, \dots, \vec{c}_l) \end{aligned}$$

$$\begin{aligned} \mathcal{E}_\downarrow[\text{count}(e)](\vec{c}_1, \dots, \vec{c}_l) := \\ \text{count}^\diamond(\mathcal{E}_\downarrow[e](\vec{c}_1, \dots, \vec{c}_l)) \end{aligned}$$

$$\begin{aligned} \mathcal{E}_\downarrow[e_1 \text{ RelOp } e_2](\vec{c}_1, \dots, \vec{c}_l) := \\ /* \text{ where } e_1 \text{ and } e_2 \text{ are of type num } */ \\ \mathcal{E}_\downarrow[e_1](\vec{c}_1, \dots, \vec{c}_l) \text{ RelOp}^\diamond \mathcal{E}_\downarrow[e_2](\vec{c}_1, \dots, \vec{c}_l) \end{aligned}$$

**Example 7.2** Consider the XPath query

$$\begin{aligned} / \text{descendant}::\text{a}[\text{count}(\text{descendant}::\text{b}/\text{child}::\text{c}) \\ + \text{position}() < \text{last}()]/\text{child}::\text{d} \end{aligned}$$

Let  $L = \langle \langle y_1, 1, l \rangle, \dots, \langle y_l, l, l \rangle \rangle$ , where the  $y_i$  are those nodes reachable from the root node through the descendant axis and which are labeled “a”. The query is evaluated top-down as

$$\mathcal{S}_\downarrow[\text{child}::\text{d}](\mathcal{S}_\downarrow[\text{descendant}::\text{a}[e]](\{\text{root}\}))$$

where  $\mathcal{E}_\downarrow[e](L)$  is computed as

$$(\text{count}^\diamond(\pi) + \mathcal{E}_\downarrow[\text{position}()](L)) <^\diamond \mathcal{E}_\downarrow[\text{last}()](L)$$

and

$$\pi = \mathcal{S}_\downarrow[\text{child}::c](\mathcal{S}_\downarrow[\text{descendant}::b](\text{ss}^\downarrow(\text{proj}_1^\downarrow(L))))).$$

Note that the arity of the tuples used to compute the outermost location path is one, while it is  $l$  for  $e$ .  $\square$

The correctness of the top-down semantics follows immediately from the corresponding result in the bottom-up case and from the definition of  $\mathcal{S}_\downarrow$  and  $\mathcal{E}_\downarrow$ .

**Theorem 7.3** (*Correctness of  $\mathcal{E}_\downarrow$* ) *Given an arbitrary XPath expression  $e$ ,  $\langle v_1, \dots, v_l \rangle = \mathcal{E}_\downarrow[e](\bar{c}_1, \dots, \bar{c}_l)$  iff  $\langle \bar{c}_1, v_1 \rangle, \dots, \langle \bar{c}_l, v_l \rangle \in \mathcal{E}_\downarrow[e]$ .*

$\mathcal{S}_\downarrow$  and  $\mathcal{E}_\downarrow$  can be immediately transformed into function definitions in a top-down algorithm. We thus have to define one evaluation function for each case of the definition of  $\mathcal{S}_\downarrow$  and  $\mathcal{E}_\downarrow$ , respectively. The functions corresponding to the various cases of  $\mathcal{S}_\downarrow$  have a location path and a list of node sets of variable length  $(X_1, \dots, X_k)$  as input parameter and return a list  $(R_1, \dots, R_k)$  of node sets of the same length as result. Likewise, the functions corresponding to  $\mathcal{E}_\downarrow$  take an arbitrary XPath expression and a list of contexts as input and return a list of XPath values (which can be of type num, str, bool or nset). Moreover, the recursions in the definition of  $\mathcal{S}_\downarrow$  and  $\mathcal{E}_\downarrow$  correspond to recursive function calls of the respective evaluation functions. Analogously to Theorem 6.5, we get

**Theorem 7.4** *The immediate functional implementation of  $\mathcal{E}_\downarrow$  evaluates XPath queries in polynomial time (combined complexity).*

Finally, note that using arguments relating the top-down method of this section with (join) optimization techniques in relational databases, one may argue that the context-value table principle is also the basis of the polynomial-time bound of Theorem 7.4.

## 8 Linear-time Fragments of XPath

### 8.1 Core XPath

In this section, we define a fragment of XPath (called Core XPath) which constitutes a clean logical core of XPath (cf. [8, 9]). The only objects that are manipulated in this language are sets of nodes (i.e., there are no arithmetical or string operations). Besides from these restrictions, the full power of location paths is supported, and so is the matching of such paths in condition predicates (with an “exists” semantics), and the closure of such condition expressions with respect to boolean operations “and”, “or”, and “not”.

We define a mapping of each query in this language to a simple algebra over the set operations  $\cap$ ,  $\cup$ , ‘-’,  $\chi$  (the axis functions from Definition 4.1), and an operation  $\frac{\text{dom}}{\text{root}}(S) := \{x \in \text{dom} \mid \text{root} \in S\}$ , i.e.  $\frac{\text{dom}}{\text{root}}(S)$  is dom if  $\text{root} \in S$  and  $\emptyset$  otherwise.

Note that each XPath axis has a natural *inverse*:  $\text{self}^{-1} = \text{self}$ ,  $\text{child}^{-1} = \text{parent}$ ,  $\text{descendant}^{-1}$

$= \text{ancestor}$ ,  $\text{descendant-or-self}^{-1} = \text{ancestor-or-self}$ ,  $\text{following}^{-1} = \text{preceding}$ , and  $\text{following-sibling}^{-1} = \text{preceding-sibling}$ .

**Lemma 8.1** *Let  $\chi$  be an axis. For each pair of nodes  $x, y \in \text{dom}$ ,  $x\chi y$  iff  $y\chi^{-1}x$ .*

(Proof by a very easy induction.)

**Definition 8.2** Let the (abstract) syntax of the Core XPath language be defined by the EBNF grammar

```

exp:          locationpath | ‘/’ locationpath
locationpath: locationstep (‘/’ locationstep)*
locationstep:  $\chi$  ‘::’  $t$  |  $\chi$  ‘::’  $t$  [‘] pred ‘]’
pred:         pred ‘and’ pred | pred ‘or’ pred
              | ‘not’ (‘(’ pred ‘)’) | exp | (‘(’ pred ‘)’)

```

“exp” is the start production,  $\chi$  stands for an axis (see above), and  $t$  for a “node test” (either an XML tag or “\*”, meaning “any label”). The semantics of Core XPath queries is defined by a function  $\mathcal{S}_\rightarrow$

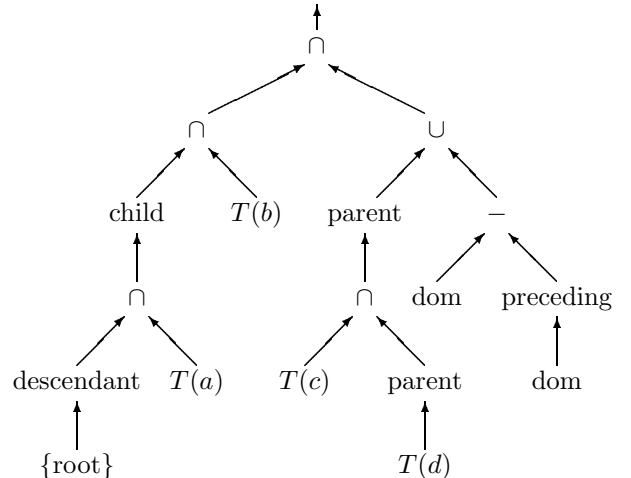
$$\begin{aligned} \mathcal{S}_\rightarrow[\chi::t[e]](N_0) &:= \chi(N_0) \cap T(t) \cap \mathcal{E}_1[e] \\ \mathcal{S}_\rightarrow[/\chi::t[e]](N_0) &:= \chi(\{\text{root}\}) \cap T(t) \cap \mathcal{E}_1[e] \\ \mathcal{S}_\rightarrow[\pi/\chi::t[e]](N_0) &:= \chi(\mathcal{S}_\rightarrow[\pi](N_0)) \cap T(t) \cap \mathcal{E}_1[e] \\ \mathcal{S}_\leftarrow[\chi::t[e]] &:= \chi^{-1}(T(t) \cap \mathcal{E}_1[e]) \\ \mathcal{S}_\leftarrow[\chi::t[e]/\pi] &:= \chi^{-1}(\mathcal{S}_\leftarrow[\pi] \cap T(t) \cap \mathcal{E}_1[e]) \\ \mathcal{S}_\leftarrow[/\chi::t[e]] &:= \frac{\text{dom}}{\text{root}}(\mathcal{S}_\leftarrow[\pi]) \\ \mathcal{E}_1[e_1 \text{ and } e_2] &:= \mathcal{E}_1[e_1] \cap \mathcal{E}_1[e_2] \\ \mathcal{E}_1[e_1 \text{ or } e_2] &:= \mathcal{E}_1[e_1] \cup \mathcal{E}_1[e_2] \\ \mathcal{E}_1[\text{not}(e)] &:= \text{dom} - \mathcal{E}_1[e] \\ \mathcal{E}_1[\pi] &:= \mathcal{S}_\leftarrow[\pi] \end{aligned}$$

where  $N_0$  is a set of context nodes or dom and a query  $\pi$  evaluates as  $\mathcal{S}_\rightarrow[\pi](N_0)$ .  $\square$

**Example 8.3** The Core XPath query

/descendant::a/child::b[child::c/child::d or  
not(following::\*)]

is evaluated as specified by the query tree



(There are alternative but equivalent query trees due to the associativity and commutativity of some of our operators.)  $\square$

The semantics of XPath and Core XPath (defined using  $\mathcal{S}_\leftarrow$ ,  $\mathcal{S}_\rightarrow$ , and  $\mathcal{E}_1$ ) coincide in the following way:

**Theorem 8.4** *Let  $\pi$  be a Core XPath query and  $N_0 \subseteq \text{dom}$  be a set of context nodes. Then,*

$$\begin{aligned}\mathcal{S}_\leftarrow[\pi] &= \{x \mid \mathcal{S}_\downarrow[\pi](\{x\}) \neq \emptyset\} \\ \mathcal{E}_1[e] &= \{x \mid \mathcal{E}_\downarrow[e](\{\langle x, 1, 1 \rangle\})\} \\ \langle \mathcal{S}_\rightarrow[\pi](N_0) \rangle &= \mathcal{S}_\downarrow[\pi](\langle N_0 \rangle).\end{aligned}$$

This can be shown by easy induction proofs. Thus, Core XPath (evaluated using  $\mathcal{S}_\rightarrow$ ) is a fragment of XPath, both syntactically and semantically.

**Theorem 8.5** *Core XPath queries can be evaluated in time  $O(|D| * |Q|)$ , where  $|D|$  is the size of the data and  $|Q|$  is the size of the query.*

**Proof** Given a query  $Q$ , it can be rewritten into an algebraic expression  $E$  over the operations  $\chi$ ,  $\cup$ ,  $\cap$ ,  $'\_'$ , and  $\frac{\text{dom}}{\text{root}}$  using  $\mathcal{S}_\rightarrow$ ,  $\mathcal{S}_\leftarrow$ , and  $\mathcal{E}_1$  in time  $O(|Q|)$ . Each of the operations in our algebra can be carried out in time  $O(|D|)$ . Since at most  $O(|Q|)$  such operations need to be carried out to process  $E$ , the complexity bound follows.  $\square$

## 8.2 XPatterns

We extend our linear-time fragment Core XPath by the operation  $\text{id}: \text{nset} \rightarrow \text{nset}$  of Table 4 by defining “id” as an axis relation

$$\text{id} := \{\langle x_0, x \rangle \mid x_0 \in \text{dom}, x \in \text{deref\_ids}(\text{strval}(x_0))\}$$

Queries of the form  $\pi_1/\text{id}(\pi_2)/\pi_3$  are now treated as  $\pi_1/\pi_2/\text{id}/\pi_3$ .

**Lemma 8.6** *Let  $\pi_1/\text{id}(\pi_2)/\pi_3$  be an XPath query s.t.  $\pi_1/\pi_2/\text{id}/\pi_3$  is a query in Core XPath with the “id” axis. Then, the semantics of the two queries relative to a set of context nodes  $N_0 \in \text{dom}$  coincide,  $\mathcal{S}_\downarrow[\pi_1/\text{id}(\pi_2)/\pi_3](\langle N_0 \rangle) = \mathcal{S}_\rightarrow[\pi_1/\pi_2/\text{id}/\pi_3](N_0)$ .  $\square$*

**Theorem 8.7** *Queries in Core XPath with the “id” axis can be evaluated in time  $O(|D| * |Q|)$ .*

**Proof Sketch.** The hard part of this proof is to define a function  $\text{id}: 2^{\text{dom}} \rightarrow 2^{\text{dom}}$  and its inverse consistent with the functions of Definition 4.1 which is computable in linear time. We make use of a binary auxiliary relation “ref” which contains a tuple of nodes  $\langle x, y \rangle$  iff the text belonging to  $x$  in the XML document, but which is directly inside it and not further down in any of its descendants, contains a whitespace-separated string referencing the identifier of node  $y$ .

“@n”, “@*”, “text()”, “comment()”, “pi(n)”, and “pi()” (where n is a label) are simply sets provided with the document (similar to those obtained through the node test function $T$ ).
“=s” (s is a string) can be encoded as a unary predicate whose extension can be computed using string search in the document before the evaluation of our query starts. Clearly, this can be done in linear time.
first-of-any := $\{y \in \text{dom} \mid \exists x : \text{nextsibling}(x, y)\}$
last-of-any := $\{x \in \text{dom} \mid \exists y : \text{nextsibling}(x, y)\}$
“id(s)” is a unary predicate and can easily be computed (in linear time) before the query evaluation.

Table 5: Some unary predicates of XLST Patterns [18].

**Example.** Let  $\text{id}(i) = n_i$ . For the XML document  $\langle \text{t id=1} \rangle 3 \langle \text{t id=2} \rangle 1 \langle / \text{t} \rangle \langle \text{t id=3} \rangle 1 2 \langle / \text{t} \rangle \langle / \text{t} \rangle$ , we have  $\text{ref} := \{\langle n_1, n_3 \rangle, \langle n_2, n_1 \rangle, \langle n_3, n_1 \rangle, \langle n_3, n_2 \rangle\}$ .  $\square$

“ref” can be efficiently computed in a preprocessing step. It does not satisfy any functional dependencies, but it is guaranteed to be of linear size w.r.t. the input data (however, not in the tree nodes). Now we can encode  $\text{id}(S)$  as those nodes reachable from  $S$  and its descendants using “ref”.

$$\begin{aligned}\text{id}(S) &:= \{y \mid x \in \text{descendant-or-self}(S), \langle x, y \rangle \in \text{ref}\} \\ \text{id}^{-1}(S) &:= \text{ancestor-or-self}(\{x \mid \langle x, y \rangle \in \text{ref}, y \in S\})\end{aligned}$$

This computation can be performed in linear time.  $\square$

We may define XPatterns as the smallest language that subsumes Core XPath and the XSLT Pattern language of [18] (see also [15] for a good and formal overview of this language) and is (syntactically) contained in XPath. Stated differently, it is obtained by extending the language of [18] without the first-of-type and last-of-type predicates (which do not exist in XPath) to support all of the XPath axes. As pointed out in the introduction, XPatterns is an interesting and practically useful query language. Surprisingly, XPatterns queries can be evaluated in linear time.

**Theorem 8.8** *Let  $D$  be an XML document and  $Q$  be an XPatterns query. Then,  $Q$  can be evaluated on  $D$  in time  $O(|D| * |Q|)$ .*

**Proof (Rough Sketch).** XPatterns extends Core XPath by the “id” axis and a number of features which are definable as unary predicates, of which we give an overview in Table 5. It becomes clear by considering the semantics definition of [15] that after parsing the query, one knows of a fixed number of predicates to populate, and this action takes time  $O(|D|)$  for each. Thus, since this computation precedes the query evaluation – which has a time bound of  $O(|D| * |Q|)$  – this does not pose a problem. “id(s)” (for some fixed string  $s$ ) may only occur at the beginning of a path, thus in a query of the form  $\text{id}(s)/\pi$ ,  $\pi$  is evaluated relative to the set  $\text{id}(s)$  just as, say,  $\{\text{root}\}$  is for query  $/\pi$ .  $\square$

Note that the unary first-of-type and last-of-type predicates can be computed in time  $O(|D| * |\Sigma|)$  when

$ Q $	IE6 10	IE6 20	IE6 200	New 10	New 20	New 200
1				0	0	0.02
2			2	0	0	0.05
3			346	0	0	0.06
4		1	-	0	0	0.07
5		21	-	0	0	0.10
6	5	406	-	0	0.01	0.11
7	42	-	-	0.01	0.01	0.13
8	437	-	-	0	0.01	0.16
⋮						
16	-	-	-	0.01	0.02	0.30

Figure 8: Benchmark results in seconds for IE6 vs. our implementation (“New”), on the queries of Experiment 2 and document sizes 10, 20, and 200.

parsing the document, but are of size  $O(|D|)$ :

$$\text{first-of-type}() := \bigcup_{l \in \Sigma} (T(l) - \text{nextsibling}^+(T(l)))$$

$$\text{last-of-type}() := \bigcup_{l \in \Sigma} (T(l) - (\text{nextsibling}^{-1})^+(T(l)))$$

where  $R^+ = R.R^*$ .

## 9 Conclusions

In this paper, we presented the first XPath query evaluation algorithm that runs in polynomial time with respect to the size of both the data and of the query. Our results will empower XPath engines to be able to deal efficiently with very sophisticated queries.

We have made a main-memory implementation of the top-down algorithm of Section 7. Figure 8 compares it to IE6 along the assumptions made in Experiment 2 (i.e., the queries of which were strictly the most demanding of all three experiments). It shows that our algorithm scales linearly in the size of the queries and quadratically (for this class of queries) in the size of the data. Our implementation is still an early, naive prototype without any optimizations, and which strictly coheres to the specification given in this paper. We plan to significantly improve on its real-world runtime in terms of data in the future. Resources and further benchmarks that become available in the course of this effort will be made accessible at

<http://www.xmltaskforce.com>

Apart from full proofs of our theorems, the long version of this paper will discuss further large XPath fragments which can be processed in improved time and space bounds. Due to lack of space, no treatment of these fragments was possible in this paper. In the future, we intend to work on algorithms for processing XPath with disk access and with streaming XML data.

## Acknowledgments

We thank G. Moerkotte and the anonymous reviewers of VLDB 2002, whose constructive comments

have helped to considerably improve this paper, and P. Fankhauser for a wealth of pointers to the literature.

## References

- [1] S. Abiteboul, P. Buneman, and D. Suciu. *Data on the Web*. Morgan Kaufmann Publishers, 2000.
- [2] G. J. Bex, S. Maneth, and F. Neven. A Formal Model for an Expressive Fragment of XSLT. In *Proc. CL 2000*, LNCS 1861, pages 1137–1151. Springer, 2000.
- [3] N. Bruno, D. Srivastava, and N. Koudas. “Holistic Twig Joins: Optimal XML Pattern Matching”. In *ACM SIGMOD 2002*, June 3-6 2002.
- [4] F. Bry, D. Olteanu, H. Meuss, and T. Furche. Symmetry in XPath. Technical Report PMS-FB-2001-16, LMU München, 2001. Short version.
- [5] J. Clark. XT. A Java Implementation of XSLT <http://www.jclark.com/xml/xt.html/>.
- [6] A. Deutsch and V. Tannen. Containment and Integrity Constraints for XPath. In *KRDB 2001*, CEUR Workshop Proceedings 45, 2001.
- [7] P. Fankhauser. A Mapping of XPath 1.0 to the XML Query Algebra (with J. Clark, M. Fernandez, and J. Siméon), Nov. 2001. Personal Communication.
- [8] G. Gottlob and C. Koch. “Monadic Datalog and the Expressive Power of Web Information Extraction Languages”. In *Proc. PODS’02*, Madison, Wisconsin, 2002. to appear.
- [9] G. Gottlob and C. Koch. “Monadic Queries over Tree-Structured Data”. In *Proc. LICS’02*, Copenhagen, Denmark, July 22–25 2002. to appear.
- [10] P. Kilpeläinen. *Tree Matching Problems with Applications to Structured Text Databases*. PhD thesis, Department of Computer Science, University of Helsinki, Nov. 1992. Report A-1992-6.
- [11] G. Miklau and D. Suciu. “Containment and Equivalence for an XPath Fragment”. In *Proc. PODS’02*, pages 65–76, Madison, Wisconsin, 2002.
- [12] T. Milo, D. Suciu, and V. Vianu. “Typechecking for XML Transformers”. In *Proc. PODS’00*, pages 11–22, 2000.
- [13] D. Shasha, J. T. L. Wang, and R. Giugno. “Algorithmics and Applications of Tree and Graph Searching”. In *Proc. PODS’02*, June 3 – 5 2002.
- [14] P. Wadler. “Two Semantics for XPath”, 2000. Draft paper available at <http://www.research.avayalabs.com/user/wadler/topics/recent.html>.
- [15] P. Wadler. “A Formal Semantics of Patterns in XSLT”. In *Markup Technologies*, Philadelphia, December 1999. Revised version in *Markup Languages*, MIT Press, June 2001.
- [16] P. T. Wood. On the Equivalence of XML Patterns. In *Proc. CL 2000*, LNCS 1861, pages 1152–1166. Springer, 2000.
- [17] World Wide Web Consortium. XPath Recommendation. <http://www.w3c.org/TR/xpath/>.
- [18] World Wide Web Consortium. XSL Working Draft <http://www.w3.org/TR/1998/WD-xsl-19981216>.
- [19] Xalan-Java version 2.2.D11. <http://xml.apache.org/xalan-j/>.