

XPath: Looking Forward

Dan Olteanu¹, Holger Meuss², Tim Furche¹, and François Bry¹

¹ Insitute for Computer Science, University of Munich, Germany
{olteanu, timfu, bry}@informatik.uni-muenchen.de

² Center for Information and Language Processing, University of Munich, Germany
meuss@cis.uni-muenchen.de

Abstract. The location path language XPath is of particular importance for XML applications since it is a core component of many XML processing standards such as XSLT or XQuery. In this paper, based on axis symmetry of XPath, equivalences of XPath 1.0 location paths involving reverse axes, such as **anc** and **prec**, are established. These equivalences are used as rewriting rules in an algorithm for transforming location paths with reverse axes into equivalent reverse-axis-free ones. Location paths without reverse axes, as generated by the presented rewriting algorithm, enable efficient SAX-like streamed data processing of XPath.

1 Introduction

Query languages for XML and semistructured data rely on location paths for selecting nodes in data items. In particular, XQuery [1] and XSLT [2] are based on XPath [3]. XPath takes a navigational approach for specifying the nodes to be selected, described by a large number of navigational axes, e.g. **child**, **descendant**, **preceding**. The number as well as the relevance of these navigational axes for querying XML has been challenged in [4,1,5].

The random access to XML data that is enabled by the various navigational axes of XPath has proven particularly difficult for an efficient stream-based processing of XPath queries. Processing of XML has seen the widespread use of the W3C document object model (DOM) [6], where an in-memory representation of the entire XML data is used. As DOM has been developed with focus on document processing in user agents (e.g. browsers), this approach has several shortcomings for other application areas:

First, a considerable amount of XML applications, in particular data-centric applications, handle documents too large to be processed in memory. Such documents are often encountered in natural language processing [7], in biological [8] and astronomical [9] projects.

Second, the need for progressive processing (also referred to as sequential processing) of XML has emerged: Stream-based processing generating partial results as soon as they are available gives rise to a more efficient evaluation in certain contexts, e.g.:

- For selective dissemination of information (SDI), documents have to be filtered according to complex requirements specified as XPath queries before

being distributed to the subscribers [10,11]. The routing of data to selected receivers is also becoming increasingly important in the context of web service architectures.

- To integrate data over the Internet, in particular from slow sources, it is desirable to progressively process the input before the full data is retrieved [12,13].
- As a general processing scheme for XML, several solutions for pipelined processing have been suggested, where the input is sent through a chain of processors each of which taking the output of the preceding processor as input, e.g. Apache Cocoon [14].
- For progressive rendering of large documents, e.g. by means of XSL(T), cf. Requirement 19 of [5]. There have been attempts to solve this problem [15].

There is a great interest in the identification of a subset of XPath that allows efficient streamed and progressive processing, cf. [4] and Requirement 19 of [5].

For stream-based processing of XML data, the Simple API for XML (SAX) [16] has been specified. Of particular concern for progressive SAX-like processing are the reverse axes of XPath, i.e. those navigational axes (e.g. `par`, `prec`) that select nodes occurring before the context node in document order. A restriction to forward axes (i.e. axes selecting only nodes after the context node) in location paths is a straightforward consideration for an efficient stream-based evaluation of XPath queries [4].

There are three principal options how to evaluate reverse axes in a stream-based context:

- Storing in memory sufficient information that allows to access past events when evaluating a reverse axis. This amounts to keeping in memory a (possibly pruned) DOM representation of the data [15].
- Evaluating an XPath expression in more than one run. With this approach, it is also necessary to store additional information to be used in successive runs. This information can be considerably smaller than what is needed in the first approach.
- Replacing XPath expressions by equivalent ones without reverse axes.

In this paper it is shown that the third approach is possible. It is less time consuming than the second approach and does not require the in-memory storage of fragments of the input as the first approach does. Hence, XPath can be evaluated without restriction on the use of reverse axes.

Section 2 specifies the location path language considered in the rest of the paper. Then, the notion of equivalence between location paths is defined in Section 3 using a formal model and a denotational semantics for XPath based on [17,18]. Furthermore, two sets of equivalences with rather different properties are established. These equivalences are used as rewriting rules in an algorithm, called “rare”, for transforming absolute XPath location paths with reverse axes into equivalent reverse-axis-free ones, as presented in Section 4. Two rewritings, based on the two rule sets, are considered. In Section 5, related work is discussed and Section 6 concludes the paper.

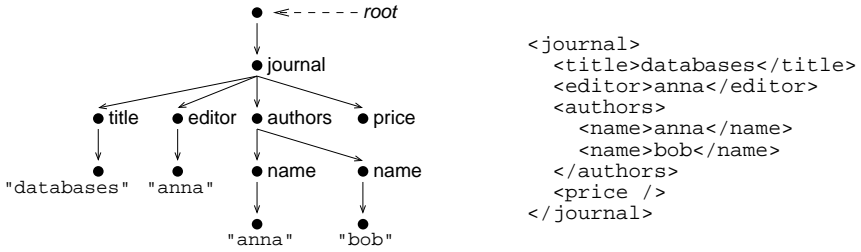


Fig. 1. Tree and XML data it represents

Due to space limitations, parts of this work have been omitted, most notably the equivalence proofs. They can be found in the full version [19] of this paper.

2 Preliminaries

In this paper, specificities of XML that are irrelevant to the issue of concern are left out. Thus, namespaces, comments, processing instructions, attributes, schema types, and references are not considered. The results given in this paper extend straightforwardly to unrestricted XML documents.

The root node of a document corresponds to the document node of DOM and of the XQuery 1.0 and XPath 2.0 Data Model [20] – i.e. it is none of the document elements. A leaf is an empty element or a text node – cf. Figure 1.

The mathematical model used in this paper is adapted from [17,18]. The full formal model as well as the denotational semantics can be found in the full version [19] of this paper. It consists of mathematical functions that can be seen as (formal specifications of) elementary procedures.

2.1 Location Path Language

The location path language considered in the following is XPath without constructs, such as attribute selection and functions, irrelevant to the issue of concern. Note that while functions are not considered in the following sections, the results almost immediately apply to XPath 1.0 [3] location paths with functions. The only class of functions that needs special treatment are functions for accessing the context position or size of a node.

For convenience, this language will be referred to as *xPath*, as defined in Figure 2. Note that the axes specified in *xPath* are shorthands of XPath axes.

\perp is used as a canonical equivalent path to the *xPath* expressions that select no nodes whatever the context node and document are, e.g. `/par::*`.

$p_1 == p_2$ expresses node equality based on identity. Thus, if p_1 and p_2 are two paths, then $p_1 == p_2$ holds if there is a node selected by p_1 which is *identical* to a node selected by p_2 . $==$ corresponds to built-in node equality operator (`==`) in XPath 2.0 and XQuery 1.0, but it can also be used for comparing

```

path ::= path | path / path | path / path | path [qualif] | axis :: nodetest | ⊥ .
qualif ::= qualif and qualif | qualif or qualif | (qualif) |
        path = path | path == path | path .
axis ::= reverse_axis | forward_axis .
reverse_axis ::= par | anc | anc-or-self | prec | prec-sibl .
forward_axis ::= self | child | desc | desc-or-self | foll | foll-sibl .
nodetest ::= tagname | * | text() | node() .

```

Fig. 2. Grammar for XPath

node sets similar to general comparisons in XPath 2.0. As XPath 1.0 has built-in support only for equality based on node values, the XPath 1.0 expression $\text{count}(p_1 | p_2) < \text{count}(p_1) + \text{count}(p_2)$ can be used for expressing $=$.

A *path* expression will be called a “location path”, or “path” for short. A *qualif* expression is a “qualifier” (or condition). Expressions $\text{axis} :: \text{nodetest}$ and $\text{axis} :: \text{nodetest}[\text{qualif}]$ are “location steps”, or “steps”. The length of a location path is the number of location steps it contains outside and inside qualifiers. Note that every location path is a qualifier, but the converse is false.

Absolute location paths are recursively defined as follows: A disjunctive path, i.e. a path of the form $p_1 | \dots | p_i | \dots | p_k$, is an absolute path if for all $i = 1, \dots, k$, p_i is an absolute path. A non-disjunctive path is an absolute path if it is of the form $/p$, where p is a path. A location path, which is not an absolute path, is a “relative path”. A step is a “forward step”, if its axis is a forward axis, or a “reverse step”, if its axis is a reverse axis.

The axes of the following pairs are “symmetrical” of each other: **par** – **child**, **anc** – **desc**, **desc-or-self** – **anc-or-self**, **prec** – **foll**, **prec-sibl** – **foll-sibl**, and (useful in proofs) **self** – **self**.

[17] and [18] give a denotational semantics for XPath, which is slightly modified for our purpose in [19]. The semantics defines a function \mathcal{S} that assigns a set of nodes to a location path and a context node: $\mathcal{S}[[p]]x$ is the set of nodes selected by p from node x .

3 Location Path Equivalences

A set of simple equivalences is first established. These are then used to prove equivalences of paths with reverse axes. We distinguish between general equivalences that can be applied to remove any reverse axis, and specific equivalences, each of them being applicable to a certain case. Making use of the semantics of XPath given in the full version [19] of this paper, the equivalence of location paths can be formally defined as follows.

Definition 1 (Path equivalence). *Two location paths p_1 and p_2 are equivalent, noted $p_1 \equiv p_2$, if $\mathcal{S}[[p_1]] = \mathcal{S}[[p_2]]$, i.e. if $\mathcal{S}[[p_1]]x = \mathcal{S}[[p_2]]x$ for all nodes x (from any document).*

Intuitively, two location paths are equivalent if they select the same set of nodes for every document and every context node in this document.

Lemma 1. *Let p , p_1 , and p_2 be location paths, q , q_1 , and q_2 qualifiers, n a node test, and $\theta \in \{=, \neq\}$.*

1. Right step adjunction: *If $p_1 \equiv p_2$ and p relative, then $p_1/p \equiv p_2/p$.*
2. Left step adjunction: *If $p_1 \equiv p_2$ and p_1, p_2 relative, then $p/p_1 \equiv p/p_2$.*
3. Qualifier adjunction: *If $p_1 \equiv p_2$, then $p_1[q] \equiv p_2[q]$ and $p[p_1] \equiv p[p_2]$.*
4. Relative/absolute path conversion: *If $p_1 \equiv p_2$, then $/p_1 \equiv /p_2$.*
5. Qualifier flattening: *$p[p_1/p_2] \equiv p[p_1[p_2]]$.*
6. **anc-or-self** decomposition: *$\mathbf{anc-or-self}::n \equiv \mathbf{anc}::n \mid \mathbf{self}::n$.*
7. **desc-or-self** decomposition: *$\mathbf{desc-or-self}::n \equiv \mathbf{desc}::n \mid \mathbf{self}::n$.*
8. Qualifiers with joins: *$p[p_1 \theta /p_2] \equiv p[p_1[\mathbf{self}::\mathbf{node}() \theta /p_2]]$.*

Recall that \perp is a location path never selecting any node whatever the context node and document are. Since the root node has no parents and no siblings, the following holds:

Lemma 2. *Let m and n be node tests, i.e. m and n are tag names or one of the XPath constructs $*$, $\mathbf{node}()$, or $\mathbf{text}()$.*

- *Let a be one of the axes \mathbf{par} , \mathbf{anc} , \mathbf{prec} , $\mathbf{prec-sibl}$, \mathbf{self} , \mathbf{foll} , or $\mathbf{foll-sibl}$. The following holds:*

$$/a::n \equiv \begin{cases} / & \text{if } a = \mathbf{self}, n = \mathbf{node}() \\ \perp & \text{otherwise} \end{cases} \quad (1)$$

- *Let a be the \mathbf{prec} or \mathbf{anc} axis. The following holds:*

$$/\mathbf{child}::m/a::n \equiv \begin{cases} / \mathbf{self}::\mathbf{node}() [\mathbf{child}::m] & \text{if } a = \mathbf{anc}, n = \mathbf{node}() \\ \perp & \text{otherwise} \end{cases} \quad (2)$$

$$/\mathbf{child}::m[a::n] \equiv \begin{cases} / \mathbf{child}::m & \text{if } a = \mathbf{anc}, n = \mathbf{node}() \\ \perp & \text{otherwise} \end{cases} \quad (3)$$

3.1 General Equivalences

The nodes selected by a reverse step within a location path are necessarily descendants of the document root. The following equivalences show how for any reverse axis only those descendants of the root can be selected that are also matched by the original reverse step.

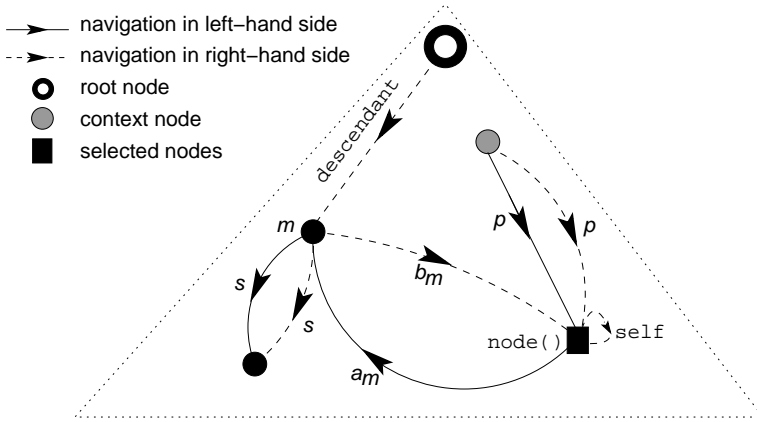


Fig. 3. Tree navigation used in Equivalence (4)

Proposition 1. *Let p and s be relative location paths, n and m node tests, a_m a reverse axis, a_n a forward axis, and b_m the symmetrical axis of a_m . Then the following holds*

$$p[a_m::m/s] \equiv p[/desc::m[s]/b_m::node() == self::node()] \quad (4)$$

$$/p/a_n::n/a_m::m \equiv /desc::m[b_m::n == /p/a_n::n] \quad (5)$$

$$/a_n::n/a_m::m \equiv /desc::m[b_m::n == /a_n::n] \quad (5a)$$

Equivalence (4) shows that it is possible to remove the first step in a location path within a qualifier. With help of Lemma 1.5 this result is generalized to reverse steps having an arbitrary position within a qualifier.

Figure 3 illustrates the key idea of Equivalence (4), where solid arrows are used for the navigation in the left-hand side of the equivalence, and dashed for the right-hand side: Instead of looking back from the context node specified by path p for matching a certain node ($a_m::m$), one can look forward from the beginning of the document for matching the node ($/desc::m$) and then, still forward, for reaching the initial context node ($b_m::node()$). Hence, e.g. instead of checking whether the context node specified by path p has a preceding m ($p[prec::m]$), one rather looks for an m node and then for a following node that is identical to the context node:

$$p[/desc::m/foll::node() == self::node()] .$$

Equivalence (5) removes the first reverse step from an absolute location path using the same underlying idea.

Note that the equality occurring in these equivalences is based on node identity. The equivalent paths might remain expensive to evaluate, but no evaluation of the $a_m::m$ reverse step is needed anymore.

Example 1. Consider the example of Figure 1 and a query asking for all **names** that appear before a **price**. A way to select these nodes is using the following location path:

$$/desc::price/prec::name .$$

By Equivalence (5a), the previous path can be translated into the following equivalent location path:

$$/desc::name[follow::price == /desc::price] .$$

While the initial location path selects all **name** nodes preceding a **price** node, the equivalent location path selects all **name** nodes, that have a following **price** node, if that node is also a descendant of the root. It is obvious, that there is a considerably simpler equivalent location path (dropping the join), $/desc::name[follow::price]$. The need for the join arises, as the location path selecting the context nodes, relative to which the reverse step is evaluated (in this case the **price** nodes), can be arbitrarily complex:

Consider a slightly modified case of the previous one, where only **prices** that are inside a **journal** with a **title** should be considered. A possible location path for this query with reverse axis is:

$$/desc::journal[child::title]/desc::price/prec::name .$$

Again, by Equivalence (5) this is equivalent to

$$/desc::name[follow::price == /desc::journal[child::title]/desc::price] .$$

In this case it is impossible to remove the introduced join. Note that the join in the first example can be removed by additional equivalence rules for simplifying location paths, rules that are outside the scope of this paper.

Using the equivalences above, it is possible to replace reverse steps in XPath expressions. Nonetheless, in the following section specific equivalences for reverse axes are given, that yield to location paths without joins.

3.2 Specific Equivalences

In this section the interaction of the reverse axes (**anc**, **anc-or-self**, **par**, **prec**, and **prec-sibl**) with forward axes is treated, i.e. equivalences are given, that (if read as rewriting rules from left to right), depending on the location step L_f before a reverse location step L_r , either replace the reverse location step L_r or rewrite the location path into one, where the reverse step L_r is “pushed leftwise”. For every reverse step the interaction with every forward step is shown.

In general, the equivalences have the following structure

$$p/L_f/L_r \equiv p' \quad \text{or} \quad p/L_f[L_r] \equiv p',$$

where p is an absolute path, L_f a forward location step, L_r a reverse location step, and p' the equivalent location path. Sometimes the equivalences can be formulated without the leading path p .

Note that interaction with reverse axes, e.g. interaction of `par` with `prec-sibl`, is not necessary to investigate in these equivalences due to the way our algorithm works (removing reverse steps from left to right of the location path in question). Also, equivalences involving `anc-or-self` and `desc-or-self` are not necessary since these location steps can be replaced using Equivalences (1.6) and (1.7).

Some of the following equivalences do still contain reverse steps on the right-hand side, but these reverse steps are either more on the left of the location path, or the right-hand side is of a form, where other equivalences can be applied to fully remove the reverse location steps as elaborated in Section 4.

Parent. The equivalences in the following proposition are divided in two sets. The first set (Equivalences (6) to (10)) covers the case of `par` location steps outside, the second inside a qualifier. Note that there is a strong structural similarity between the equivalences of the two sets.

Proposition 2 (par axis). *Let m and n be node tests and p a location path. Then the following holds:*

$$\text{desc}::n/\text{par}::m \equiv \text{desc-or-self}::m[\text{child}::n] \quad (6)$$

$$\text{child}::n/\text{par}::m \equiv \text{self}::m[\text{child}::n] \quad (7)$$

$$p/\text{self}::n/\text{par}::m \equiv p[\text{self}::n]/\text{par}::m \quad (8)$$

$$p/\text{foll-sibl}::n/\text{par}::m \equiv p[\text{foll-sibl}::n]/\text{par}::m \quad (9)$$

$$p/\text{foll}::n/\text{par}::m \equiv p/\text{foll}::m[\text{child}::n] \quad (10)$$

$$| p/\text{anc-or-self}::*[\text{foll-sibl}::n]/\text{par}::m$$

$$\text{desc}::n[\text{par}::m] \equiv \text{desc-or-self}::m/\text{child}::n \quad (11)$$

$$\text{child}::n[\text{par}::m] \equiv \text{self}::m/\text{child}::n \quad (12)$$

$$p/\text{self}::n[\text{par}::m] \equiv p[\text{par}::m]/\text{self}::n \quad (13)$$

$$p/\text{foll-sibl}::n[\text{par}::m] \equiv p[\text{par}::m]/\text{foll-sibl}::n \quad (14)$$

$$p/\text{foll}::n[\text{par}::m] \equiv p/\text{foll}::m/\text{child}::n \quad (15)$$

$$| p/\text{anc-or-self}::*[\text{par}::m]/\text{foll-sibl}::n$$

Example 2. Consider the data of Figure 1. The following location path selects all editors of journals:

`/desc::editor[par::journal].`

According to Equivalence (11), this path is equivalent to:

`/desc-or-self::journal/child::editor.`

Ancessor. The following proposition gives equivalences that either move an **anc** step to the left of a path or remove it completely. Equivalences (16a) and (21a) are special cases of Equivalences (16) and (21), respectively.

Proposition 3 (anc axis). *Let m and n be node tests and p a location path. Then the following holds:*

$$p/\text{desc}::n/\text{anc}::m \equiv p[\text{desc}::n]/\text{anc}::m \quad (16)$$

$$| p/\text{desc-or-self}::m[\text{desc}::n]$$

$$/\text{desc}::n/\text{anc}::m \equiv /\text{desc-or-self}::m[\text{desc}::n] \quad (16a)$$

$$p/\text{child}::n/\text{anc}::m \equiv p[\text{child}::n]/\text{anc-or-self}::m \quad (17)$$

$$p/\text{self}::n/\text{anc}::m \equiv p[\text{self}::n]/\text{anc}::m \quad (18)$$

$$p/\text{foll-sibl}::n/\text{anc}::m \equiv p[\text{foll-sibl}::n]/\text{anc}::m \quad (19)$$

$$p/\text{foll}::n/\text{anc}::m \equiv p/\text{foll}::m[\text{desc}::n] \quad (20)$$

$$| p/\text{anc-or-self}::*[\text{foll-sibl}::*/\text{desc-or-self}::n]$$

$$/\text{anc}::m$$

$$p/\text{desc}::n[\text{anc}::m] \equiv p[\text{anc}::m]/\text{desc}::n \quad (21)$$

$$| p/\text{desc-or-self}::m/\text{desc}::n$$

$$/\text{desc}::n[\text{anc}::m] \equiv /\text{desc-or-self}::m/\text{desc}::n \quad (21a)$$

$$p/\text{child}::n[\text{anc}::m] \equiv p[\text{anc-or-self}::m]/\text{child}::n \quad (22)$$

$$p/\text{self}::n[\text{anc}::m] \equiv p[\text{anc}::m]/\text{self}::n \quad (23)$$

$$p/\text{foll-sibl}::n[\text{anc}::m] \equiv p[\text{anc}::m]/\text{foll-sibl}::n \quad (24)$$

$$p/\text{foll}::n[\text{anc}::m] \equiv p/\text{foll}::m/\text{desc}::n \quad (25)$$

$$| p/\text{anc-or-self}::*[\text{anc}::m]$$

$$/\text{foll-sibl}::*/\text{desc-or-self}::n$$

Preceding-sibling. In the following proposition the **prec-sibl** axis is treated. Note that the right-hand side of equivalences for **prec-sibl** (and **prec**) contains more union terms than the other equivalences, since there is no **-or-self** variant of these axes.

Proposition 4 (prec-sibl axis). *Let m and n be node tests and p a location path. Then the following holds:*

$$\text{desc}::n/\text{prec-sibl}::m \equiv \text{desc}::m[\text{foll-sibl}::n] \quad (26)$$

$$\text{child}::n/\text{prec-sibl}::m \equiv \text{child}::m[\text{foll-sibl}::n] \quad (27)$$

$$p/\text{self}::n/\text{prec-sibl}::m \equiv p[\text{self}::n]/\text{prec-sibl}::m \quad (28)$$

$$p/\text{foll-sibl}::n/\text{prec-sibl}::m \equiv p[\text{self}::m/\text{foll-sibl}::n] \quad (29)$$

$$| p[\text{foll-sibl}::n]/\text{prec-sibl}::m$$

$$| p/\text{foll-sibl}::m[\text{foll-sibl}::n]$$

$$p/\text{foll}::n/\text{prec-sibl}::m \equiv p/\text{foll}::m[\text{foll-sibl}::n] \quad (30)$$

$$| p/\text{anc-or-self}::*[\text{foll-sibl}::n]$$

$$/\text{prec-sibl}::m$$

$$| p/\text{anc-or-self}::m[\text{foll-sibl}::n]$$

$$\text{desc}::n[\text{prec-sibl}::m] \equiv \text{desc}::m/\text{foll-sibl}::n \quad (31)$$

$$\text{child}::n[\text{prec-sibl}::m] \equiv \text{child}::m/\text{foll-sibl}::n \quad (32)$$

$$p/\text{self}::n[\text{prec-sibl}::m] \equiv p[\text{self}::n]/\text{foll-sibl}::m \quad (33)$$

$$p/\text{foll-sibl}::n[\text{prec-sibl}::m] \equiv p[\text{self}::m]/\text{foll-sibl}::n \quad (34)$$

$$| p/\text{foll-sibl}::m/\text{foll-sibl}::n$$

$$| p[\text{prec-sibl}::m]/\text{foll-sibl}::n$$

$$p/\text{foll}::n[\text{prec-sibl}::m] \equiv p/\text{foll}::m/\text{foll-sibl}::n \quad (35)$$

$$| p/\text{anc-or-self}::*[\text{prec-sibl}::m]$$

$$/\text{foll-sibl}::n$$

$$| p/\text{anc-or-self}::/\text{foll-sibl}::n$$

Preceding. The following proposition describes the interaction of *prec* with other axes.

Proposition 5 (prec axis). *Let m and n be node tests and p a location path. Then the following holds:*

$$p/\text{desc}::n/\text{prec}::m \equiv p[\text{desc}::n]/\text{prec}::m \quad (36)$$

$$| p/\text{child}::*[\text{foll-sibl}::*/\text{desc-or-self}::n]$$

$$/\text{desc-or-self}::m$$

$$/\text{desc}::n/\text{prec}::m \equiv / \text{desc}::m[\text{foll}::n] \quad (36a)$$

$$p/\text{child}::n/\text{prec}::m \equiv p[\text{child}::n]/\text{prec}::m \quad (37)$$

$$| p/\text{child}::*[\text{foll-sibl}::n]/\text{desc-or-self}::m$$

$$p/\text{self}::n/\text{prec}::m \equiv p[\text{self}::n]/\text{prec}::m \quad (38)$$

$$p/\text{foll-sibl}::n/\text{prec}::m \equiv p[\text{foll-sibl}::n]/\text{prec}::m \quad (39)$$

$$| p/\text{foll-sibl}::*[\text{foll-sibl}::n]/\text{desc-or-self}::m$$

$$| p[\text{foll-sibl}::n]/\text{desc-or-self}::m$$

$$p/\text{foll}::n/\text{prec}::m \equiv p[\text{foll}::n]/\text{prec}::m \quad (40)$$

$$| p/\text{foll}::m[\text{foll}::n]$$

$$| p[\text{foll}::n]/\text{desc-or-self}::m$$

$$\begin{aligned}
p/\text{desc}::n[\text{prec}::m] &\equiv p[\text{prec}::m]/\text{desc}::n & (41) \\
&| p/\text{child}::*[\text{desc-or-self}::m] \\
&| \text{foll-sibl}::*/\text{desc-or-self}::n
\end{aligned}$$

$$/\text{desc}::n[\text{prec}::m] \equiv / \text{desc}::m/\text{foll}::n \quad (41a)$$

$$p/\text{child}::n[\text{prec}::m] \equiv p[\text{prec}::m]/\text{child}::n \quad (42)$$

$$| p/\text{child}::*[\text{desc-or-self}::m]/\text{foll-sibl}::n$$

$$p/\text{self}::n[\text{prec}::m] \equiv p[\text{prec}::m]/\text{self}::n \quad (43)$$

$$p/\text{foll-sibl}::n[\text{prec}::m] \equiv p[\text{prec}::m]/\text{foll-sibl}::n \quad (44)$$

$$| p/\text{foll-sibl}::*[\text{desc-or-self}::m]/\text{foll-sibl}::n$$

$$| p[\text{desc-or-self}::m]/\text{foll-sibl}::n$$

$$p/\text{foll}::n[\text{prec}::m] \equiv p[\text{prec}::m]/\text{foll}::n \quad (45)$$

$$| p/\text{foll}::m/\text{foll}::n$$

$$| p[\text{desc-or-self}::m]/\text{foll}::n$$

Example 3. Consider the location path

$$/\text{desc}::\text{price}/\text{prec}::\text{name}$$

of Example (1). With Rule 36a it can be rewritten to

$$/\text{desc}::\text{name}[\text{foll}::\text{price}].$$

This result is more compact and closer to the original than the result of Example (1) using Equivalence (5a).

4 Location Path Rewriting

Each Equivalence (i) $p_1 \equiv p_2$ of Section 3 gives rise to a rewriting rule: A path matching with the left-hand side p_1 can be rewritten into a path corresponding to the right-hand side p_2 . In the following, Rule (i) denotes the rewriting rule $p_1 \rightarrow p_2$ induced by Equivalence (i) $p_1 \equiv p_2$. The equivalences of Lemma 1 induce rewriting rules, denoted Rules (1.1) to (1.8).

The equivalences of Section 3 are split in two sets of rules for use in a rewriting algorithm:

1. RuleSet₁, containing Rules (1) to (3) and the general Rules (4) to (5a).
2. RuleSet₂, containing Rules (1) to (3) and the specific Rules (3) to (42).

A rule can be applied to a location path in the following manner:

Definition 2 (Rule application). *Let p be a non-disjunctive location path, and let $p_l \rightarrow p_r$ be a rule either from RuleSet₁ or RuleSet₂. If p is of the form p_l/p' , then let q denote the path p_r/p' . If p_l is a relative path and if p is of the form $p_1/p_l/p_2$, then let q denote the path $p_1/p_r/p_2$. In both cases q is called the result of the application of rule $p_l \rightarrow p_r$ to p .*

An algorithm, called “rare” (sketched in Figure 4) for computing a reverse-axis-free path equivalent to an absolute path is considered below. The input for the algorithm is restricted to paths without qualifiers containing “RR joins”:

Definition 3 (RR join). *An RR join is an XPath expression of the form $p_1 \theta p_2$ where $\theta \in \{==, =\}$, and both p_1 and p_2 are **R**elative paths such that at least one of them contains a **R**everse step.*

For the consideration of termination and correctness of the algorithm, some important properties of the application of the rewriting rules to a location path are required:

Lemma 3 (Properties of rule application). *Let p be an absolute location path with no qualifier containing RR joins.*

1. *If p contains a reverse step, then a rule from RuleSet_1 and a rule from RuleSet_2 is applicable to p . Possibly, Rules (1.1) to (1.8) have to be applied first.*
2. *The result of a rule application to the first reverse step in p is an absolute path with no qualifiers containing RR joins.*
3. *If q is the result of a rule application to p , then $p \equiv q$.*

Proof. (1): Let L be the first reverse location step.

First consider RuleSet_1 : If L occurs outside a qualifier, Rules (5), (5a) or (1) to (3) can be applied, since p is an absolute location path. If L occurs as the first location step inside a qualifier, Rule (4) can be applied. If L appears at any other position inside a qualifier, Rule (1.5) can be applied in order to construct a qualifier with L as first location step. Rule (4) can be applied now.

RuleSet_2 provides rules for interaction between each reverse step and an arbitrary forward step s , there is always a rule, that can be applied to the first reverse step in p .

(2) Only Rules (4), (5), and (5a) introduce a binary relation (namely $==$), if they are applied to a location path. But always one of the operands is an absolute path. Hence, in any case the result of the rule application contains no RR join. Furthermore, since p is an absolute path, the result of applying a rule to p is also an absolute path.

(3) This holds due to Lemma 1.1–4.

“rare” Algorithm. The “rare” Algorithm, outlined in Figure 4, can be used for RuleSet_1 as well as for RuleSet_2 . The algorithm takes as input a location path which is absolute, since some rules from RuleSet_1 and RuleSet_2 are only applicable to absolute location paths.

Theorem 1 (Removal of reverse location steps using RuleSet_1). *Let p be an absolute path without qualifiers in which RR joins occur. There exists an absolute path p' with no reverse steps such that $p \equiv p'$. Using “rare” and RuleSet_1 , this path p' has a length and can be computed in a time linear in the length of p .*

Let $\xi = \text{RuleSet}_1$ or RuleSet_2 .

Auxiliary functions:

match(p): returns the result of a rule application from ξ to the first reverse location step in p .

apply-lemmas(p): returns p if Rules (1.1-8) are not applicable to p . Otherwise, it returns the result of the repeated application of Rules (1.1-8) to p .

union-flattening(p): returns a path equivalent to p with unions at top level only.

rare(p)

Input: p {absolute location path without qualifiers containing RR joins}.

$p \leftarrow \text{apply-lemmas}(p)$.

$p \leftarrow \text{union-flattening}(p) = U_1 \mid \dots \mid U_n$ ($n \geq 1$).

$S \leftarrow$ empty stack.

for $i \leftarrow 1$ to n **do**

$\text{push}(U_i, S)$.

end for

$p' \leftarrow \perp$. {initialization}

while $\text{not}(\text{empty}(S))$ **do**

$U \leftarrow \text{pop}(S)$.

while U contains a reverse step **do**

$U \leftarrow \text{match}(U)$.

$U \leftarrow \text{apply-lemmas}(U)$.

$U \leftarrow \text{union-flattening}(U) = V_1 \mid \dots \mid V_n$ ($n \geq 1$).

for $i \leftarrow 2$ to n **do**

$\text{push}(V_i, S)$.

end for

$U \leftarrow V_1$.

end while

$p' \leftarrow p' \mid U$.

end while

Output: p' {location path without reverse axes equivalent to p }.

Fig. 4. Algorithm *rare* (reverse axis removal)

Proof. A path equivalent to p is constructed as sketched in Figure 4. All reverse location steps are rewritten, one after another. Lemma 3 guarantees that a rule of RuleSet_1 can be applied to any path containing a reverse location step. The resulting path p' contains no reverse location steps and is equivalent to p .

The location path p' is of linear size and constructed in linear time, since each rule application removes one reverse step, adds at most two forward location steps and no reverse ones.

Theorem 2 (Removal of reverse location steps using RuleSet_2). *Let p be an absolute path with no qualifiers in which RR joins occur. There exists an absolute path p' with no reverse steps such that $p \equiv p'$. Using “rare” and*

RuleSet₂, this path p' has a length and can be computed in a time exponential in the length of p .

Proof. An application of a rule from *RuleSet₂* can have three different result types:

1. removes completely a reverse step, e.g. Rules (1) to (3)) or (6);
2. pushes the reverse step from right to left in the path, e.g. Rule (8);
3. for the interaction between a `fo:ll` and a reverse step L_r , i.e. `fo:ll::n/Lr` or `fo:ll::n[Lr]`, a union of several other paths is obtained, e.g. Rule (10); the resulting union terms have reverse steps at positions less than or equal in the original path and they do not contain anymore the interaction between the initial `fo:ll` step and a reverse step.

Since the path has a finite length, the procedure of pushing reverse steps leftwise terminates. Also, the number of interactions between `fo:ll` and reverse steps is finite. Hence, the algorithm terminates.

Each rule application having the first result type removes a reverse step and does not change the number of union terms. Hence, in the best case, i.e. using only rule applications with the first result type, the algorithm has a linear time complexity in the length of the input path p .

The last two result types are significant for the worst-case complexity of the algorithm, since each rule application can produce intermediate rewritten paths with more than one union term (up to three union terms). Hence, each rule application can increase the order of the input for the next rule application, yielding an exponential time complexity in the length of the input path p .

Example runs of the algorithm for both set of rules are presented in Figure 5.

Comparison. Both *RuleSet₁* and *RuleSet₂* have advantages and it is still an open issue which one is preferable. The path rewriting using *RuleSet₂* has in the worst case an exponential time complexity and output size in the length of the input location path. As location paths are in practice small (less than ten steps), the exponential worst-case complexity of *RuleSet₂* does not necessarily generate longer paths than *RuleSet₁*. In addition, since they do not contain joins, the location paths generated using *RuleSet₂* are simpler (as can be seen in the examples), hence more convenient to evaluate, than those generated using *RuleSet₁*, which contain the same number of joins as there are reverse steps in the input location path.

For further comparison, practical tests of the rewritten location paths, as generated by the two rule sets, have been performed using a Java prototype of a streamed XPath processor [21] against the MONDIAL geographical database [22], a highly structured XML document of comparatively small size.

Three types of location paths with reverse steps have been considered for comparing the efficiency of the rule sets:

Q1: simple location paths without qualifiers, e.g.

`/desc::name/anc::country.`

Consider the example of Figure 1 and a query asking for all titles that appear before a name and are inside journals. This query can be expressed as the following location path:

`/desc::name/prec::title[anc::journal]`

Note that p is an absolute path without qualifiers containing RR-joins.

RuleSet₁

- Step 1:** $p \leftarrow \text{apply-lemmas}(p) = p.$
- Step 2:** $U_1 \leftarrow \text{/desc::name/prec::title[anc::journal].}$
- Step 3:** $\text{push}(U_1, S).$
- Step 4:** $p' \leftarrow \perp, U \leftarrow \text{pop}(S).$
- Step 5:** U contains a reverse step (`prec::title`).
- Step 6:** $U \leftarrow \text{match}(U) =$
 $\text{/desc::title[anc::journal]}$
 $[\text{foll::name} == \text{/desc::name}] \text{ {Rule (5)}}$
- Step 7:** $U \leftarrow \text{apply-lemmas}(U) = U.$
- Step 8:** U contains a reverse step (`anc::journal`).
- Step 9:** $U \leftarrow \text{match}(U) =$
 $\text{/desc::title[/desc::journal/desc::node() == self::node()]}$
 $[\text{foll::name} == \text{/desc::name}]. \text{ {Rule (4)}}$
- Step 10:** $U \leftarrow \text{apply-lemmas}(U) = U.$
- Step 11:** U does not contain reverse steps.
- Step 12:** $p' \leftarrow U, S$ is empty.
- Output:** $p = \text{/desc::title[/desc::journal/desc::node() == self::node()]}$
 $[\text{foll::name} == \text{/desc::name}].$

RuleSet₂

- Step 1:** $p \leftarrow \text{apply-lemmas}(p) = p.$
- Step 2:** $U_1 \leftarrow \text{/desc::name/prec::title[anc::journal].}$
- Step 3:** $\text{push}(U_1, S).$
- Step 4:** $p' \leftarrow \perp, U \leftarrow \text{pop}(S).$
- Step 5:** U contains a reverse step (`prec::title`).
- Step 6:** $U \leftarrow \text{match}(U) =$
 $\text{/desc-or-self::title[foll::name] [anc::journal]} \text{ {Rule (36a)}}$
- Step 7:** $U \leftarrow \text{apply-lemmas}(U) = \text{/desc::title[foll::name] [anc::journal]}$
 $|\text{/self::title[foll::name] [anc::journal].}$
- Step 8:** U contains a reverse step (`anc::journal`).
- Step 9:** $U \leftarrow \text{match}(U) =$
 $\text{/desc::title[foll::name] [anc::journal] | \perp. \text{ {Rules (1) to (3)}}$
- Step 10:** U contains a reverse step (`anc::journal`).
- Step 11:** $U \leftarrow \text{match}(U) =$
 $\text{/desc::journal/desc::title[foll::name]}. \text{ {Rule (21a)}}$
- Step 12:** $U \leftarrow \text{apply-lemmas}(U) = U.$
- Step 13:** U does not contain reverse steps.
- Step 14:** $p' \leftarrow U, S$ is empty.
- Output:** $p' = \text{/desc::journal/desc::title[foll::name]}.$

Fig. 5. Example runs of *rare* algorithm

Q2: location paths with a single qualifier, e.g.

`/desc::name/par::city[anc::country] .`

Q3: complex location paths with several qualifiers, potentially rewritten by the second RuleSet₂ to location paths with exponential size, e.g.

`/desc::province[par::country]/child::city[anc::mondial] .`

The queries are rewritten using RuleSet₁ and RuleSet₂ respectively. Figure 6 shows the results of processing the given queries against the MONDIAL database on a Pentium III 1 GHz, 512 MB system running SuSE Linux 7.3.

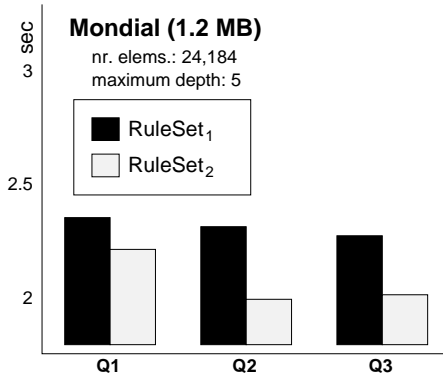


Fig. 6. Experimental evaluation of the rule sets

The results show that the two rule sets are very competitive for all three types of queries. As expected, RuleSet₂ generating simpler rewritten queries performs slightly better than RuleSet₁.

Rewriting location paths using variables. There are two classes of location paths not covered by the rules given so far: relative location paths and location paths with RR joins (cf. Definition 3), e.g. $p[\text{self}::* = \text{prec}::*]$. Any attempt to remove the reverse location steps in these cases results in losing the context given by p .

In the full version [19] of this paper an approach is proposed to solve this problem by remembering the context in a variable. It is based on a **for-return** construct for variable binding, as provided by XPath 2.0, XQuery, and XSLT. Using this approach *every* location path can be rewritten to an equivalent reverse-axis-free one.

5 Related Work

Several methods have been proposed for rewriting XPath expressions taking integrity constraints or schemas into account [23,24], and the equivalence and containment problems for XPath expressions have been investigated [25,26]. Furthermore, a growing interest in query optimization for XML databases, including optimization of XPath expressions, recently emerged. To the best of our knowledge, however, no other approach has been proposed for removing reverse steps from XPath expressions relying upon XPath symmetry. Note that using

equivalence preserving rewriting rules for removing reverse steps from XPath expressions, as it is proposed in the present paper, is not closely related to the general equivalence problem for XPath expressions.

In [27] redundancies in XPath expressions based on a “model-oriented” approach are investigated. Such an approach relies on an abstract model of XPath that views XPath expressions as tree patterns. [27] shows that redundant branches of a tree pattern can be eliminated in polynomial time. Tree patterns are more abstract than XPath expressions in a way which is relevant to the work described in the present paper: A same tree pattern represents multiple equivalent XPath expressions. In particular, the symmetries in XPath exploited in the present paper are absent from tree patterns. Tree patterns do not consider the document order and therefore the concept of forward and reverse steps. In some sense the present work shows in which cases this simplified view upon an XPath expression can be justified.

Stream-based query processing has gained considerable interest in the past few years, e.g. due to its application in data integration [13,12] and in publish-subscribe architectures [11,10]. They all consider a navigational approach (XML-QL or XPath) consisting of a restricted subset of forward axes from XPath. This fact contrasts with the present work, which enables the use of the unrestricted set of XPath axes in a stream-based context.

6 Conclusion

The main result of this paper consists in two rule sets, RuleSet_1 and RuleSet_2 , used in an algorithm for transforming XPath 1.0 expressions containing reverse axes into reverse-axis-free equivalents. Both RuleSet_1 and RuleSet_2 have advantages and it is still an open issue which one is preferable. The location paths generated using RuleSet_2 do not contain joins, in contrast with those generated using RuleSet_1 , which contain the same number of joins as there are reverse steps in the input location path. However, the path rewriting using RuleSet_2 has an exponential complexity in the length of the input location path, in contrast with rewriting using RuleSet_1 which has only a linear complexity. Preliminary experimental results show that the rewritten queries generated by RuleSet_2 are evaluated slightly faster than those generated by RuleSet_1 .

Closely related to the comparison of RuleSet_1 and RuleSet_2 we plan to investigate the notion of “minimality” or “simplicity” of XPath expressions. We are focusing on defining a notion of a minimal XPath expression that can be evaluated more efficiently in a stream-based context than its equivalents. A notion of minimality will allow for well-founded optimization techniques for XPath expressions.

The equivalences proposed in this paper are drawn from and represent prerequisites for an efficient streamed evaluation of unrestricted XPath, as considered in [21].

References

- [1] W3C, “XQuery 1.0: An XML query language,” W3C Working Draft, 2002.
- [2] W3C, “XSL Transformations (XSLT) Version 1.0,” W3C Recommendation, 1999.
- [3] W3C, “XML Path Language (XPath) Version 1.0,” W3C Recommendation, 1999.
- [4] A. Desai, “Introduction to Sequential XPath,” in *Proc. of IDEAlliance XML Conference*, 2001.
- [5] W3C, “XSL Transformations (XSLT) Version 2.0,” W3C Working Draft, 2002.
- [6] W3C, “Document Object Model (DOM) Level 2 Core Specification,” W3C Recommendation, 2000.
- [7] N. Ide, P. Bonhomme, and L. Romary, “XCES: An XML-based standard for linguistic corpora,” in *Proc. of the Second Annual Conference on Language Resources and Evaluation*, 2000.
- [8] F. Bry and P. Krüger, “A Computational Biology Database Digest: Data, Data Analysis, and Data Management,” Tech. Rep. PMS-FB-2002-8, University of Munich, 2002.
- [9] “Astronomical Data Center,” homepage <http://adc.gsfc.nasa.gov>.
- [10] C. Chan, P. Felber, M. Garofalakis, and R. Rastogi, “Efficient Filtering of XML Documents with XPath Expressions,” in *Proc. of International Conference on Data Engineering (ICDE)*, 2002.
- [11] M. Altinel and M. Franklin, “Efficient Filtering of XML Documents for Selective Dissemination of Information,” in *Proc. of 26th Conference on Very Large Databases (VLDB)*, 2000.
- [12] A. Levy, Z. Ives, and D. Weld, “Efficient Evaluation of Regular Path Expressions on Streaming XML Data,” Tech. Rep., University of Washington, 2000.
- [13] T. J. Green, M. Onizuka, and D. Suciu, “Processing XML Streams with Deterministic Automata and Stream Indexes,” Tech. Rep., University of Washington, 2001.
- [14] Apache Project, “Cocoon 2.0: XML publishing framework,” available at <http://xml.apache.org/cocoon/index.html>.
- [15] Apache Project, “Xalan-Java Version 2.2,” available at <http://xml.apache.org/xalan-j/index.html>.
- [16] D. Megginson, “SAX: The Simple API for XML,” 1998.
- [17] P. Wadler, “A formal semantics of patterns in XSLT,” in *Proc. of Conference on Markup Technologies*, 1999.
- [18] P. Wadler, “Two semantics of XPath,” Tech. Rep., 2000.
- [19] D. Olteanu, H. Meuss, T. Furche, and F. Bry, “XPath: Looking Forward,” Tech. Rep. PMS-FB-2001-17, University of Munich, 2001.
- [20] W3C, “XQuery 1.0 and XPath 2.0 data model,” W3C Working Draft, 2001.
- [21] “XPath Evaluation Project,” University of Munich, homepage <http://www.pms.informatik.uni-muenchen.de/forschung/xpath-eval.html>.
- [22] W. May, “Information Extraction and Integration with FLORID: The MONDIAL Case Study,” Tech. Rep. 131, University of Freiburg, Institut für Computer Science, 1999, available at www.informatik.uni-freiburg.de/~may/Mondial/.

- [23] K. Boehm, K. Gayer, T. Oezsu, and K. Aberer, "Query Optimization for Structured Documents Based on Knowledge on the Document Type Definition," in *Proc. of the Advances in Digital Libraries Conference*, 1998.
- [24] P. T. Wood, "Optimising Web Queries Using Document Type Definitions," in *2nd ACM Workshop on Web Information and Data Management (WIDM'99)*, 1999.
- [25] A. Deutsch and V. Tannen, "Containment for Classes of XPath Expressions Under Integrity Constraints," in *Knowledge Representation meets Databases (KRDB)*, 2001.
- [26] P. T. Wood, "On the Equivalence of XML Patterns," in *Proc. 6th Int. Conf. on Rules and Objects in Databases (DOOD)*, 2000.
- [27] Sihem Amer-Yahia, SungRan Cho, Laks V. S. Lakshmanan, and Divesh Srivastava, "Minimization of Tree Pattern Queries," in *SIGMOD*, 2001.