

# Optimizing XML querying using type-based document projection

VÉRONIQUE BENZAKEN, GIUSEPPE CASTAGNA, DARIO COLAZZO, and KIM NGUYỄN

XML data projection (or pruning) is a natural optimization for main memory query engines: given a query  $Q$  over a document  $D$ , the subtrees of  $D$  that are not necessary to evaluate  $Q$  are pruned, thus producing a smaller document  $D'$ ; the query  $Q$  is then executed on  $D'$ , hence avoiding to allocate and process nodes that will never be reached by  $Q$ .

In this article, we propose a new approach, based on types, that greatly improves current solutions. Besides providing comparable or greater precision and far lesser pruning overhead, our solution—unlike current approaches—takes into account backward axes, predicates, and can be applied to multiple queries rather than just to single ones. A side contribution is a new type system for XPath able to handle backward axes. The soundness of our approach is formally proved. Furthermore, we prove that the approach is also complete (i.e., yields the best possible type-driven pruning) for a relevant class of queries and Schemas. We further validate our approach using the XMark and XPathMark benchmarks and show that pruning not only improves the main memory query engine's performances (as expected) but also those of state of the art native XML databases.

Categories and Subject Descriptors: H.2.4 [DATABASE MANAGEMENT]: Systems—*Query Processing*; F.3.2 [LOGICS AND MEANINGS OF PROGRAMS]: Semantics of Programming Languages—*Program analysis*

General Terms: Languages, Theory

Additional Key Words and Phrases: XML, Type projection, Document pruning, Query optimization, Type systems

## 1. INTRODUCTION

Main-memory XML query engines are often the primary choice for applications that do not wish or cannot afford to build secondary storage indexes or load a database before query processing. One of the main optimisation techniques recently adopted in this context is XML data projection (or pruning) [Marian and Siméon 2003; Bressan et al. 2005].

The idea behind document projection is very simple yet very powerful. Given a query  $Q$  and a document  $D$ , sub-trees of  $D$  that are not necessary to evaluate  $Q$  are pruned, thus yielding a smaller document  $D'$ . Then  $Q$  is executed over  $D'$ , hence avoiding to allocate and process nodes that will never be reached by navigational specifications in  $Q$ . This ensures that evaluation over  $D'$  is equivalent to and more efficient than the evaluation over  $D$ .

As shown in [Marian and Siméon 2003; Bressan et al. 2005], XML navigation specifications expressed in queries tend to be very selective, especially in terms of document structure. Therefore, pruning may yield significant improvements both in terms of execution time and in terms of memory usage: as a matter of facts, for main-memory XML query engines, very large documents can not be queried without pruning.

### 1.1. State of the art

Marian and Siméon [Marian and Siméon 2003] propose that the actual data-needs of an XQuery query  $Q$  (that is, the part of data that is necessary to the execution of the query) is determined by

---

Authors' addresses: V. Benzaken, D. Colazzo, and K. Nguyễn, LRI, Bât 650, Université Paris Sud, 91405 Orsay Cedex, France. G. Castagna, PPS, Université Paris Diderot, rue Albert Einstein, 75013 Paris, France.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© YYYY ACM 0362-5915/YYYY/01-ARTA \$15.00

DOI 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

Véronique Benzaken, Giuseppe Castagna, Dario Colazzo, and Kim Nguyễn. Optimizing XML querying using type-based document projection. *ACM Transactions on Database Systems*, TO APPEAR.

statically extracting all paths in  $Q$ . These paths are then applied to  $D$  at load time, in a streaming fashion, in order to prune unneeded parts of data. The technique is powerful since: (i) it applies to most of XQuery *core*, (ii) it can be applied to a set of queries over the same document, and (iii) it does not require any *a priori* knowledge of the structure of  $D$ . However, this technique suffers some limitations. First, the document loader-pruner can manage neither *backward axes* nor path expressions with predicates which, especially the latter, can contain precious information to optimise pruning. Second, the advantage described in point (iii) becomes a big drawback when “//” occurs in paths since, in that case, the technique does not behave efficiently in terms of loading time and pruning precision (hence, memory allocation). Indeed, when a // is present in a projection path, the pruning process requires to visit all descendants of a node in order to decide whether the node contains a useful descendant. Worst, pruning time tends to be quite high and it drastically increases (together with memory consumption) when the number of // augments in the pruning path-set. In their technique, pruning actually corresponds to computing a further query whose time and memory occupation may be comparable to those required to compute the original query. In particular, in this technique every occurrence of // may yield a full exploration of the tree (e.g., see in [Marian and Siméon 2003] the test for the XMark [Schmidt et al. 2002] query Q7 which only contains three // steps and for which just computing the pruning takes longer than executing the query on the original document). Therefore, pruning execution overhead and its high memory footprint may jeopardise the gains obtained by using the pruned document. Third and finally, the precision of pruning drastically degrades (down to being nullified) for queries containing the XPath expressions `descendant::node[cond]`, which are often used in practice.

Motivated by efficient XML stream processing, Green *et al.* [Green et al. 2004] introduced a framework for discarding sequences of SAX events in an XML data stream. Although their approach allows them to prune an input stream with respect to sets of queries, the language they handle is restricted to forward linear XPath expressions (that is, XPath expressions with only *child* and *descendant* axes and without predicates).

Bressan *et al.* [Bressan et al. 2005] introduce a different and quite precise XML pruning technique for a subset of XQuery FLWR expressions. The technique is based on the *a priori* knowledge of a data-guide for  $D$ . The document  $D$  is first matched against an abstract representation of  $Q$ . Pruning is then performed at run time, it is very precise, and, thanks to the use of some indexes over the data-guide, it ensures good improvements in terms of query execution time. However, the technique is one-query oriented—in the sense that it cannot be applied to multiple queries—it does not handle XPath predicates, and cannot handle backward axes (recall that the encodings of [Olteanu et al. 2002] are defined for XPath, and no extension to XQuery-like languages is known). Also, the approach requires the construction and management of the data-guide and of adequate indexes.

Improving the work by [Marian and Siméon 2003], [Benzaken et al. 2006] proposes a type-based approach to document pruning where a query is run over a DTD to statically determine the query need. Unlike [Marian and Siméon 2003], their approach remains precise and efficient in the presence of backward axes. It is, however, limited to DTD’s and uses approximations for nested or conjunctive predicates which jeopardize the precision of their pruning algorithm.

## 1.2. Our contribution

In this article, we improve and extend the work of [Benzaken et al. 2006] in several important aspects. First and foremost we generalized their definition of type projectors by using as projectors sets of production rules (as opposed to the sets of non-terminals used in [Benzaken et al. 2006]) of regular tree grammars (as opposed to the DTDs used in [Benzaken et al. 2006]). This generalization was far from being straightforward. In particular, we had to prove the applicability of our technique to the more general framework under consideration (cf. Section 3.2). However the result is worth the effort since the advantages of this generalization are twofold. On the one hand using regular tree grammars allows us to compute type projectors for every kind of XML schema formalism we are aware of as, for instance, DTDs, XMLSchemas, CDuce and XDuce types, Relax-Core and TREX schemas. On the other hand, inferring grammar production rules rather than grammar non-terminals

allows us to compute context-aware and, thus, more precise projectors. More precisely, the new type projectors introduced in this work can prune a subtree not only based on its tag (as it was done in [Benzaken et al. 2006]), but also on any structural condition expressible by a regular tree language. So for instance our pruning process may decide to prune just one of two trees generated by the same non-terminal, because they appear in different contexts (in [Benzaken et al. 2006] either both trees were pruned or they were both preserved). Therefore these new projectors are both more general and can perform much finer-grained pruning. Second, although we develop the theory of type projection for a simplified data-model and restricted forms of XPath expressions, we thoroughly detail how to tackle many of the peculiarities of the XML [XML 2008] and XPath [XPath 1999] specifications, including the handling of attributes, the presence of absolute axes in XPath predicates or a wide range of predefined XPath functions (all absent in [Benzaken et al. 2006]). The path language we formally study extends the one in [Benzaken et al. 2006] with top-level unions of paths, predicate conjunctions (“and”) and arbitrarily nested predicates (their previous work formally treated only non-nested predicates and resorted to an approximation in the case of nested predicates). Third, we provide an extensive list of experiments showing the overall benefits of type projection for a wide range of queries and query engines. These experiments supersede the early benchmarks realised in [Benzaken et al. 2006] and show that despite the advances in XML query technologies in the recent years, our static analysis can significantly improve the performances (both in time and memory consumption) of many different XML query engines.

Our technique combines the advantages of the previously mentioned works while relaxing their limitations. Unlike [Marian and Siméon 2003; Bressan et al. 2005; Green et al. 2004], our approach accounts for backward axes, performs a fine-grained analysis of predicates, allows (unlike [Bressan et al. 2005]) for dealing with bunches of queries, and (unlike [Marian and Siméon 2003]) cannot be jeopardised by pruning overhead. Our solution provides in all cases comparable or greater precision than the other approaches, while it requires always negligible or no pruning overhead. Moreover, contrary to [Marian and Siméon 2003; Bressan et al. 2005], our approach is formally proven to be *sound* (i.e., pruning does not alter the result of queries) and, furthermore, is also *complete* (i.e., it produces the best possible type-driven pruning) for a substantial class of queries and DTDs.

We introduce our framework in three steps. In the first step, we consider a simplified version of XPath, dubbed XPath<sup>ℓ</sup>, which correspond to *structural* queries (that is path expressions whose (possibly nested) predicates only contains disjunctions and conjunctions of paths). We define for XPath<sup>ℓ</sup> a static analysis that determines a set of typing rules, a *type projector*, that is then used to prune the document(s). One of the particular features of this approach is that our pruning algorithm is characterised by a constant (and low) memory consumption and by an execution time linear in the size of the document to prune. More precisely, a pruning based on type projectors is equivalent to a single buffer-less one-pass traversal of the parsed document (it simply discards elements not generated by any of the rules in the projector). So if embedded in query processors, pruning can be executed during parsing and/or validation and brings no overhead at all, while if used as an external tool it requires a time always smaller than or equal to the time used to parse the queried document. Soundness and (partial) completeness results for the static analysis are stated.

The second step consists of extending the analysis to full XPath (more precisely, to XPath 1.0), that is, we need to show how to deal with missing axes as well as the *Core Library Functions* that may occur in predicates. This is done by associating to each XPath query  $Q$  a XPath<sup>ℓ</sup> query  $P$  that soundly approximates  $Q$ , in the sense that the projector inferred for  $P$  by the static analysis developed at the first step is also a sound projector for  $Q$ .

The final step of our process is to extend the approach to XQuery (hence, to XPath 2.0). This is obtained in the same way as done in [Marian and Siméon 2003], by defining a path extraction algorithm. Our path extraction algorithm improves and extends in several aspects (in particular, in terms of extracted paths’ selectivity) the one of [Marian and Siméon 2003]. It also computes the XPath<sup>ℓ</sup> approximation of the extracted paths so that the static analysis of the first step can be directly applied to them.

We prove some important closure properties that guarantee that type projections can always be

performed at load time during the validation process, and this without any overhead. In particular for XML documents typed with DTDs or XML Schemas the document can be pruned in streaming.

We gauged and validated our approach by testing it both on the XPathMark [Franceschet 2005] and on the XMark [Schmidt et al. 2002] benchmarks. The result of this validation confirmed what was expected: thanks to the handling of backward axes and of predicates the precision of our pruning is in general noticeably higher than that of current approaches; the pruning time is linear in the size of the queried document and has a very low memory footprint; the time of the static analysis is always negligible (lower than half a second on the hardware we used for our benchmarks described in Section 9) even for complex queries and DTDs. But benchmarks also brought unexpected (and quite pleasant) results. In particular, they showed that type-based pruning brings benefits that go beyond those of the reduced size of the pruned document: by excluding a whole set of data structures (those whose type names are not included in the type projector), the pruning may drastically reduce the resources that must be allocated at run-time by the query processor. For instance, our benchmarks show that for several XMark and XPathMark queries our pruning yields a document whose size is two thirds of the size of the original document, but the query can then be processed using three times less memory than when processed on the original document. This is a very important gain, especially for DOM-based processors, or memory sensitive processors. Not only our approach is relevant in the case of main memory query engines such as Saxon but it is also shown to be useful for native query engines as efficient as MonetDB [Boncz et al. 2006]. Even in the latter case our experiments demonstrate the relevance of type projection as a complementary optimisation technique. Indeed, this is not totally surprising as type projection can be thought of as a way of defining clustering policies in the same line as what was done in the context of object-oriented databases [Benzaken et al. 1992; Benzaken 1990; Benzaken and Delobel 1990]. Clustering and indexing are well-known complementary tools used in the context of query optimisation.

While the present work tackles the problem of query optimization by the mean of type projection, it also features a more general and directly reusable result: the definition of a new type system for XPath. In particular, this type system is able to handle *precisely* backward axes. This, on its own, constitutes a contribution of this work. In particular the precision of type inference for backward axes goes beyond what is proposed in the XQuery Static Semantic recommendation ([Draper et al. 2007], which essentially leaves the result of a backward step untyped).

### 1.3. Plan of the article

The article is organised as follows. Section 2 introduces basic definitions and notations: data model, types, validation. Section 3 presents type projectors, type-based projection, and several theoretical (closure) properties. In Section 4 we define XPath<sup>ℓ</sup> and its semantics. In Section 5 we present our type projectors inference algorithm for XPath<sup>ℓ</sup> and state its formal properties. In Section 6 we extend our approach to full XPath and in Section 7 to XQuery. In Section 8 we discuss how to apply our technique to other typing policies as well as to un-typed documents. Section 9 presents our implementation and reports the results of our benchmarks. We finally conclude in Section 10 by presenting the perspectives of this work. All the proofs for the stated results can be found in the electronic appendix that can be accessed in the ACM Digital Library. .

## 2. NOTATIONS

### 2.1. Data Model

For the sake of concision and clarity we present our solution for a simplified version of the XQuery data model where we do not consider node attributes. However, attributes are fully supported in our implementation through a trivial encoding, documented in Section 6. An instance of the XQuery data model can then be generated by the following grammar:

*Definition 2.1 (Data model).*

$$\mathbf{Tree} \ t ::= s_i \mid l_i[f] \qquad \mathbf{Forest} \ f ::= () \mid f, f \mid t$$

Essentially, an instance of the XQuery data model is an ordered sequence of labelled ordered *trees* (ranged over by  $t$ ). That is, an ordered *forest* (ranged over by  $f$ ), where each node has a unique *identifier* (ranged over by  $\mathbf{i}$ ) and where  $()$  denotes the empty forest. Tree nodes are labelled by *element tags* (ranged over by  $l$ ) while, without loss of generality, we consider only leaves that are text nodes (that is, strings, ranged over by  $s$ ) or empty trees (that is, elements that label the empty forest).

We define a complete partial order  $\preceq$  on forests (and thus on trees) by relating a forest with the forests obtained either by adding or by deleting subforests:

*Definition 2.2 (Projection ( $\preceq$ )).* Given two forests  $f$  and  $f'$  we say that  $f'$  is a *projection* of  $f$ , noted as  $f' \preceq f$ , if  $f'$  is obtained by replacing some subforests of  $f$  by the empty forest. In other terms  $\preceq$  is the smallest pre-congruence on forests that contains  $() \preceq f$  for all  $f$ .

We also define a notion of good formation, with respect to the data model given in Definition 2.1:

*Definition 2.3 (Good formation).* A forest is *well formed* if every identifier  $\mathbf{i}$  occurs in it at most once. Given a well-formed forest  $f$  and an identifier  $\mathbf{i}$  occurring in it, we denote by  $f@i$  the unique subtree  $t$  of  $f$  such that  $t = s_i$  or  $t = l_i[f']$ . The set of identifiers of a forest  $f$  is then defined as  $\mathbf{Ids}(f) = \{\mathbf{i} \mid \exists t. f@i = t\}$

Henceforth we will consider only well-formed forests and confound the notions of a node with that of the identifier of the node.

*Definition 2.4 (Root id).* Let  $t$  be a tree. If  $t = s_i$  or  $t = l_i[f]$ , we define  $\mathbf{RootId}(t) = \mathbf{i}$ .

## 2.2. Types and validation

In this work, we present our approach for an abstract model of types, namely *regular tree grammars*. It is well known that regular tree grammars encompass most of the features of well established schema specifications such as DTDs, XMLSchemas, RelaxNG definitions, XDuce and CDuce's regular expression types. This is for instance documented in [Murata et al. 2005], from where we borrow the definition of regular tree grammar:

*Definition 2.5 (Regular tree grammar).* A regular tree grammar is a pair  $(\mathcal{S}, E)$  where  $\mathcal{S}$  is a set of distinguished names (actually, non-terminal meta-variables) and  $E$  is a set of production rules of the form  $\{X_1 \rightarrow R_1, \dots, X_n \rightarrow R_n\}$  such that:

(1) each  $R_i$  is either the terminal *String*—denoting string content—, or the terminal *Any*—denoting any tree—, or  $l[r]$  where  $l$  ranges over valid element names and  $r$  is a regular expression on the non-terminal symbols  $X_1, \dots, X_n$ , that is:

$$\mathbf{RegExp} \quad r ::= \varepsilon \mid r r \mid r \mid r^* \mid X_i$$

(henceforth, we use  $r+$  for  $r r^*$  and  $r?$  for  $\varepsilon \mid r$ );

(2)  $\mathcal{S} \subseteq \{X_1, \dots, X_n\}$  is the set of start symbols;

(3) for any two production rules with the same left hand side  $X_i \rightarrow l[r]$  and  $X_i \rightarrow l'[r']$ , we have  $l \neq l'$ .

The intuition is that a regular tree grammar describes (i.e., it “types”) a set of trees of the data-model. Notice that the left-hand sides of the rules in  $E$  do not need to be pairwise distinct. Allowing two rules to have the same left-hand side allows us to freely take the union of two sets of rules and also simplifies some definitions. Furthermore, given a regular tree grammar, it is always possible to equivalently rewrite it so that condition 3 holds: if there are two rules  $X_i \rightarrow l[r]$  and  $X_i \rightarrow l'[r']$  then they can be merged into a single rule,  $X_i \rightarrow l[r|r']$ .

*Definition 2.6 (Names of a regular expression).* Given a regular expression  $r$  we denote by  $\mathbf{Names}(r)$  the set of non-terminals occurring in it, namely:

$$\begin{aligned}
Names(\varepsilon) &= \emptyset & Names(r^*) &= Names(r) \\
Names(r_1 r_2) &= Names(r_1) \cup Names(r_2) & Names(X) &= \{X\} \\
Names(r_1 | r_2) &= Names(r_1) \cup Names(r_2)
\end{aligned}$$

By extension, given a set  $E = \{X_0 \rightarrow R_0, \dots, X_n \rightarrow R_n\}$ , we define

$$Names(E) = \bigcup_{i \in \{0, \dots, n\}} Names(R_i)$$

*Definition 2.7 (Defined name).* Given a rule  $X \rightarrow R$ , we call  $X$  the defined name of the rule and we note  $\mathbf{Dn}(X \rightarrow R)$ . By extension, given a set  $E = \{X_0 \rightarrow R_0, \dots, X_n \rightarrow R_n\}$  we define  $\mathbf{Dn}(E) = \{X_0, \dots, X_n\}$ .

Note that in general,  $Names(E) \subseteq \mathbf{Dn}(e)$ . We also say that  $r$  is a regular expression over  $(\mathcal{S}, E)$ , if  $r$  is a regular expression over names in  $\mathbf{Dn}(E)$ . We will denote by  $\mathcal{L}(r)$  the language recognized by the regular expression  $r$ . We will use  $W, X, Y, Z$  to range over *names*. We use Greek letters to range over sets of rules. As  $(\mathcal{S}, E)$  represents a regular tree grammar we shall use  $\pi$  to stress that the set of rules is a *type projector* (cf. Definition 3.1) and  $\kappa$  and  $\tau$  to stress that the set is used as a context or as a type, respectively (cf. Section 5.1). Last, we use  $S$  to range over sets of (node) identifiers.

We illustrate the syntax of regular tree grammars with the following example:

EXAMPLE 2.8 (REGULAR TREE GRAMMAR FOR THE BIBLIOGRAPHY DTD). The bibliography DTD (taken from the XML Query use cases [Chamberlin et al. 2003]) can be written as a regular tree grammar  $(\{X\}, E)$ , with unique start symbol  $X$  and the following set  $E$  of rules:

$$\begin{aligned}
X &\rightarrow \text{bib}[\text{Book}^*] & Author &\rightarrow \text{author}[\text{String}] \\
Book &\rightarrow \text{book}[\text{Title}, (\text{Author} + | \text{Editor} +), \text{Publ}] & Editor &\rightarrow \text{editor}[\text{String}] \\
Title &\rightarrow \text{title}[\text{String}] & Publ &\rightarrow \text{publisher}[\text{String}]
\end{aligned}$$

This regular tree grammar “types” all XML documents (i.e., trees of the data model) that are rooted in a `bib` element, that contains a possibly empty list of `book` elements, each one containing a list starting with a `title` element containing a string, followed by a non-empty homogeneous list formed either by `author` elements or `editor` elements, and ended by a `publisher` element.

The concept of typing an XML document by a regular tree grammar is formalized by the notion of *validity* defined as follows:

*Definition 2.9 (Valid Trees).*

A tree  $t$  is *valid* with respect to a type  $(\mathcal{S}, E)$ , if there exists a mapping (interpretation)  $\mathcal{J}$  from  $\mathbf{Ids}(t)$  to  $E$  such that:

- (1)  $\mathbf{Dn}(\mathcal{J}(\mathbf{RootId}(t))) \in \mathcal{S}$
- (2) for each  $\mathbf{i}$  in  $\mathbf{Ids}(t)$ , if  $t@\mathbf{i} = s_i$  then  $\mathcal{J}(\mathbf{i}) = X \rightarrow \text{Any}$  or  $\mathcal{J}(\mathbf{i}) = X \rightarrow \text{String}$
- (3) for each  $\mathbf{i}$  in  $\mathbf{Ids}(t)$ , if  $t@\mathbf{i} = l_i[t_1, \dots, t_n]$ , then either:
  - $\mathcal{J}(\mathbf{i}) = X \rightarrow \text{Any}$  and  $\forall 1 \leq i \leq n, \exists X_i$  such that  $\mathcal{J}(\mathbf{RootId}(t_i)) = X_i \rightarrow \text{Any}$
  - or  $\mathcal{J}(\mathbf{i}) = X \rightarrow l[r]$  and  $\mathbf{Dn}(\mathcal{J}(\mathbf{RootId}(t_1))), \dots, \mathbf{Dn}(\mathcal{J}(\mathbf{RootId}(t_n))) \in \mathcal{L}(r)$ .

In this case, we say that  $t$  is  $\mathcal{J}$ -valid with respect to  $(\mathcal{S}, E)$  and denote it by  $t \in_{\mathcal{J}} (\mathcal{S}, E)$ .

For instance the following tree (in which we omit the node identifiers)

```

bib[
  book[
    title["Divina Commedia"], author["Dante"], publisher["Ludovico Dolce"]
  ]
]

```

is valid with respect to the type  $(\{X\}, E)$  defined in Example 2.8. There exist various techniques and algorithms to validate XML trees against regular tree grammars (for instance, by using tree automata: cf. Algorithm 4.4 in [Murata et al. 2005]). Note however that due to our use of regular tree grammars, the interpretation  $\mathcal{J}$  might not be unique and that a validating algorithm will generate—for a document  $t$  and a type  $(\mathcal{S}, E)$ —one possible interpretation such that  $t$  is  $\mathcal{J}$ -valid with respect to  $(\mathcal{S}, E)$ .

Given a tree  $t$  valid with respect to a type  $(\mathcal{S}, E)$ , we can use subsets of  $E$  to project that tree. Essentially, from the rules in  $E$  we compute another set of “simpler” rules which denotes only the nodes to be kept. At the beginning of the section we defined the projection of a forest as a forest obtained by replacing some subforests by the empty tree. Here we define an analogous concept for types, called *erasure* according to which a type is obtained from another by replacing some non-terminals by the empty regular expression.

**Definition 2.10 (Erasure of a regular expression).** Let  $r$  be a regular expression and  $N$  a set of names. We define the erasure of  $r$  with respect to  $N$  and we write  $r|_N$  the regular expression inductively defined as:

$$\begin{array}{ll} \varepsilon|_N = \varepsilon & (r^*)|_N = (r|_N)^* \\ (r_1 r_2)|_N = r_1|_N r_2|_N & X|_N = X \quad \text{if } X \in N \\ (r_1 | r_2)|_N = r_1|_N | r_2|_N & X|_N = \varepsilon \quad \text{if } X \notin N \end{array}$$

We generalize this notion to production rules of a grammar:

**Definition 2.11 (Erasure of a rule).** Let  $X \rightarrow R$  be a production rule, and  $N$  a set of names. We define the *erasure* of  $X \rightarrow R$  with respect to  $N$ , noted  $(X \rightarrow R)|_N$ , as:

$$\begin{array}{l} (X \rightarrow l[r])|_N = X \rightarrow l[r|_N] \\ (X \rightarrow \text{String})|_N = X \rightarrow \text{String} \\ (X \rightarrow \text{Any})|_N = X \rightarrow \text{Any} \end{array}$$

We recall that *String* and *Any* are special *terminals* denoting string and any content, respectively. We can finally define the erasure of a grammar:

**Definition 2.12 (Erasure of a tree grammar).** Let  $(\mathcal{S}, E)$  and  $(\mathcal{S}', E')$  be tree grammars. We say that  $(\mathcal{S}', E')$  is an *erasure* of  $(\mathcal{S}, E)$ , noted  $(\mathcal{S}', E') <: (\mathcal{S}, E)$ , if and only if all the following conditions hold

- (1)  $\mathcal{S}' \subseteq \mathcal{S}$ ;
- (2) if  $X \rightarrow \text{String} \in E'$ , then  $X \rightarrow \text{String} \in E$ ;
- (3) if  $X \rightarrow \text{Any} \in E'$ , then  $X \rightarrow \text{Any} \in E$ ;
- (4) for all rules  $X \rightarrow l[r'] \in E'$ , there exists a rule  $X \rightarrow l[r] \in E$  such that  $r' = r|_N$  for some  $N \subseteq \text{Names}(r)$ .

Note that in the definition above, the set  $N$  can be different for different rules.

We conclude this section by recalling some definitions taken from [Murata et al. 2005] that will be useful for establishing further results.

**Definition 2.13 (Competing non-terminals).** Let  $(\mathcal{S}, E)$  be a tree grammar. Let  $A, B \in \text{Names}(E)$  be two non-terminals such that  $A \neq B$ .  $A$  and  $B$  are *competing* if and only if there exist  $A \rightarrow l[r] \in E$  and  $B \rightarrow l'[r'] \in E$  such that  $l = l'$ .

The definitions that are actually interesting are those of *local* and *single-type tree* grammars, which can be defined in terms of competing non-terminals:

**Definition 2.14 (local tree grammar).** A regular tree grammar  $(\mathcal{S}, E)$  is a *local* tree grammar if and only if:

$$\text{— } |\mathcal{S}| \leq 1$$

- $E$  does not contain any competing non-terminals
- For all  $Y \in \mathbf{Names}(E)$  there is exactly one rule in  $E$  whose left-hand-side is  $Y$ .

*Definition 2.15 (single-type tree grammar).* A regular tree grammar  $(\mathcal{S}, E)$  is a *single-type tree grammar* if and only if:

- (1) For all  $X \rightarrow l[r] \in E$ , if  $A, B \in \mathbf{Names}(r)$  and  $A \neq B$ , then  $A$  and  $B$  are not competing
- (2) no pair of distinct non-terminals in  $\mathcal{S}$  is competing.

The interest of these two definitions is that—as shown in [Murata et al. 2005]—they characterize the structural constraints that can be expressed by the two most widespread schema formalisms, namely DTDs (which roughly correspond to local tree grammars) and XML Schemas (which are, essentially, single-type tree grammars).

### 3. TYPE PROJECTORS

In this section we shall first precisely define what type projectors are and then establish some useful closure results on type projectors.

#### 3.1. Definition

*Definition 3.1 (Type Projector).* Given a type  $(\mathcal{S}, E)$ , a (possibly empty) set of rules  $\pi \subseteq E$  is a *type projector* if and only if  $(\mathcal{S} \cap \mathbf{Names}(\pi), \pi)$  is a regular tree grammar erasure of  $(\mathcal{S}, E)$ .

A type projector is thus a set of rules obtained from the type  $(\mathcal{S}, E)$  by erasing some rules and some non-terminals in the remaining rules.

A type projector for a given type describes a particular pruning for XML documents of that type, that is, a *type driven projection*:

*Definition 3.2 (Type Driven Projections).* Let  $\pi$  be a type projector for  $(\mathcal{S}, E)$  and  $t$  a forest such that  $t \in_{\mathcal{J}}(\mathcal{S}, E)$ . The  $\pi$ -projection of  $t$ , noted as  $t \setminus_{\mathcal{J}} \pi$ , is defined as follows:

$$\begin{aligned}
 () \setminus_{\mathcal{J}} \pi &= () \\
 s_{\mathbf{i}} \setminus_{\mathcal{J}} \pi &= s_{\mathbf{i}} && \text{if } \mathcal{J}(\mathbf{i}) \rightarrow \text{String} \in \pi \text{ or } \mathcal{J}(\mathbf{i}) \rightarrow \text{Any} \in \pi \\
 s_{\mathbf{i}} \setminus_{\mathcal{J}} \pi &= () && \text{if } \mathcal{J}(\mathbf{i}) \rightarrow \text{String} \notin \pi \text{ and } \mathcal{J}(\mathbf{i}) \rightarrow \text{Any} \notin \pi \\
 l_{\mathbf{i}}[f] \setminus_{\mathcal{J}} \pi &= l_{\mathbf{i}}[f] && \text{if } \mathcal{J}(\mathbf{i}) \rightarrow \text{Any} \in \pi \\
 l_{\mathbf{i}}[f] \setminus_{\mathcal{J}} \pi &= l_{\mathbf{i}}[f \setminus_{\mathcal{J}} \pi] && \text{if } \mathcal{J}(\mathbf{i}) \rightarrow l[r] \in \pi \text{ and } \mathcal{J}(\mathbf{i}) \rightarrow \text{Any} \notin \pi \\
 l_{\mathbf{i}}[f] \setminus_{\mathcal{J}} \pi &= () && \text{if } \mathcal{J}(\mathbf{i}) \rightarrow l[r] \notin \pi \text{ and } \mathcal{J}(\mathbf{i}) \rightarrow \text{Any} \notin \pi \\
 (f, f') \setminus_{\mathcal{J}} \pi &= (f \setminus_{\mathcal{J}} \pi, f' \setminus_{\mathcal{J}} \pi)
 \end{aligned}$$

In words, pruning erases (by replacing it by an empty forest) every node that cannot be derived by a rule in  $\pi$ .

**LEMMA 3.3.** *Let  $\pi$  be a type projector for  $(\mathcal{S}, E)$ . Then for every tree  $t \in_{\mathcal{J}}(\mathcal{S}, E)$  it holds  $(t \setminus_{\mathcal{J}} \pi) \preceq t$ .*

As the knowledgeable reader might have already noticed, validation (as in Definition 2.9) and type-driven projection are quite similar. Given a tree  $t$  and a type  $(\mathcal{S}, E)$ , a validation algorithm builds an interpretation  $\mathcal{J}$  of  $t$  with respect to that type. More precisely, the algorithm associates to each node of  $t$  a non terminal of  $E$ . If it cannot find at least one, validation fails and the tree is not valid with respect to  $(\mathcal{S}, E)$ . A type-driven projecting algorithm works *exactly in the same way* but when a node cannot be associated with a name it is simply discarded together with the associated subtree. Projecting a document can be seen as an instance of validation. This observation allows us to determine the complexity of type-driven projection, given a particular type projector  $\pi$ . If  $\pi$  is a local tree grammar or a single-type tree grammar (that is, a DTD or an XML-Schema, see [Murata et al. 2005]) then projection can be performed in a streaming fashion. On the contrary, if  $\pi$  ends up being a general tree grammar, then projection might require in the worst case to keep the whole tree



in memory (see our remark at the end of Section 3.2, for how to use type projection in this particular setting).

### 3.2. Closure properties

The fact that *if* a type projector is a DTD or an XMLSchema, then type-driven projection can be done efficiently is already a good thing. However, we can show a stronger result: a type projector inherits the properties of the type it was deduced from. This is important since in practice if someone chooses to use DTDs or XML-Schemas to specify their documents, the projection process should not be more expensive than the validation process.

Indeed, a nice property of the *erasure* of a type is that it preserves both the local tree and single type property. In other words, the erasure of a DTD remains a DTD and the erasure of an XML-Schema remains an XML-Schema.

**PROPOSITION 3.4 (ERASURE PRESERVES LOCALITY).** *Let  $(\mathcal{S}, E)$  be a local tree grammar and  $(\mathcal{S}', E')$  a regular tree grammar. If  $(\mathcal{S}', E') <: (\mathcal{S}, E)$  then  $(\mathcal{S}', E')$  is a local tree grammar.*

**PROPOSITION 3.5 (ERASURE PRESERVES SINGLE-TYPEDNESS).** *Let  $(\mathcal{S}, E)$  be a single-type tree grammar and  $(\mathcal{S}', E')$  a regular tree grammar. If  $(\mathcal{S}', E') <: (\mathcal{S}, E)$  then  $(\mathcal{S}', E')$  is a single-type tree grammar.*

Last but not least, we show that if two projectors coming from the same type enjoy the local (resp. single-type) property, then their union is also local (resp. single-type). This property of type projectors is instrumental to our approach. Indeed, given a set of paths, we will compute a type projector for it by taking the union of all the type projectors of the individual paths. However, if taking the union of type projectors caused the loss of local or single-type properties, the interest of extending our approach to sets of paths (and thus to XQuery or to bunches of queries) would be quite limited. The key observation here is that, while in general local and single-type tree grammars are not closed under union, two type-projectors that *come from the same type* share a common structure and therefore are not completely independent from one another. In particular computing the union of two type projectors for the same type cannot introduce competing non-terminals. In terms of term-rewrite systems, we can say that the union of two type projectors does not introduce a critical pair (of non-terminals).

**PROPOSITION 3.6 (UNION CLOSURE OF LOCAL TYPE PROJECTORS).** *Let  $(\mathcal{S}, E)$  be a local tree grammar. Let  $(\mathcal{S}_1, E_1)$  and  $(\mathcal{S}_2, E_2)$  be two tree grammars such that  $(\mathcal{S}_1, E_1) <: (\mathcal{S}, E)$  and  $(\mathcal{S}_2, E_2) <: (\mathcal{S}, E)$ . Then  $(\mathcal{S}_1 \cup \mathcal{S}_2, E_1 \cup E_2)$  is a local tree grammar.*

**PROPOSITION 3.7 (UNION CLOSURE OF SINGLE-TYPE TYPE PROJECTORS).** *Let  $(\mathcal{S}, E)$  be a single-type tree grammar. Let  $(\mathcal{S}_1, E_1)$  and  $(\mathcal{S}_2, E_2)$  be two tree grammars such that  $(\mathcal{S}_1, E_1) <: (\mathcal{S}, E)$  and  $(\mathcal{S}_2, E_2) <: (\mathcal{S}, E)$ . Then  $(\mathcal{S}_1 \cup \mathcal{S}_2, E_1 \cup E_2)$  is a single-type tree grammar.*

To conclude this presentation of the formal properties of type-projectors we could note that a third category of deterministic regular tree grammars, namely *restrained-competition tree grammars* (see [Murata et al. 2005]), is not closed under erasure. However, pruning may still be computed in a streaming fashion. Indeed, it is known (see for instance [Martens et al. 2005]) that restrained-competition tree grammars allow one to type a document using a one pass pre-order traversal. Therefore, when meeting the opening tag of an element, the pruning process knows its type and can therefore use the type-projector to decide whether to keep it or not. All the other schema specifications that we are aware of (XDuce and CDuce regular expression types, TREX, Relax Core, ...) possess the full expressive power of regular tree languages which, as it is well-known, are closed under erasure and union (see for instance [Comon et al. 1997]). This means that type driven projection proposed here can be applied to these kinds of schemas, as well. However, projection remains as expensive as validation which, for these particular schemas, implies that the whole document might need to be loaded into memory to actually decide which subtrees must be pruned. Practical solutions to this problem are discussed in Section 8.2

#### 4. XPATH<sup>ℓ</sup>

In XPath, queries are expressed by defining a path of steps separated by “/”. For instance,

$$Q = /descendant::author/child::text[self::node = "Dante"]/parent::book/child::title$$

is the query that returns all titles of books whose author is "Dante". First, the navigational part instructs to descend to all text nodes whose parent is an author (by following the path `/descendant::author/child::text`), then the predicate selects those nodes that are the string "Dante" (with the test `self::node = "Dante"`), and finally the navigation ascends to the `book` element and descends to the `title`.

The inference rules we define in Section 5 do not work directly on queries such as  $Q$ . The rules are defined for a subset of XPath that we dub XPath<sup>ℓ</sup> and introduce in this section. XPath<sup>ℓ</sup> (for XPath *light*) includes forward and backward axes and a special kind of predicates. In order to statically analyse  $Q$  (or any other XPath query that is not in XPath<sup>ℓ</sup>), we will find an XPath<sup>ℓ</sup> query that approximates  $Q$  soundly with respect to the pruning inferred (Section 6), and use it to deduce the pruning for  $Q$ . Of course, these approximations, as well as those we introduce later on, will only be used to determine the pruning: the pruned document will be queried by the original query. Therefore we are going to proceed as follows. In this section we define XPath<sup>ℓ</sup>, which is roughly equivalent to the structural subset of positive core XPath, without absolute paths. Then in Section 5, we introduce our type and type-projector inference algorithms, which work on XPath<sup>ℓ</sup> queries. To complete the treatment of XPath we show in Section 6 how to compute a sound approximation of a query  $Q$  with respect to type projection. In other words, given a (full) XPath query  $Q$ , we will compute an XPath<sup>ℓ</sup> query  $Q'$  such that the type projector inferred from  $Q'$  preserves the semantics of  $Q$ .

Let us start with defining XPath<sup>ℓ</sup> paths and their semantics. From now on, “path” refers to an XPath<sup>ℓ</sup> query as defined hereafter unless otherwise specified (our XPath<sup>ℓ</sup> queries are sometimes called “twig queries” or “tree patterns” in the literature).

*Definition 4.1 (XPath<sup>ℓ</sup> path).* An XPath<sup>ℓ</sup> path is a term inductively generated by the following grammar:

$$\begin{aligned} Path &::= Step \mid Path/Path \mid Path \mid Path \\ Step &::= Axis::Test \mid Axis::Test[Cond] \\ Axis &::= self \mid child \mid descendant \mid parent \mid ancestor \\ Test &::= tag \mid node \mid text \\ Cond &::= Cond \text{ or } Cond \mid Cond \text{ and } Cond \mid Path \end{aligned}$$

where *tag* is a meta-variable ranging over element tags.

As customary, “and” takes precedence over “or” and the path delimiter “/” takes precedence over the top-level union “ $\mid$ ”. We will also use the (possibly indexed) meta-variables  $P$  and  $C$  to range over paths and conditions, respectively.

The formal semantics of paths is inductively defined on the productions of Definition 4.1. First, we formalise *Test* filtering as the set of nodes that satisfy a given test. Then *Axis* selection as the set of nodes reachable from some context nodes by following some *Axis*. Finally, we combine these notions to define the semantics of paths. The definitions comply with the semantics of XPath 1.0 (see [XPath 1999]).

*Definition 4.2 (Node test semantics).* Given a tree  $t$  and a set of nodes  $S \subseteq \mathbf{Ids}(t)$  we define:

$$\begin{aligned} S:::l &= S \cap \{\mathbf{i} \in \mathbf{Ids}(t) \mid t@\mathbf{i} = l_i[f]\} \\ S:::node &= S \\ S:::text &= S \cap \{\mathbf{i} \in \mathbf{Ids}(t) \mid \exists s, t@\mathbf{i} = s_i\} \end{aligned}$$

*Definition 4.3 (Axes selection).* Given a tree  $t$  and a set of nodes  $S \subseteq \mathbf{Ids}(t)$  (called context nodes), we define  $\llbracket Axis \rrbracket_t(S)$  as the set of nodes obtained by applying *Step* to each node in  $S$ :

$$\begin{aligned}
\llbracket \text{self} \rrbracket_r(S) &= S \\
\llbracket \text{child} \rrbracket_r(S) &= \bigcup_{i \in S} \{i' \mid (i, i') \in \text{Edg}(t)\} \\
\llbracket \text{parent} \rrbracket_r(S) &= \bigcup_{i \in S} \{i' \mid (i', i) \in \text{Edg}(t)\} \\
\llbracket \text{descendant} \rrbracket_r(S) &= \bigcup_{i \in S} \{i' \mid (i, i') \in \text{Edg}(t)^+\} \\
\llbracket \text{ancestor} \rrbracket_r(S) &= \bigcup_{i \in S} \{i' \mid (i', i) \in \text{Edg}(t)^+\}
\end{aligned}$$

where  $\text{Edg}(t)$  is the *edge* relation of  $t$ , that is

$$\text{Edg}(t) = \{(i, i') \mid t@i = l_i[f, t', f'] \wedge \text{RootId}(t') = i'\}$$

and  $\text{Edg}(t)^+$  denotes its transitive closure.

Since predicates may contain paths and conversely, path and predicate semantics are mutually defined.

*Definition 4.4 (XPath<sup>ℓ</sup> semantics).* Given  $t$ , a set  $S \subset \text{Ids}(t)$  and a path  $P$ , we define the evaluation of path  $P$  over the set of context nodes  $S$  as the function  $\llbracket P \rrbracket_r(S)$  defined as:

$$\begin{aligned}
\llbracket \text{Axis} :: \text{Test} \rrbracket_r(S) &= (\llbracket \text{Axis} \rrbracket_r(S)) :: {}_r \text{Test} \\
\llbracket \text{Axis} :: \text{Test}[C] \rrbracket_r(S) &= (\llbracket \text{Axis} \rrbracket_r(S)) :: {}_r \text{Test} \cap \{i \in S \mid \text{Check}_r[C](i)\} \\
\llbracket \text{Path}_1 / \text{Path}_2 \rrbracket_r(S) &= \llbracket \text{Path}_2 \rrbracket_r(\llbracket \text{Path}_1 \rrbracket_r(S)) \\
\llbracket \text{Path}_1 \downarrow \text{Path}_2 \rrbracket_r(S) &= \llbracket \text{Path}_2 \rrbracket_r(S) \cup \llbracket \text{Path}_1 \rrbracket_r(S)
\end{aligned}$$

where  $\text{Check}_r[\_](\_)$  is the Boolean function defined as:

$$\begin{aligned}
\text{Check}_r[\text{Path}](i) &= \llbracket \text{Path} \rrbracket_r(\{i\}) \neq \emptyset \\
\text{Check}_r[C_1 \text{ or } C_2](i) &= \text{Check}_r[C_1](i) \vee \text{Check}_r[C_2](i) \\
\text{Check}_r[C_1 \text{ and } C_2](i) &= \text{Check}_r[C_1](i) \wedge \text{Check}_r[C_2](i)
\end{aligned}$$

It is easy to see that the last definition is well founded since terms are inductively generated by the productions of the grammar in Definition 4.1.

Although the paths in XPath<sup>ℓ</sup> are quite simple, defining rules for their static analysis is challenging: the simultaneous presence in a single step of axes, tests, and predicates can cause a case explosion in the definition of the analysis. This is not a problem for a static analyzer, but it is a problem for a human reader. Fortunately, for the human reader, XPath<sup>ℓ</sup> paths can be further simplified and transformed into equivalent normal forms in which all non trivial axes, tests and predicates are distributed over different steps. The idea is then to normalize paths before passing them to the static analyzer so that the definition of the latter can result much simpler. The normal forms that will be analyzed by the static analysis of Section 5 are defined as follows

*Definition 4.5 (Single step normal form).* Let  $P$  be an XPath<sup>ℓ</sup> query. The *single step normal form* of  $P$ , noted  $\text{Snf}(P)$ , is defined as:

$$\begin{aligned}
\text{Snf}(\text{Axis} :: \text{node}) &= \text{Axis} :: \text{node} \\
\text{Snf}(\text{self} :: \text{Test}) &= \text{self} :: \text{Test} \\
\text{Snf}(\text{self} :: \text{node}[C]) &= \text{self} :: \text{node}[\text{Snf}_C(C)] \\
\text{Snf}(\text{Axis} :: \text{Test}) &= \text{Axis} :: \text{node} / \text{self} :: \text{Test} \quad (\text{if } \text{Axis} \neq \text{self} \wedge \text{Test} \neq \text{node}) \\
\text{Snf}(\text{Axis} :: \text{Test}[C]) &= \text{Axis} :: \text{node} / \text{self} :: \text{Test} / \text{self} :: \text{node}[\text{Snf}_C(C)] \\
&\hspace{15em} (\text{if } \text{Axis} \neq \text{self} \wedge \text{Test} \neq \text{node}) \\
\text{Snf}(P_1 / P_2) &= \text{Snf}(P_1) / \text{Snf}(P_2) \\
\text{Snf}_C(C_1 \text{ or } C_2) &= \text{Snf}_C(C_1) \text{ or } \text{Snf}_C(C_2) \\
\text{Snf}_C(C_1 \text{ and } C_2) &= \text{Snf}_C(C_1) \text{ and } \text{Snf}_C(C_2) \\
\text{Snf}_C(P) &= \text{Snf}(P)
\end{aligned}$$

It is clear from this definition that  $P$  and  $\text{Snf}(P)$  have the same semantics. Indeed, if we have a step

$$\text{Axis} :: \text{Test}[C]$$

then its single step normal form

$$\text{Axis}::\text{node}/\text{self}::\text{Test}/\text{self}::\text{node}[\text{Snf}_C(C)]$$

only makes the order of node selection more explicit<sup>1</sup>. For a given set of context nodes  $S$ , we first select all nodes that can be reached by the *Axis*. Then we keep only nodes that match the *Test*. Finally, we normalize the predicate  $C$  by putting every path it contains in single-step normal form (through the use of the auxiliary function  $\text{Snf}_C$ ).

## 5. STATIC ANALYSIS

In this section we define deduction rules to statically infer from an XPath<sup>l</sup> path  $P$  and a type  $(\mathcal{S}, E)$  a type-projector for any input document valid with respect to  $(\mathcal{S}, E)$ . We show that the analysis is sound, and that it enjoys completeness for a large class of queries when  $E$  is a \*-guarded and non-recursive local tree grammar (see Definition 5.6 later on). Soundness means that executing the query on the original document and on the document pruned by the inferred projector yields the same result. Completeness means that the analysis infers the best correct projector, that is, that if we take a type projector smaller (i.e., more selective) than the inferred one, then there exists a document validating  $(\mathcal{S}, E)$  for which the result of the two executions is not the same. When the conditions on schemas or on queries are relaxed, then the analysis is still sound but it may be not complete. Nevertheless, as we will formally illustrate, it is still very precise.

In order to define our static type-projector inference algorithm we proceed in two steps.

(1) Given a path  $P$  and a regular expression grammar  $(\mathcal{S}, E)$  the rough idea is to use a type system to associate  $P$  with the set of all trees that may appear in the result of applying  $P$  to a document validating  $(\mathcal{S}, E)$ . In order to achieve a great precision, we then “type”  $P$  by the set of all rules of  $E$  that *validate* any tree in the result.<sup>2</sup> This is done in Section 5.1.

(2) Next, we use the type system defined in the previous point to define inference of type projectors. In particular we use the cases in which the previous type system returns an empty set of rules to determine the points in which pruning must be performed. This is done in Section 5.2.

### 5.1. Type inference

Given a path  $Path$  and a schema  $(\mathcal{S}, E)$  we want to find a subset of rules in  $E$  that can type all the nodes that can occur in the result of  $Path$  when applied to the *RootId* of a tree validating  $(\mathcal{S}, E)$ . Formally, we want to infer a set  $\tau \subseteq E$  such that

$$\forall t \in \mathcal{T}(\mathcal{S}, E), \mathcal{J}(\llbracket Path \rrbracket_t, (\text{RootId}(t))) \subseteq \tau \quad (1)$$

The equation above states the soundness of the analysis. In words it says that if we take any tree  $t$  valid for  $(\mathcal{S}, E)$  and we apply the path  $Path$  to its *RootId*, then the type  $\tau$  inferred in the type systems defines every rule interpreting a node in the result. As usual, soundness alone is not interesting since there always are sets that trivially satisfy it (notably, the set of all rules in  $E$ ). What we aim at is an analysis that is as selective as possible, that is, an analysis that is precise enough to guarantee, on a large class of types and for a large class of queries, that whenever the path semantics is empty over all possible instances of the input type, then the inferred type  $\tau$  is empty, as well:

$$\forall t \in \mathcal{T}(\mathcal{S}, E), \mathcal{J}(\llbracket Path \rrbracket_t, (\text{RootId}(t))) = \emptyset \implies \tau = \emptyset \quad (2)$$

(the converse is a consequence of (1)). In other terms we want that if there does not exist any instance of the type that matches the path, then the path is typed by the empty set.

The precision described by (2) will then be used during the inference of type-projectors to discard elements that are useless in the evaluation of  $Path$ , that is, all the sub-trees of the original document that cannot be matched by  $Path$ .

<sup>1</sup>As an aside, note that this kind of equivalence does not hold for full XPath because of the `position()` function. Indeed, `descendant::a[position()=1]` and `descendant::node/self::a[position()=1]` do not return, in general, the same result. The former returns the first “a”-node in pre-order while the latter returns all the “a”-nodes of the document.

<sup>2</sup>This yields a finer-grained analysis since different rules may generate the same tree but in different contexts.

We start by inferring types for single-step paths without predicates.

*Definition 5.1 (Unconditional Single Step Typing).* The type of an unconditional single-step query  $Axis :: Test$  for the schema  $(\mathcal{S}, E)$  is given by:

$$\mathbf{T}_E(\mathbf{A}_E(\mathcal{S}, Axis), Test)$$

where axes are typed as:

$$\begin{aligned} \mathbf{A}_E(\tau, \text{self}) &= \tau \\ \mathbf{A}_E(\tau, \text{child}) &= \bigcup_{Y \rightarrow R \in \tau} \{Z \rightarrow R' \in E \mid Z \in \text{Names}(R)\} \\ \mathbf{A}_E(\tau, \text{parent}) &= \bigcup_{Y \rightarrow R \in \tau} \{Z \rightarrow R' \in E \mid Y \in \text{Names}(R')\} \\ \mathbf{A}_E(\tau, \text{descendant}) &= \tau' \cup \mathbf{A}_E(\tau', \text{descendant}) \quad \text{where } \tau' = \mathbf{A}_E(\tau, \text{child}) \\ \mathbf{A}_E(\tau, \text{ancestor}) &= \tau' \cup \mathbf{A}_E(\tau', \text{ancestor}) \quad \text{where } \tau' = \mathbf{A}_E(\tau, \text{parent}) \end{aligned}$$

and tests are typed as:

$$\begin{aligned} \mathbf{T}_E(\tau, \text{node}) &= \tau \\ \mathbf{T}_E(\tau, a) &= \{Y \rightarrow R \in \tau \mid R = a[R'] \text{ or } R = \text{Any}\} \\ \mathbf{T}_E(\tau, \text{text}) &= \{Y \rightarrow R \in \tau \mid R = \text{String} \text{ or } R = \text{Any}\} \end{aligned}$$

This definition introduces two typing operators. Firstly,  $\mathbf{A}_E(\tau, Axis)$  returns all the rules that can be reached from names in  $\tau$  following  $Axis$ . If  $Axis$  is `self`, `child` or `descendant`, our definition coincide with the static semantics of XQuery and XPath, as defined by Draper *et al.* in [Draper et al. 2007]. However, Draper *et al.*'s static semantics is much less precise than ours when dealing with backward axes. Translated in our formalism, the type of `parent` and `ancestor` for any  $\tau$  would be  $\{X \rightarrow \text{Any}\}$  for some name  $X$ <sup>3</sup>. Secondly,  $\mathbf{T}_E(\tau, test)$  restricts the rules in  $\tau$  to only rules which type elements compatible with  $test$ .

The soundness of this definition, that is, the property stated by Formula (1) is given by the following lemma.

**LEMMA 5.2.** *Let  $t$  be a tree  $\mathcal{J}$ -valid with respect to the schema  $(\mathcal{S}, E)$ . For every  $S \subseteq \mathbf{Ids}(t)$  and type  $\tau$ , if  $\mathcal{J}(S) \subseteq \tau$ , then*

$$\mathcal{J}(\llbracket Axis \rrbracket_t(S)) \subseteq \mathbf{A}_E(\tau, Axis) \quad \text{and} \quad \mathcal{J}(S ::_t Test) \subseteq \mathbf{T}_E(\tau, Test)$$

It is easy to check that the property stated by Formula (1) is a direct consequence of Definition 4.4 and the composition of the two properties of the lemma above.

The presence of upward axes makes the typing of composed paths much more difficult. To ensure precision, that is the property stated by Formula (2), we have to be careful in dealing with types in which an element may occur in the content of different elements. The crux of the matter is that a given name  $X$  might be reachable by two distinct paths from a start symbol (even when very simple types such as DTDs are considered). When navigating upward from  $X$  with single step typing alone, the two possible sets of ancestor may end-up in the result type. This behaviour does not take into account the fact that on a tree (that is, an instance of the type) there is only one path from the root to a node of type  $X$ . This can be illustrated by the following grammar rooted at  $A$ :

$$A \rightarrow a[B C], B \rightarrow b[D], C \rightarrow c[D], D \rightarrow d[]$$

and observe that  $D$  is a child of both  $B$  and  $C$ . Now consider the path

$$\text{self} :: a / \text{child} :: b / \text{child} :: d / \text{parent} :: \text{node}$$

applied to documents of the above type, then the precise type that this path should have is  $\{B \rightarrow b[D]\}$ . However if we naively iterate Definition 5.1, we obtain at the first step  $\{A \rightarrow a[B C]\}$ , onto which we apply `child::b`, which yields  $\{B \rightarrow b[D]\}$  to which we apply `child::d`, which yields

<sup>3</sup>More precisely `parent::test` and `ancestor::test` return the union type `element() | document()` independently of  $test$ ; where `element()` is the type of any element node and `document()` is the type of the document node which we don't consider in our data-model.

$\{D \rightarrow d[\ ]\}$ , to which we finally apply `parent::` which returns  $\{B \rightarrow b[D], C \rightarrow c[D]\}$ , which is sound but imprecise. This is due to the fact that the single step typing blindly selects all rules associated with a name which can generate  $Z$ , here all the rules associated with  $Y$ .

To solve this problem we introduce particular sets of rules, called *contexts*, to be updated at each step and containing rules already encountered in previous steps. We then use them to refine type inference for upward axes. In the previous example, when typing the first two steps we build a *context*

$$\{A \rightarrow a[B\ C], B \rightarrow b[D]\}$$

indicating that for the moment the two rules are the only ones visited by the traversal. Then, we use Definition 5.1 to type `parent::node` thus obtaining  $\{B \rightarrow b[D], C \rightarrow c[D]\}$ , as before, but this time we intersect it with the context, thus obtaining the precise answer  $\{B \rightarrow b[D]\}$ . Intuitively a context has the grammar productions that type ancestors of nodes that are (currently) selected by the query.

We now formalize this idea:

*Definition 5.3 (Type inference).* Let  $(\mathcal{S}, E)$  be a type and  $P$  an XPath<sup>ℓ</sup> query in single step normal form. Let  $\tau$  and  $\kappa$  be two sets of rules of  $E$ . If the deduction system in Figure 1 deduces for a path  $P$  the judgment,

$$(\tau, \kappa) \vdash_E P : (\tau', \kappa')$$

then we say that  $P$  has type  $(Dn(\tau'), \tau')$ .

The idea underlying the judgments of the definition is that if the system proves

$$(\tau, \kappa) \vdash_E P : (\tau', \kappa')$$

then from an input set of rules  $\tau$  and an input context  $\kappa$  the application of  $P$  returns an output set of rules  $\tau'$  and an updated context  $\kappa'$ . In other terms  $\tau$  is (the production part of) a type that approximates the current nodes,  $\kappa$  is the context that was visited to type them,  $\tau'$  is (the production part of) a type that approximates the set of nodes reachable from the current ones by following  $P$ , and  $\kappa'$  is the additional context visited to reach them. In Figure 1 environments—that is pairs of sets of rules—are ranged over by  $\Sigma$  for concision.  $\Sigma$  being a pair  $(\tau, \kappa)$ , we use  $\Sigma_{\text{typ}}$  to denote its first projection (i.e., the “type” component  $\tau$ ) and  $\Sigma_{\text{ctx}}$  to denote the second projection (i.e., the “context” component  $\kappa$ ).

The first two rules implement our main idea: when we follow an axis *Axis*, we compute the type by  $\mathbf{A}_E(\Sigma_{\text{typ}}, \text{Axis})$ ; if the axis is a downward one, then we add this type to the current context, otherwise if the axis is an upward one, then we intersect it with the current context (both for the type part and for the context part). The rule for (**test**) discards from the current set of rules those that do not satisfy the test: the type is computed by  $\mathbf{T}_E(\Sigma_{\text{typ}}, \text{Test})$ , while the context is obtained by removing all the rules that were in there just because they generated one of the discarded nodes; to do so it generates (the type of) all ancestors of the nodes satisfying the test, and intersects them with the current context. The fourth rule, (**predicate**), is the most difficult one. It uses the auxiliary Boolean judgement  $\_ \vdash^{\text{Cond}} \_ : \_$ . This judgment simulates the evaluation of XPath conditions at the level of types: if a path yields a non empty-type (rule (**path**)) then its Boolean value is `true`. Rule (**boolean**) combines the truth value of the sub-conditions according to the operator at issue. This auxiliary judgement is then used in the (**predicate**) rule to discard from  $\Sigma_{\text{typ}}$  all rules for which the predicate never holds. The context is then computed as in the deduction rule (**test**), by discarding from it all rules that generated only rules discarded from  $\Sigma_{\text{typ}}$ . The deduction rule (**seq**) chains the result of one step to the following one. Lastly, the rule (**union**) handles the top-level union operator “[ ]”.

Let us illustrate how the algorithm works on an example. Consider the grammar

$$\{A \rightarrow a[B|C|E], B \rightarrow b[D], C \rightarrow c[\ ], D \rightarrow d[\ ], E \rightarrow b[\ ]\}$$

rooted in  $A$ , and the path

## Single Step

$$\begin{array}{l}
\text{(down-axis)} \frac{Axis \in \{\text{self}, \text{child}, \text{descendant}\}}{\Sigma \vdash_E Axis :: \text{node} : (\mathbf{A}_E(\Sigma_{\text{typ}}, Axis), \Sigma_{\text{ctx}} \cup \mathbf{A}_E(\Sigma_{\text{typ}}, Axis))} \\
\text{(up-axis)} \frac{Axis \in \{\text{parent}, \text{ancestor}\}}{\Sigma \vdash_E Axis :: \text{node} : (\mathbf{A}_E(\Sigma_{\text{typ}}, Axis)) \cap \Sigma_{\text{ctx}}, \mathbf{A}_E(\Sigma_{\text{ctx}}, Axis) \cap \Sigma_{\text{ctx}}} \\
\text{(test)} \frac{Test \neq \text{node}}{\Sigma \vdash_E \text{self} :: Test : (\tau, (\Sigma_{\text{ctx}} \cap \mathbf{A}_E(\tau, \text{ancestor})) \cup \tau)} (*) \\
\hspace{15em} (*) \text{ where } \tau = T_E(\Sigma_{\text{typ}}, Test) \\
\text{(predicate)} \frac{(\forall X_i \rightarrow R_i \in \Sigma_{\text{typ}}) \quad (\{X_i \rightarrow R_i\}, \Sigma_{\text{ctx}}) \vdash_E^{Cond} C : b_i}{\Sigma \vdash_E \text{self} :: \text{node}[C] : (\tau, (\Sigma_{\text{ctx}} \cap \mathbf{A}_E(\tau, \text{ancestor})) \cup \tau)} (***) \\
\hspace{15em} (***) \text{ where } \tau = \{X_i \rightarrow R_i \mid b_i = \text{true}\}
\end{array}$$

## Predicates

$$\text{(boolean)} \frac{\Sigma \vdash_E^{Cond} C_i : b_i \quad op \in \{\text{or}, \text{and}\}}{\Sigma \vdash_E^{Cond} C_1 \quad op \quad C_2 : b_1 \quad op \quad b_2} \quad \text{(path)} \frac{\Sigma \vdash_E P : \Sigma'}{\Sigma \vdash_E^{Cond} P : \Sigma'_{\text{typ}} \neq \emptyset}$$

## Composed paths

$$\text{(seq)} \frac{\Sigma \vdash_E Step : \Sigma'' \quad \Sigma'' \vdash_E Path : \Sigma'}{\Sigma \vdash_E Step/Path : \Sigma'} \quad \text{(union)} \frac{\Sigma \vdash_E Path_1 : (\tau_1, \kappa_1) \quad \Sigma \vdash_E Path_2 : (\tau_2, \kappa_2)}{\Sigma \vdash_E Path_1 \mid Path_2 : (\tau_1 \cup \tau_2, \kappa_1 \cup \kappa_2)}$$

Fig. 1. Inference rules for XPath<sup>l</sup> queries
$$\text{child}::\text{node}/\text{self}::b/\text{self}::\text{node}[\text{child}::\text{node}/\text{self}::d]$$

Notice that the path above is nothing but the single step normal form of

$$\text{child}::b[\text{child}::d]$$

We start from an initial environment

$$\Sigma = (\{A \rightarrow a[B|C|E]\}, \{A \rightarrow a[B|C|E]\})$$

in which both the context and the type component contain all the rules whose left hand side is a root of the grammar (in this case we have just one rule). The first step is typed with the **(down-axis)** rule, giving the result  $\Sigma^1$  where

$$\Sigma^1_{\text{typ}} = \{B \rightarrow b[D], C \rightarrow c[], E \rightarrow b[]\}$$

and

$$\Sigma^1_{\text{ctx}} = \{A \rightarrow a[B|C|E], B \rightarrow b[D], C \rightarrow c[], E \rightarrow b[]\}$$

The second step is typed by applying the rule **(test)**, which returns  $\Sigma^2$ :

$$\Sigma^2_{\text{typ}} = \{B \rightarrow b[D], E \rightarrow b[]\}$$

and more interestingly, the context

$$\Sigma^2_{\text{ctx}} = \{A \rightarrow a[B|C|E], B \rightarrow b[D], E \rightarrow b[]\}$$

Indeed, the intersection of  $\Sigma^1_{\text{ctx}}$  with the name generated by the ancestors of  $B$ , namely  $A$  yields exactly  $\{A \rightarrow a[B|C|E]\}$  to which we add the result of the current step:

$$\{B \rightarrow b[D], E \rightarrow b[\ ]\}$$

As we said, this intersection ensures that we only keep in the context rules from which we can derive the current type. In this example, the rules for  $C$  which was introduced by the wildcard step `child::node` is removed by the typing of the more restrictive step `self::b`. The third step is typed by the **(predicate)** rule. Intuitively, this rule types independently the path `child::node/self::d` and keeps in the result only the input rules for which the path yields a non-empty result which, in this case, is the rule for  $B$ :

$$\Sigma_{\text{typ}}^3 = \{B \rightarrow b[D]\}$$

As before, the context is purged from rules that do not generate the current type:

$$\Sigma_{\text{ctx}}^3 = \{A \rightarrow a[B|C|E], B \rightarrow b[D]\}$$

Before stating the main theorems of type inference, namely soundness and completeness, let us first show that the inference rules of Figure 1 form indeed an algorithm.

**LEMMA 5.4 (TERMINATION OF TYPE INFERENCE).** *Let  $(\mathcal{S}, E)$  be a type,  $P$  a path, and  $\Sigma$  and  $\Sigma'$  two environments. If there is a derivation for the judgment  $\Sigma \vdash_E P : \Sigma'$ , then this derivation is unique and finite. Furthermore, the algorithm runs in time  $O(|E|^d \times |P|)$ , where  $d$  is the maximum number of nested predicates.*

The complexity result of type inference is consistent with those of related problems. For instance [Benedikt et al. 2008] studies the problem of the satisfiability of XPath queries over a DTD for various fragments of XPath. For instance, satisfiability of forward XPath (`child`, `descendant`) without predicates is polynomial, which corresponds for us to the case where  $d = 1$  (no nested predicates). For even slightly richer fragments, with either `child` and `parent` or `child`, `descendant` and `predicates`, satisfiability is NP-Hard (see again [Benedikt et al. 2008]). Since *exact* type inference would imply satisfiability (a query is satisfiable if and only if its output type is not empty), one may rightfully wonder why for instance our algorithm runs in polynomial time for the fragment with `child` and `parent` (without predicates). The reason is that our algorithm is not complete in the general case (thus not *exact*) and that it might return a non-empty output type for an unsatisfiable query. We argue that our algorithm seems to be a good compromise between complexity and precision. This claim is further supported by our experimental results (Section 9).

We can now proceed to prove the soundness of the type system.

**THEOREM 5.5 (SOUNDNESS OF TYPE INFERENCE).** *Let  $(\mathcal{S}, E)$  be a type and  $P$  a path. Let  $E_0 = \{X \rightarrow R \mid X \rightarrow R \in E, X \in \mathcal{S}\}$ . If  $(E_0, E_0) \vdash_E P : (\tau, \kappa)$  then:*

$$\tau \supseteq \bigcup_{t \in \mathcal{J}(\mathcal{S}, E)} \mathfrak{J}(\llbracket P \rrbracket_t(\mathbf{RootId}(t)))$$

The type system is sound. It is also complete for a particular class of schemas, namely local tree grammars that are *\*-guarded*, *non-recursive*, and *parent-unambiguous*. Formally, we have the following definition:

*Definition 5.6.* Let  $(\mathcal{S}, E)$  be a local tree grammar.

- (1)  $E$  is *\*-guarded* if for each  $Y \rightarrow l[r]$  in  $E$ , the regular expression is a product  $r = r_1 \cdots r_n$  and whenever  $r_i$  contains a union, then  $r_i = (r')^*$ ;
- (2)  $E$  is *non-recursive* if for all  $Y \rightarrow R \in E$ ,  $Y \notin \mathbf{Dn}(\mathbf{A}_E(\{Y \rightarrow R\}, \text{descendant}))$ ;
- (3)  $E$  is *parent-unambiguous* if for all  $Y \rightarrow R \in E$ :

$$\mathbf{A}_E(\{Y \rightarrow R\}, \text{parent}) \cap \mathbf{A}_E(\mathbf{A}_E(\{Y \rightarrow R\}, \text{parent}), \text{ancestor}) = \emptyset$$

Non-recursive and *\*-guardedness* are properties enjoyed by a large number of commonly used DTDs. As an example, the reader can consider the DTDs of the XML Query Use Cases [Chamberlin et al. 2003]: among the ten DTDs defined in the Use Cases, seven are both non-recursive and



\*-guarded, one is only \*-guarded, one is only non-recursive, and just one does not satisfy either property. Beyond the XML Query use cases, it seems that \*-guardedness is commonly enjoyed by DTDs (e.g., Table 3 in [Choi 2002] shows that in most cases where a union is used, it is \*-guarded). Concerning the parent-unambiguous property, although DTDs satisfying this property are less frequent (five on the ten DTDs in [Chamberlin et al. 2003]), its absence is in practice not very problematic since, as we will see, only the presence of the `parent` axis may hinder completeness in that case.

Before stating the completeness of type inference, we illustrate on simple examples what happens when the conditions do not hold. For \*-guardedness, consider the grammar

$$X \rightarrow a[B|C], B \rightarrow b[], C \rightarrow c[]$$

rooted in  $X$ , together with the path:

$$\text{child::node/self::b/parent::node/child::node}$$

For the first two steps, our algorithm would determine the exact type and context:

$$\Sigma^2 = (\{B \rightarrow b[]\}, \{X \rightarrow a[B|C], B \rightarrow b[]\})$$

For the parent step, the type and context are:

$$\Sigma^3 = (\{X \rightarrow a[B|C]\}, \{X \rightarrow a[B|C]\})$$

which are also exact. However, the last step induces the final type:

$$\Sigma_{\text{typ}}^4 = \{B \rightarrow b[], C \rightarrow c[]\}$$

and the context:

$$\Sigma_{\text{ctx}}^4 = \{X \rightarrow a[B|C], B \rightarrow b[], C \rightarrow c[]\}$$

This is not exact because a document matching the first part, `child::node/self::b` does not have any “ $c$ ” tag and therefore the rule  $C$  in the output type is superfluous: this query will never return a node with type  $C$  for a document of the considered type. Note that the condition that unions in regular expressions must be guarded must also hold for rules, namely that there must not be two rules  $Y \rightarrow l[r_1]$  and  $Y \rightarrow l'[r_2]$  in the input type. Indeed these two rules behave like an un-guarded union and therefore jeopardize completeness. Local tree grammars forbid such rules and are thus an essential condition of the input type for completeness to hold.

The recursiveness of the schema also interacts with the `parent` axis in a way that prevents completeness of type inference. Consider the grammar:

$$\{A \rightarrow a[B], B \rightarrow b[B?]\}$$

and the path expression:

$$\text{child::node/self::b/child::node/self::b/parent::node}$$

Our type inference algorithm deduces on the second `self::b` step that the output type is be  $\{B \rightarrow b[B?]\}$ . However, the last step, `parent::node` is typed with a type

$$\{A \rightarrow a[B], B \rightarrow b[B?]\}$$

this is because in the grammar,  $A$  is a name reachable from  $B$  with a parent axis. However, consider any document valid with respect to this grammar. Either it has only one  $b$  element, in which case the result is empty, since we try to match two levels of  $b$ 's with the query. Or it has at least two  $b$ 's and then the output is always a  $b$  node (the topmost one). Therefore, an  $a$  node is never part of the result, while the type  $A$  is returned by our algorithm. Lastly, with the following parent-ambiguous grammar

$$\{A \rightarrow a[B | C], B \rightarrow b[], C \rightarrow c[B]\}$$

the algorithm fails to type *exactly* (but the output type is still sound) the query:

$$\text{child::node/self::c/child::node/self::b/parent::node}$$

By a similar reasoning, we can see that the algorithm returns the rules

$$\{A \rightarrow a[B|C], C \rightarrow c[B]\}$$

while only nodes with tag  $c$  can be returned by this query.

Intuitively, the reason why completeness does not hold in the three previous examples is that there are chains of derivations in the grammar that may not reflect actual paths in a document. For instance in the last example, in a document “a[ b[ ] ]”, the derivation “ $ACB$ ” has no interpretation (since there are no  $c$ -nodes). In this case, there exists a *valid* document which does not contain all the paths described by the possible derivation in its type. Therefore, the type inference algorithm will use rules which are actually not part of the interpretation of some documents of the type at issue. Fortunately, if a local tree grammar is  $*$ -guarded, non-recursive, and parent-unambiguous, then there always exists a document in which *all* the sequences of derivations in the grammar are instantiated by some path. We call such a document a *witness* of the grammar. We prove the existence of such a witness before stating the completeness theorem.

**LEMMA 5.7 (WITNESS OF A GRAMMAR).** *Let  $(\mathcal{S}, E)$  be a non-recursive,  $*$ -guarded, parent-unambiguous local tree grammar. There exists a document  $t$ ,  $\mathcal{J}$ -valid with respect to  $(\mathcal{S}, E)$  such that:*

$$\forall X \rightarrow R \in E, \exists \mathbf{i} \in \mathbf{Ids}(t) \text{ such that } \mathcal{J}(\mathbf{i}) = X \rightarrow R$$

*we call such a document a witness of the schema  $(\mathcal{S}, E)$ .*

We are now equipped to state (and prove) the completeness theorem:

**THEOREM 5.8 (COMPLETENESS OF TYPE INFERENCE).** *Let  $(\mathcal{S}, E)$  be a  $*$ -guarded non-recursive and parent unambiguous local tree grammar, and  $P$  a path. Let*

$$E_0 = \{X \rightarrow R \mid X \rightarrow R \in E, X \in \mathcal{S}\}.$$

*If  $(E_0, E_0) \vdash_E P : (\tau, \kappa)$  then:*

$$\tau \subseteq \bigcup_{t \in \mathcal{J}E} (\llbracket P \rrbracket_t(\mathbf{RootId}(t)))$$

One of the main reasons why completeness does not hold in general is because the intersections performed by the type rule for `parent` are not powerful enough to guarantee precision for a recursive or parent-ambiguous grammar. In a nutshell, this happens because in the presence of a parent-ambiguous grammar the type analysis may produce contexts containing false parent types (with respect to the current type  $\tau$ ). This suggests that in order to be more precise, instead of sets of rules, contexts should rather be sets of *chains* of names, computed and opportunely managed by the type analysis. However, we feel that (i) managing sets of chains instead of simple sets of rules dramatically complicates the treatment, due to the interaction of recursive axes like `descendant` and recursive grammars, (ii) the problem may arise only for queries that use parent axis and the concomitance of parent-ambiguity may make the event rare in practice, (iii) the loss of precision seems to be negligible in most cases, (iv) even though it would be possible to obtain more precise results for a larger class of grammars, it is well known that exact type-inference for XPath routinely escapes regular tree languages and therefore all existing formalisms to type XML: at some point, an approximation in the type inference process is necessary to remain in the realm of regular types. Therefore we considered that such a small gain did not justify the dramatic increase in complexity needed to relax the condition on the type for completeness to hold. Indeed, completeness is just some icing on the cake since, while it helps to gauge the precision of the approach, its absence does not hinder its application.

Of course, the completeness theorem is only stated for XPath<sup>ℓ</sup> queries and does not account for full XPath queries. Yet it illustrates how precise our type system is in the best case. We will show on various examples that on less favorable cases for schemas or for XPath queries which need to be approximated, type inference still remains very precise.

## 5.2. Type-Projection inference

In this section we use the type inference of the previous section to infer type-projectors. Once more, naive solutions do not work. For instance, for simple paths  $Step_1/\dots/Step_n$ , we may consider as type projector with respect to  $(\mathcal{S}, E)$  the set

$$\bigcup_{i \in \{1, \dots, n\}} \tau_i \cup \{X_i \rightarrow R_i \mid X_i \in \mathcal{S}\}$$

where for  $i \in \{1, \dots, n\}$ :

$$\Sigma \vdash_E Step_1/\dots/Step_i : (\tau_i, -)$$

with  $\Sigma = \{X_i \rightarrow R_i \mid X_i \in \mathcal{S}\}, \{X_i \rightarrow R_i \mid X_i \in \mathcal{S}\}$  (we use “-” as a placeholder for uninteresting parameters). This definition is sound but not precise at all, as can be seen by considering the query  $descendant :: node/Path$ : the use of the above union yields a set containing  $\tau_1$  defined as

$$\Sigma \vdash_E descendant :: node : (\tau_1, -)$$

that is, all descendants of the root start symbols in  $\mathcal{S}$  (no pruning is performed). Instead, we would like to discard, at least, all rules that are descendants of  $\mathcal{S}$  but that are not ancestors of a node matching  $Path$ . These are the rules

$Y \rightarrow R \in T_E(\mathbf{A}_E(\mathcal{S}, descendant), node)$  such that

$$(\{Y \rightarrow R\}, \kappa) \vdash_E descendant :: node/Path : (\emptyset, -)$$

for some appropriate context  $\kappa$ . A similar reasoning applies to ancestor.

As for the type inference, we define type-projector inference by a judgment and associated inference rules:

*Definition 5.9 (Type-projector inference).* Let  $(\mathcal{S}, E)$  be a type and  $P$  an XPath<sup>l</sup> query in simple step normal form, and  $\tau$  and  $\kappa$  be subsets of  $E$ . If the deduction system in Figure 2 proves the judgment

$$(\tau, \kappa) \Vdash_E P : \pi$$

then the *type-projector* induced by  $\pi$  is the grammar:

$$(\mathcal{S} \cap \mathbf{Dn}(\pi), \{(X \rightarrow R) \mid \mathbf{Dn}(\pi) \mid X \rightarrow R \in \pi\})$$

Obtaining a type projector from a set of rules returned by the judgment is straightforward. In essence, the derivation collects in  $\pi$  the rules of  $E$  that are sufficient to answer the query. Since in general not all rules in  $E$  are kept, then the rule in  $\pi$  may use names that are not defined in  $\pi$ . Therefore, the erasure operation (defined in Definition 2.10) simply removes references to names not defined by any rule in  $\pi$  (the definition of  $R|_{\mathcal{S}}$  is straightforward: it is  $R$  where every occurrence of a name in  $\mathcal{S}$  is replaced by  $\varepsilon$ ).

The rules in Figure 2 reflect the intuition we gave earlier. At each step, we execute the type inference algorithm on the current set of rules and accumulate only those for which the resulting type is not empty. Informally, each rule preserves the following properties:

**well-formedness:** if a rule  $Y \rightarrow R$  is added to the type projector  $\pi$  then there must be a rule  $X \rightarrow R' \in \pi$  such that  $Y \in \mathbf{Names}(R')$ .

**precision:** given a path  $P$  and a rule  $Y \rightarrow R$ . If  $(\{Y \rightarrow R\}, -) \vdash_E P : (\emptyset, -)$  then  $Y \rightarrow R$  must not be added to the projector.

Let us explain how the different rules preserve these properties. The easiest case is the one of a query consisting of a single step, handled by the Rule **(p-step)**. In this rule, we just apply the type inference algorithm to determine the output type of the results. The resulting projector is the set of rules in the results to which we add their upward context  $\kappa$ , that is the rules linking the results to a start symbol. The rules **(p-union)** and **(p-iterate)** are only inductive cases which allows us to handle respectively top-level union and projectors applied to a set of rules. In particular, **(p-iterate)** splits the checking of a path over all the possible rules specified in the type component of the environment (each one identifies a different set of current nodes). This allows us to define the

## Base and induction

$$\begin{array}{l}
\text{(p-step)} \frac{\Sigma \vdash_E \text{Step} : (\tau, \kappa)}{\Sigma \Vdash_E \text{Step} : \tau \cup \kappa} \text{ if } \tau \neq \emptyset \quad \text{(p-union)} \frac{\Sigma \Vdash_E P_1 : \tau_1 \quad \Sigma \Vdash_E P_2 : \tau_2}{\Sigma \Vdash_E P_1 \mid P_2 : \tau_1 \cup \tau_2} \\
\text{(p-erase)} \frac{}{\Sigma \Vdash_E P : \emptyset} \text{ if no other rule applies} \\
\text{(p-iterate)} \frac{(\{X_1 \rightarrow R_1\}, \kappa) \Vdash_E P : \tau_1 \quad \dots \quad (\{X_n \rightarrow R_n\}, \kappa) \Vdash_E P : \tau_n}{(\{X_1 \rightarrow R_1, \dots, X_n \rightarrow R_n\}, \kappa) \Vdash_E P : \bigcup_{i=1..n} \tau_i} \text{ if } n \geq 2
\end{array}$$

## Path Rules

$$\begin{array}{l}
\text{(p-test)} \frac{(\{Y \rightarrow R\}, \kappa) \vdash_E \text{self} :: \text{Test} : \Sigma \quad \Sigma \Vdash_E P : \tau}{(\{Y \rightarrow R\}, \kappa) \Vdash_E \text{self} :: \text{Test} / P : \{Y \rightarrow R\} \cup \tau} \text{ if } \Sigma_{\text{typ}} \neq \emptyset \\
\text{(p-predicate)} \frac{\Sigma \Vdash_E P : \tau \quad (\{Y \rightarrow R\}, \kappa) \vdash_E \text{self} :: \text{node} [C] : \Sigma \quad \forall P_i \in C \quad \Sigma \Vdash_E P_i : \tau_i \text{ if } \Sigma_{\text{typ}} \neq \emptyset}{(\{Y \rightarrow R\}, \kappa) \Vdash_E \text{self} :: \text{node} [C] / P : \{Y \rightarrow R\} \cup \tau \cup \bigcup_i \tau_i} \\
\text{(p-single)} \frac{(\{Y \rightarrow R\}, \kappa) \vdash_E \text{Axis} :: \text{node} : (\tau, \kappa') \quad (\text{for } i = 1..n) \quad (\{X_i \rightarrow R_i\}, \kappa') \vdash_E P : \Sigma^i \quad (\tau', \kappa') \Vdash_E P : \tau''}{(\{Y \rightarrow R\}, \kappa) \Vdash_E \text{Axis} :: \text{node} / P : \{Y \rightarrow R\} \cup \tau' \cup \tau''} (*) \\
(*) \text{ where } \text{Axis} \in \{\text{parent}, \text{child}\}, \tau = \{X_1 \rightarrow R_1, \dots, X_n \rightarrow R_n\}, \\
\tau' = \{X_i \rightarrow R_i \mid i = 1..n, \Sigma_{\text{typ}}^i \neq \emptyset\}, \tau \neq \emptyset, \text{ and } \tau' \neq \emptyset \\
\text{(p-many)} \frac{(\{Y \rightarrow R\}, \kappa) \vdash_E \text{Axis} :: \text{node} : (\tau, \kappa') \quad (\text{for } i = 1..n) \quad (\{X_i \rightarrow R_i\}, \kappa') \vdash_E \text{Axis} :: \text{node} / P : \Sigma^i \quad (\tau', \kappa') \Vdash_E \text{s}(\text{Axis}) :: \text{node} / P : \tau''}{(\{Y \rightarrow R\}, \kappa) \Vdash_E \text{Axis} :: \text{node} / P : \{Y \rightarrow R\} \cup \tau' \cup \tau''} (**) \\
(**) \text{ where } \text{Axis} \in \{\text{ancestor}, \text{descendant}\}, \\
\tau = \{X_1 \rightarrow R_1, \dots, X_n \rightarrow R_n\}, \\
\tau' = \{Y \rightarrow R\} \cup \{X_i \rightarrow R_i \mid i = 1..n, \Sigma_{\text{typ}}^i \neq \emptyset\}, \tau \neq \emptyset, \tau' \neq \emptyset, \\
\text{s}(\text{descendant}) = \text{child}, \text{ and } \text{s}(\text{ancestor}) = \text{parent}.
\end{array}$$

Fig. 2. Inference rules for type-projectors

so-called “Path Rules” in much a simpler way since they can be written for environments in which the type component is just a singleton. The Path Rules actually perform the projection and they all follow the same scheme. The Rule **(p-test)** handles a simple node test. If the type inference returns some *non-empty* type  $\Sigma_{\text{typ}}$  for the step, then we can compute the projector for the continuation  $P$  and add its result to the rule for the current node. The Rule **(p-predicate)** is similar: the type  $\Sigma_{\text{typ}}$  returned by the type inference is the set of nodes for which the predicate is satisfied. We then recursively compute the projector for the continuation  $P$  as well as for the paths  $P_{ij}$  occurring in the predicate. In the end, we return the union of all the computed projectors to which we add the rule for the current node. Again we only do this if the type inference returned a non-empty type. The following rules handle the actual navigation. They are split in two sets, one for the parent and child axes another for their recursive variant, ancestor and descendant. Since they are the most delicate rules let us explain them in detail. The two cases are similar. In the Rule **(p-**

**single**), the algorithm first retrieves all the rules matching the axis (`child` or `parent`). These rules are collected in  $\tau$  and the analysis yields a current context  $\kappa'$ . Then, by using  $n$  calls to the type inference algorithm ( $n$  being the number of rules in  $\tau$ ), it collects among  $\tau$  only the rules which are a suitable starting point for the rest of the path, that is all the rules yielding a non-empty result type when typed against  $P$ . These rules are collected in  $\tau'$  which, as it can be easily seen, a subset of  $\tau$ . Finally,  $\tau'$  and  $\kappa'$  are used as the environment to infer the projection with the rest of the path. The **(p-many)** rule handles the recursive axes, `descendant` and `ancestor`. The rule is almost the same as Rule **(p-single)** with the exception that it does not test whether the continuation  $P$  yields a non empty result on the node but on a `descendant` (or `ancestor`) of the node, to ensure that we put not only the correct rules in the projector but also the rules leading to them, and therefore that we maintain well-formedness. If for any of these rules one of the side conditions does not hold, then the rule **(p-erase)** is applied and returns an empty projector for the current path.

We illustrate hereafter the behavior of type-projection inference by unrolling it on an example. Consider the grammar:

$$\{A \rightarrow a[(B|C)*], B \rightarrow b[D], C \rightarrow c[E], D \rightarrow d[E], E \rightarrow e[]\}$$

with start symbol  $A$  and the path  $P$ :

$$\text{descendant::node/self::e/ancestor::node/self::b}$$

which is the single step normal form of

$$//e/ancestor::b$$

To ease the reading, we identify every rule with the non-terminal it defines. Therefore in what follows when we write, say,  $A$  in types or contexts, we actually mean  $A \rightarrow a[(B|C)*]$ . The algorithm computes the type projector for  $P$  as follows. The initial environment is  $(\{A\}, \{A\})$ . We apply the rule **(p-many)** for the first step. The first premise computes the type of `descendant::node` applied to  $A$ , which returns the type and context (these instantiate the  $(\tau', \kappa')$  of the rule):

$$(\{B, C, D, E\}, \{A, B, C, D, E\})$$

Then the second premise filters out the unwanted names and keeps only those for which the whole path may succeed. This gives us an intermediary type:  $\{B, D, E\}$  (and unchanged context) onto which we can compute the projector for the path:

$$\text{child::node/self::e/ancestor::node/self::b}$$

the final result for this rule will be the projector for the above path to which we add  $\{A\}$  (i). At this point, since the input type contains many rules, we can apply the rule **(p-iterate)** which will apply the continuation path on  $\{B\}$ ,  $\{C\}$  and  $\{D\}$ . It is easy to see that on  $\{B\}$  the side condition for the rule **(p-single)** is not fulfilled, since the type inference returns empty. The rule **(p-erase)** applies and returns an empty projector. The projection continues with only  $\{C\}$  and  $\{D\}$  left (the context is unchanged until now,  $\{A, B, C, D, E\}$ ). First let's consider the derivation for  $\{D\}$ . The current step is `child::node` which was introduced by the previous **(p-many)** rule. On this step, we apply the rule **(p-single)**. This rule adds  $\{D\}$  (ii) in the final projector and continues by computing a projector from  $\{E\}$  using the path:

$$\text{self::e/ancestor::node/self::b}$$

When we apply the same rule to  $\{C\}$  however, while the first premise returns a non empty type, the second one returns an empty result, since from a node with type  $C$  the path

$$\text{child::node/self::e/ancestor::node/self::b}$$

yields an empty result. Thus the rule is not applied and the result of the projector for the remaining path for the node type  $\{C\}$  is the empty projector. We continue with our only set,  $\{E\}$ . We compute the projector for `self::e` which adds  $\{E\}$  to the final projector (iii) and computes the projector for the path:

$$\text{ancestor::node/self::b}$$

It is easy to see that these will return  $\{B\}$  as a projector (iv). If we summarize, we obtain from (i, ii, iii, and iv) the set of rules

$$\pi = \{A \rightarrow a[(B|C)*], B \rightarrow b[D], D \rightarrow d[E], E \rightarrow e[\ ]\}$$

The actual type projector is:

$$(\mathcal{S} \cap \mathbf{Dn}(\pi), \{(X \rightarrow R) |_{\mathbf{Dn}(\pi)} \mid X \rightarrow R \in \pi\})$$

that is:

$$(\{A\}, \{A \rightarrow a[B*], B \rightarrow b[D], D \rightarrow d[E], E \rightarrow e[\ ]\})$$

This example shows how both precision and well-formedness are preserved:

**well-formedness:** what we obtained at the end is a valid type without unneeded rules.

**precision:** although the query references  $e$  nodes explicitly, we do not naively keep all the  $e$  nodes but only those that are useful to compute the query, namely those occurring *below* a  $b$  node.

We can now present the formal properties of type-projection inference

LEMMA 5.10 (TERMINATION OF TYPE-PROJECTOR INFERENCE). *Let  $(\mathcal{S}, E)$  be a type,  $P$  a path, and  $\Sigma$  and  $\Sigma'$  environments. The judgment  $\Sigma \Vdash_E P : \Sigma'$  has a unique and finite derivation. Furthermore the algorithm runs in time  $O(|E|^{o(d)} \times |P|^2)$  where  $d$  is the maximum number of nested predicates.*

The lemma above states that the rules in Figure 2 describe a terminating algorithm (the complexity result is a straightforward consequence of the one given by Lemma 5.4. If one ignores the recursive calls to the typing algorithm, it is clear that during projection inference, one perform at most  $O(|E|)$  operations and performs at most one recursive call to a strict sub-path of the argument path). We show now that they compute a type-projector by formalizing the “well-formedness” property that we outlined above. The intuition is that when the output type for a step is computed (e.g., in the first premise of the rule **(p-predicate)**), then the *context* corresponding to this computation is kept and passed as a parameter for the inference of the remainder of the path. On the last step, (rule **(p-step)**) the context is added to the type projector. There, it ensures that whenever a rule  $Y \rightarrow R$  is added to the type-projector, all the rules needed to derive  $Y \rightarrow R$  from the start symbols are added to the type-projector as well. This is what we formally state by the following lemma:

LEMMA 5.11 (WELL-FORMEDNESS OF TYPE-PROJECTOR INFERENCE). *Let  $(\mathcal{S}, E)$  be a type,  $\tau$ ,  $\tau'$ , and  $\kappa$  sets of rules, and  $P$  a path. If  $(\tau, \kappa) \Vdash_E P : \tau'$ , then  $(\tau, \kappa) \vdash_E P : (\tau'', \kappa'')$  implies  $\kappa'' \subseteq \tau'$ .*

We can now state the soundness of type-projection inference:

THEOREM 5.12 (SOUNDNESS OF TYPE-PROJECTOR INFERENCE). *Let  $(\mathcal{S}, E)$  be a type and  $P$  an XPath<sup>l</sup> query. Let  $S$  be the set of rules:  $S = \{X \rightarrow R \mid X \in \mathcal{S}\}$ . If*

$$(S, S) \Vdash_E P : \tau$$

*then  $\tau$  is a type-projector for  $(\mathcal{S}, E)$  and for every  $t \in_{\exists} (\mathcal{S}, E)$  we have:*

$$\llbracket P \rrbracket_{t \setminus_{\exists} \tau}(\mathbf{RootId}(t)) = \llbracket P \rrbracket_t(\mathbf{RootId}(t))$$

In words, if  $\tau$  is the projector inferred for a query  $P$  and a grammar  $(\mathcal{S}, E)$ , then for every tree  $t$  validating the grammar, the result of executing  $P$  on  $t$  or on its pruned version  $t \setminus_{\exists} \tau$  is the same.

Completeness requires not only completeness of the type system (\*-guarded, non-recursive, and parent-unambiguous DTDs), but also the following condition on queries:

*Definition 5.13.* An XPath query  $Q$  is *strongly-specified* if:

- i. its predicates do not use backward axes,
- ii. along  $Q$  and along each path in the predicates of  $Q$  there are no two consecutive (possibly conditional) steps whose *Test* part is node
- iii. each predicate in  $Q$  does not contain any or or and operator. Furthermore any path in a predicate does not terminate by a step whose *Test* is node.

For instance, among the following queries, only the first two are strongly-specified:

- descendant::node/self::a/ancestor::node
- descendant::node[child::b]/self::a/parent::node
- descendant::node/ancestor::node/self::a
- descendant::node[child::b/child::node]/self::a
- child::a[descendant::node/parent::b/child::c]

In the third query, there are two consecutive steps with a “node” test, which violates condition (ii). In the fourth query the predicate contains a path ending with “node”—failing to satisfy condition (iii)—and for the last query, the predicate contains backward axes, which violates condition (i).

Once more, we are in presence of a very common class of queries: for instance, almost all paths in the XMark and XPathMark benchmarks are strongly specified.

If all the conditions are met, then we can show that our algorithm is complete, in the sense that it infers the best possible sound projector. In words, if we remove any rule (and its consequences) from a projector inferred for a path  $P$  and a grammar  $(\mathcal{S}, E)$ , then we obtain a projector for which there exists a tree  $t$  validating the grammar for which the execution on  $t$  and on its pruned version yield different results. Formally:

**THEOREM 5.14 (COMPLETENESS OF PROJECTOR INFERENCE).** *Let  $(\mathcal{S}, E)$  be a non-recursive, parent-unambiguous and \*-guarded local tree grammar, and let  $P$  be a strongly-specified XPath<sup>ℓ</sup> path. Let  $S$  be the set of rules:  $S = \{X \rightarrow R \mid X \in \mathcal{S}\}$ . If*

$$(S, S) \Vdash_E P : \tau$$

*then there exists  $t \in \mathcal{T}(\mathcal{S}, E)$  such that for each  $Y \rightarrow R \in \tau$ , if  $\pi = \tau \setminus (\{Y \rightarrow R\} \cup \mathbf{A}_E(\{Y \rightarrow R\}, \text{descendant}))$ , then:*

$$\llbracket P \rrbracket_{t \setminus \pi}(\mathbf{RootId}(t)) \neq \llbracket P \rrbracket_t(\mathbf{RootId}(t))$$

The fact that completeness may not hold for not \*-guarded, non-recursive, or parent-ambiguous local tree grammar, is a consequence of the analogous property of the type system. To see that also strong-specification is a necessary condition consider documents valid with respect to the following grammar rooted at  $X$ :

$$\{X \rightarrow a[Y, W], W \rightarrow c[], Y \rightarrow b[Z], Z \rightarrow d[]\}$$

If we query a document of that type with the following non strongly-specified query

$$\text{self}::a[\text{child}::\text{node}]$$

which does not satisfy (iii), then  $\{X, Y\}$  is an optimal projector for this query (once more, we use a name to denote the rule that defines it), but the presence of the condition `child::node` forces the system to include also  $W$  in the inferred projector, thus breaking completeness. A similar reasoning applies for `self::a[child::b or child::c]`, which does not satisfy condition (iii) because of the presence of an “or” in the predicate. Concerning the presence of backward axes in predicates, consider the query

$$\text{self}::a[\text{descendant}::\text{node}/\text{ancestor}::a]$$

which does not satisfy condition (i). An optimal projector for this query on the same grammar is  $\{X, Y\}$ . However, since the `ancestor` condition is true for all descendants of  $a$  nodes,  $\{W, Z\}$  is included in the projector as well. Finally, with a similar reasoning on the same grammar, it is clear that the query

$$\text{descendant}::\text{node}/\text{ancestor}::\text{node}/\text{self}::a$$

for which condition (ii) does not hold, jeopardises completeness. The first step selects all the rules in the grammar that can be derived from the start symbol (that is, all the rules). None of these rule are discarded by the projector-inference since for none of them, the output type of “`ancestor::node/self::a`” is empty. The point here is that for the given grammar, there is no need to keep all the nodes, but only one child of the root. Indeed, having one element below the root

guarantees that the sequence `descendant::node`, `ancestor::node` is not empty and therefore that the root can be selected.

Of course, it is possible to state completeness for other classes of queries but, once more, this seems a satisfactory compromise between simplicity and generality.

## 6. EXTENSION TO FULL XPATH

The formal developments of the previous section only deal with the XPath<sup>ℓ</sup> language. This language allows one to express *structural queries*, that is, queries whose predicates contain only conjunctions or disjunctions of paths. In this section we show how to translate a (full) XPath query into a (set of) XPath<sup>ℓ</sup> queries and perform type-projection inference for the latter that is sound for the former. In other terms, we show that our translation is a sound approximation with respect to type-projection. Finally, we also show how to encode the XPath axes not present in XPath<sup>ℓ</sup> and how to extend our theoretical framework to handle most XML and XPath peculiarities (attributes, absolute paths,...

### 6.1. Handling XPath predicates

We extend Definition 4.1 to XPath 1.0 paths ([XPath 1999]):

*Definition 6.1.*

A path is a finite production of the following grammar:

$$\begin{aligned}
 \text{Path} & ::= \text{Step} \mid \text{Path}/\text{Path} \mid \text{Path} \mid \text{Path} \\
 \text{Step} & ::= \text{Axis}::\text{Test} \mid \text{Axis}::\text{Test}[\text{Cond}] \\
 \text{Axis} & ::= \text{self} \mid \text{child} \mid \text{descendant} \mid \text{parent} \mid \text{ancestor} \\
 \text{Test} & ::= \text{tag} \mid \text{node} \mid \text{text} \\
 \text{Cond} & ::= \text{Cond} \text{ or } \text{Cond} \mid \text{Cond} \text{ and } \text{Cond} \mid \text{Expr} \\
 \text{Expr} & ::= \text{Expr} \text{ cmp } \text{Expr} \mid \text{Arith} \\
 \text{Arith} & ::= \text{Arith} \text{ op } \text{Arith} \mid \text{Atom} \\
 \text{Atom} & ::= f(\text{Expr}, \dots, \text{Expr}) \mid \text{Path} \mid v
 \end{aligned}$$

where *tag* ranges over element tags,  $\text{cmp} \in \{=, \neq, <=, <, >, >=\}$ ,  $\text{op} \in \{+, -, *, \text{div}, \text{mod}\}$ , *f* ranges over a set of built-in functions of the Core Function Library and *v* ranges over values (strings, sequences, integers,...).

We wish to provide a safe translation from an XPath query *Q* to an XPath<sup>ℓ</sup> query *P* that approximates *Q* and use it to infer a type projector. By safe we mean that the type-projector inferred for *P* must not change the semantics of *Q*.

What exactly is an approximating query in this context? A naive approach to define query approximation is to consider inclusion of the results. According to it the query *P* translation of *Q* should always select more nodes than *Q*. However this works only as long as we do not use non-structural conditions (that is, predicates that make a query be non-structural). This is clear for example when we use the negation function `not`. Consider

$$\text{descendant}::a[\text{not}(\text{child}::b)]$$

For all documents, the query `descendant::a` returns more results than the query above. However, a projector inferred for `descendant::a` would discard *b* nodes not occurring before an *a* node, and therefore possibly also some *b* nodes children of an *a* node. In this way it would change the result of the original query. What the approximating query needs to reflect cannot be defined in terms of inclusion of results but rather in terms of *data-need*. We must ensure that the approximation traverses at least the same nodes as the original one to ensure that the former will not be pruned. However, we want the approximation also to be as precise as possible. For instance “`descendant::node`” is a sound approximation for any XPath query but the projector we infer from it is utterly imprecise: it performs no pruning.

As the reader will have understood, the tricky part is to approximate non-structural conditions. We do it as follows:



*Definition 6.2 (Approximation of a path).* Let  $P$  and  $S$  respectively denote a path and a set of paths of XPath<sup>l</sup>. Let  $P/S$  denote the set of XPath<sup>l</sup> paths defined as  $\bigcup_{P' \in S} \{P/P'\}$ . Given an XPath expression  $Q$ , its *approximation*  $\mathbf{P}(Q)$  is the set of XPath<sup>l</sup> paths defined as:

$$\begin{aligned} \mathbf{P}(Q_1|Q_2) &= \mathbf{P}(Q_1) \cup \mathbf{P}(Q_2) \\ \mathbf{P}(Axis::Test/Q) &= Axis::Test/\mathbf{P}(Q) \\ \mathbf{P}(Axis::Test[C]/Q) &= Axis::Test[\mathbf{C}(C)]/\mathbf{P}(Q) \cup Axis::Test/S(C) \end{aligned}$$

where:

$$\begin{aligned} \mathbf{C}(P) &= P && \text{if } \mathbf{P}(P) = \{P\} \\ \mathbf{C}(P) &= \text{self}::\text{node} && \text{if } \mathbf{P}(P) \neq \{P\} \\ \mathbf{C}(C_1 \text{ or } C_2) &= \mathbf{C}(C_1) \text{ or } \mathbf{C}(C_2) \\ \mathbf{C}(C_1 \text{ and } C_2) &= \mathbf{C}(C_1) \text{ and } \mathbf{C}(C_2) \\ \mathbf{C}(C) &= \text{self}::\text{node} && \text{otherwise} \end{aligned}$$

and:

$$\begin{aligned} \mathbf{S}(P) &= \emptyset && \text{if } \mathbf{P}(P) = \{P\} \\ \mathbf{S}(P) &= \mathbf{P}(P) && \text{if } \mathbf{P}(P) \neq \{P\} \\ \mathbf{S}(C_1 \text{ or } C_2) &= \mathbf{S}(C_1) \cup \mathbf{S}(C_2) && \mathbf{S}(C_1 \text{ op } C_2) = \mathbf{S}(C_1) \cup \mathbf{S}(C_2) \\ \mathbf{S}(C_1 \text{ and } C_2) &= \mathbf{S}(C_1) \cup \mathbf{S}(C_2) && \mathbf{S}(C_1 \text{ cmp } C_2) = \mathbf{S}(C_1) \cup \mathbf{S}(C_2) \\ &&& \mathbf{S}(f(C_1, \dots, C_n)) = \mathbf{F}(f(C_1, \dots, C_n)) \end{aligned}$$

and  $\mathbf{F}(\_)$  is the approximation of built-in functions (see Figure 3 for an excerpt).

The most technical point in the definition above is, as expected, the approximation of conditions, implemented by the auxiliary functions  $\mathbf{C}()$  and  $\mathbf{S}()$ . To be precise, we differentiate between purely structural paths and non structural paths. For a structural path,  $\mathbf{P}()$  does not introduce any approximation and returns the singleton containing the path itself. Otherwise, a non-structural path is approximated by a set of paths. The translation is non trivial when the path contains non structural conditions. Let us illustrate the rationale of the definition first by an example. The path

$$\text{descendant}::a[(\text{count}(\text{child}::b)>3 \text{ and } \text{child}::c) \text{ or } \text{descendant}::b]/\text{child}::d$$

is approximated by the following set of two paths

$$\{ \text{descendant}::a[\text{self}::\text{node} \text{ and } \text{child}::c \text{ or } \text{descendant}::b]/\text{child}::d, \\ \text{descendant}::a/\text{child}::b \}$$

The first is generated by an application of the function  $\mathbf{C}()$ , while the second derives from the application of  $\mathbf{S}()$ . As we see, the arithmetic expression  $\text{count}(\text{child}::b)>3$  is approximated by the function  $\mathbf{C}()$  into the  $\text{self}::\text{node}$  path occurring in the first path of the set. This condition is always true and therefore it is a sound approximation of the Boolean value of the expression (since the result is always true the type-inference algorithm will never be able to deduce an empty output type for this sub-path and therefore the type-projector inference algorithm will keep the rules associated with this node). However this is not sufficient to ensure the safety of type projection. Indeed for this test to be possible at run-time, the projected document must have the “ $b$ ” nodes that were below  $a$  nodes in the original document. This approximation is made via the second path by the  $\mathbf{S}()$  function and, in particular, by the  $\mathbf{F}(\text{count}(\text{child}::b))$ . Of course, what the function actually does depends on the semantics of the built-in function. For instance,  $\text{count}(P)$  returns the number of nodes selected by  $P$ , thus a projector keeping the type of the nodes selected by  $P$  is sound. On the contrary, the function  $\text{string}(P)$  when applied to a node set returns the concatenation of the string-value of all the nodes in the set. The string-value of a node is the concatenation of all the **PCDATA** elements occurring below it. Therefore a suitable approximation for  $\text{string}(P)$  is not  $P$  but rather  $P/\text{descendant}::\text{text}$ . Giving an approximation for all the functions of the XPath Core Library is a tedious task. Although our prototype implements approximation for all functions, in Figure 3 we just give an excerpt that completely covers all the different techniques we used in our prototype to approximate built-in functions.

$$\begin{aligned}
\mathbf{F}(\text{count}(Q)) &= \mathbf{P}(Q) \\
\mathbf{F}(\text{last}(Q)) &= \mathbf{P}(Q) \\
\mathbf{F}(\text{position}(Q)) &= \mathbf{P}(Q) \\
\mathbf{F}(\text{string}(Q)) &= \mathbf{P}(Q)/\text{descendant-or-self}::\text{node}[\text{text}()] \\
\mathbf{F}(\text{number}(Q)) &= \mathbf{P}(Q)/\text{descendant-or-self}::\text{node}[\text{text}()] \\
\mathbf{F}(\text{not}(Q)) &= \mathbf{P}(Q) \\
\mathbf{F}(\text{true}()) &= \{\text{self}::\text{node}\} \\
\mathbf{F}(\text{false}()) &= \{\text{self}::\text{node}[\text{self}::a \text{ and } \text{self}::b]\} \\
\mathbf{F}(f(v_1, \dots, v_n)) &= \{\text{self}::\text{node}\} \quad \text{where } v_i \text{ is a value}
\end{aligned}$$

Fig. 3. Approximation of XPath functions

## 6.2. Other XPath features

We purposely left out from our definitions some features of XPath that would have led to a much more intricate formalization process, in particular for what concerns definitions of the algorithms and the proofs of the theorems. Here we illustrate how these features can be either encoded or approximated within our framework.

descendant-or-self and ancestor-or-self axes: These axes—that we used in Figure 3—can be encoded exactly by using the “|” operator. Precisely:

$$\text{descendant-or-self}::\text{Test}[Cond]$$

can be equivalently written as:

$$(\text{descendant}::\text{Test}[Cond] \mid \text{self}::\text{Test}[Cond])$$

Sibling axes: We could have defined a sibling relation over node identifiers in the same way as we defined the edge relation *Edg* in Section 4, and used it to deal with the following-sibling and preceding-sibling axes natively. However we can also approximate these axes using only “vertical” moves. So for instance:

$$\text{following-sibling}::\text{Test}[Cond]'$$

becomes:

$$\text{parent}::\text{node}/\text{child}::\text{Test}[Cond]'$$

The transformation above approximates the following siblings of a node by all its siblings, including itself. Our experiments showed that, as far as type-projection is concerned, this kind of approximation does not yield any noticeable loss of precision in practice.

preceding and following axes For these axes, we can directly use the W3C recommendation [XPath 1999] and encode *following* accordingly. That is,

$$\text{following}::\text{Test}[Cond]$$

becomes:

$$\text{ancestor}::\text{node}/\text{following-sibling}::\text{node}/\text{descendant-or-self}::\text{Test}[Cond]$$

Document node The XPath data model enforces the presence of a *document node*, the real root of the document which has no label and is selected by the initial “/” of an XPath expression. It is of course possible to represent such documents in our framework but we preferred to omit it here since it would cause many presentation issues with little theoretical interest. In particular, the document node is never referenced by the schema of the document.

Absolute paths Absolute paths are paths with a leading /. They do not start their evaluation from the current context node but from the root of the document. Our formalism easily allows us to encode absolute paths. First, if an absolute path occurs outside of a predicate, as in:

$$P/(P_1 \mid P_2)/P'$$

then we can simply rewrite it as:

$$(P/P_1/P') \mid (P_2/P')$$

Second, if the path /P occurs in a predicate, then we can replace it with *self::node* (as if it was a non structural condition) and add  $\mathbf{P}(P)$  to the global approximation. Direct treatment of absolute paths would have further complicated Definition 6.2, where we would have had to maintain a set of

absolute approximations, modified only by absolute paths and propagated at each function call. We chose not to clutter this definition (but absolute paths are handled by our implementation).

attribute axis and attributes in the data-model and schema: Conceptually, the handling of the attribute axis is not very different from the child axis, and could be encoded as such. For instance, a possible solution would be to encode an element

$$\langle e \text{ att}="value" \text{ id}="34" \rangle \langle a \rangle \langle b \rangle \langle /e \rangle$$

as the tree:

$$e[ @[att["value"] id["34"]] a[] b[] ]$$

by introducing a phony node with label @. If such a solution were retained then we would also need to update the definitions of child and descendant to ignore @ nodes, and add an attribute axis selecting only the content of such nodes.

As far as schemas are concerned, they need to reflect the *uniqueness* and *unorderedness* of a sequence of attributes within an element node. This can be done with a union type. For instance, the document above could have type:

$$\begin{array}{ll} E \rightarrow e[ ATTS A B ] & ID \rightarrow id[String] \\ ATTS \rightarrow @[ (ATT ID) | (ID ATT) ] & A \rightarrow a[] \\ ATT \rightarrow att[String] & B \rightarrow b[] \end{array}$$

this encoding however incurs an exponential blow-up in the size of the sequence of attributes. Our implementation follows a much more pragmatic approach. Precisely, even though attributes could be encoded in our approach we preferred to add an unordered attribute construct directly at the grammar level and specialize type-inference and type-projector inference rules for attributes.

id() function: The id() function of XPath is peculiar in the sense that unlike other functions, it does not take the context node as implicit argument (e.g., the position() function returns the position of the context node within the current result set). Rather, the expression "id("foo")" returns the node whose id is "foo" if it exists (a node has id "foo" if it has an attribute named id whose value is "foo" and if this attribute has been declared with type ID in the Schema, [XPath 1999]). We choose to approximate this function in two steps. First, we rewrite it as an absolute path. Then we can let our approximation algorithm handle the absolute path (with the technique described in Section 6.2). For instance an expression such as

$$id("item34")/child::name$$

can be rewritten as

$$/descendant::*[@id="item34"]/child::name$$

This rewrite technique was used in particular to handle queries C5-C7 of the XPathMark benchmark (see Section 9).

## 7. EXTENSION TO XQUERY

In this section we extend our technique to XQuery.

*Definition 7.1* (*XQuery*).

$$\begin{aligned}
FLOWR & ::= FORLET \text{ return } ExprS \mid ExprS \\
FORLET & ::= FOR \mid LET \\
FOR & ::= \text{ for } \$x \text{ in } ExprS \\
LET & ::= \text{ let } \$x := ExprS \\
ExprS & ::= \text{ if } ExprS \text{ then } Cond \text{ else } Cond \mid Cond \\
Cond & ::= Cond \text{ or } Cond \mid Cond \text{ and } Cond \mid Expr \\
Expr & ::= Expr \text{ cmp } Expr \mid Arith \\
Arith & ::= Arith \text{ op } Arith \mid Atom \\
Atom & ::= f(Expr, \dots, Expr) \mid FLOWR/P \mid x \mid v \\
& \quad \mid FLOWR, FLOWR \mid \langle \text{tag} \rangle FLOWR \langle / \text{tag} \rangle \mid () \\
P & ::= Step[Cond]/P \mid Step/P \mid Path
\end{aligned}$$

where  $x$  ranges over identifier names, and  $v$ ,  $cmp$ ,  $op$ ,  $Path$  and  $Step$  are the same as in Definition 6.1.

For the sake of clarity and concision we only considered formally a subset of the XQuery grammar ([XQuery 2010]). In a nutshell, the definition of *Atom* (given in Section 6) is extended with two new constructs: *variables* (ranged over by  $x, y, z$  in what follows) and path applications *FLOWR/P*.

Note that XQuery constructs may occur inside a path expression (production *P*) or not (production *Path*). Also, we consider neither queries that first construct new elements and then navigate on them nor queries containing “order by” constructs. XQuery queries are ranged over by  $q$ . The absence of the former is not really a problem. For instance it is shown in [Koch 2005] that *composition-free* XQuery —where such queries are absent— is equivalent to Core XQuery, when considering only child and descendant axes.

In order to apply the previous analysis to infer a projector for an XQuery query  $q$ , we first extract a set of full XPath expressions from  $q$ . Then, we apply to each of these extracted paths the approximation function  $\mathbf{P}(\_)$  given in Definition 6.2 to obtain an XPath <sup>$l$</sup>  expression. We can finally use the projector inference algorithm of Section 5.2 on the set of approximated paths, which is a sound type projector for the original XQuery query  $q$ .

Path extraction is performed by the extraction function  $\mathbf{E}(\_, \_, \_)$ , whose definition is given in Figure 4. The extraction function has the form  $\mathbf{E}(q, \Gamma, m)$  and performs a straightforward recursive descent over its first parameter  $q$  which is the query at issue. The second parameter  $\Gamma$  is an environment, that keeps track of bindings of the form  $(x; P)$  in whose scope  $q$  occurs. Finally,  $m$  is a flag indicating whether  $q$  is a query that serves to materialise the full content of the queried elements ( $m = 1$ ) or if the query just selects a set of nodes whose descendants are not needed ( $m = 0$ ). Before explaining the rules of Figure 4 in details, we introduce two auxiliary functions. The first one is  $\mathbf{M}(\_, \_)$  (Figure 5) which given a built-in XPath function and the position of one of its arguments, returns a suitable value for the parameter  $m$  (intuitively,  $\mathbf{M}(f, i)$  returns 1 if  $f$  needs the full content of its  $i^{\text{th}}$  arguments and 0 otherwise). This function is similar to the function  $\mathbf{F}(\_)$  introduced in Section 6, Figure 3.

The second one,  $\mathbf{E}'(\_, \_, \_)$  is defined mutually, together with  $\mathbf{E}(\_, \_, \_)$  and allows to recursively traverse XQuery expressions and resolves the variable names they contain. It works similarly to  $\mathbf{E}(\_, \_, \_)$  but do not return sets of XPath paths, but sets of particular XQuery expressions which do not contain any variables.

Now that we have introduced environments and the auxiliary functions, we can easily describe the rules in Figure 4. First, rules 1 and 2 form the basic case of the recursive descent and return the empty set if the whole query consists of a constant. Rules 3 and 4 straightforwardly apply the extraction recursively for the content of sequence (Rule 3) and element (Rule 4) constructors. Rules 5 and 6 handle the case of variables bound in the environment Rules 7 and 8 add a constant path to the set of extracted paths, according to the value of the parameter  $m$ . Note that in those rules, *Path* refers to the corresponding entry in the grammar of Definition 6.1, that is it does not contain any XQuery construct and only pure XPath ones. Path containing XQuery expressions are handled in the subsequent rules. Rule 9 handles the application of a path *FLOWR* expression with a path *P*.

1.  $\mathbf{E}(\cdot, \Gamma, m) = \emptyset$
2.  $\mathbf{E}(v, \Gamma, m) = \emptyset$
3.  $\mathbf{E}((q_1, q_2), \Gamma, m) = \mathbf{E}(q_1, \Gamma, m) \cup \mathbf{E}(q_2, \Gamma, m)$
4.  $\mathbf{E}(\langle \text{tag} \rangle q \langle / \text{tag} \rangle, \Gamma, m) = \mathbf{E}(q, \Gamma, m)$
5.  $\mathbf{E}(x, \Gamma, 1) = \bigcup_{(x;P) \in \Gamma} \{P/\text{descendant-or-self} :: \text{node}\}$
6.  $\mathbf{E}(x, \Gamma, 0) = \bigcup_{(x;P) \in \Gamma} \{P\}$
7.  $\mathbf{E}(\text{Path}, \Gamma, 1) = \{\text{Path}/\text{descendant-or-self} :: \text{node}\}$
8.  $\mathbf{E}(\text{Path}, \Gamma, 0) = \{\text{Path}\}$
9.  $\mathbf{E}(\text{FLOWR}/P, \Gamma, m) = \mathbf{E}(\text{FLOWR}, \Gamma, m) / \mathbf{E}(P, \Gamma, m)$
10.  $\mathbf{E}(\text{Step}/P, \Gamma, m) = \text{Step} / \mathbf{E}(P, \Gamma, m)$
11.  $\mathbf{E}(\text{Step}[\text{Cond}]/P, \Gamma, m) = (\bigcup_{q \in \mathbf{E}'(\text{Cond}, \Gamma, m)} \text{Step}[q]) / \mathbf{E}(P, \Gamma, m)$
12.  $\mathbf{E}(\text{if } q \text{ then } q_1 \text{ else } q_2, \Gamma, m) = \mathbf{E}(q, \Gamma, 0) \cup \mathbf{E}(q_1, \Gamma, m) \cup \mathbf{E}(q_2, \Gamma, m)$
13.  $\mathbf{E}(\text{let } \$x := q_1 \text{ return } q_2, \Gamma, m) = \mathbf{E}(q_2, \Gamma \cup \Gamma', m)$   
where  $\Gamma' = \{(x;P) \mid P \in \mathbf{E}(q_1, \Gamma, 0)\}$
14.  $\mathbf{E}(\text{for } \$x \text{ in } q_1 \text{ return } q_2, \Gamma, m) = \mathbf{E}(q_1, \Gamma, 0) \cup \mathbf{E}(q_2, \Gamma \cup \Gamma', m)$   
where  $\Gamma' = \{(x;P) \mid P \in \mathbf{E}(q_1, \Gamma, 0)\}$

- 1'.  $\mathbf{E}'(\text{Cond}_1 \text{ op } \text{Cond}_2, \Gamma, m) = \bigcup_{q \in \mathbf{E}'(\text{Cond}_1, \Gamma, m)} \bigcup_{q' \in \mathbf{E}'(\text{Cond}_2, \Gamma, m)} \{q \text{ op } q'\}$   
where  $\text{op} \in \{\text{and, or}\}$
- 2'.  $\mathbf{E}'(\text{Expr}_1 \text{ cmp } \text{Expr}_2, \Gamma, m) = \bigcup_{q \in \mathbf{E}'(\text{Expr}_1, \Gamma, m)} \bigcup_{q' \in \mathbf{E}'(\text{Expr}_2, \Gamma, m)} \{q \text{ cmp } q'\}$   
where  $\text{cmp} \in \{=, !=, <, >, >=, <= \}$
- 3'.  $\mathbf{E}'(\text{Arith}_1 \text{ op } \text{Arith}_2, \Gamma, m) = \bigcup_{q \in \mathbf{E}'(\text{Arith}_1, \Gamma, m)} \bigcup_{q' \in \mathbf{E}'(\text{Arith}_2, \Gamma, m)} \{q \text{ op } q'\}$   
where  $\text{op} \in \{+, -, *, \text{div}, \text{mod}\}$
- 4'.  $\mathbf{E}'(f(\text{Expr}_1, \dots, \text{Expr}_n), \Gamma, m) = \bigcup_{q_1 \in \mathbf{E}'(\text{Expr}_1, \Gamma, \mathbf{M}(f, 1))} \dots \bigcup_{q_n \in \mathbf{E}'(\text{Expr}_n, \Gamma, \mathbf{M}(f, n))} \{f(q_1, \dots, q_n)\}$
- 5'.  $\mathbf{E}'(\text{Atom}, \Gamma, m) = \mathbf{E}(\text{Atom}, \Gamma, m) \quad \text{Atom} \neq f(q_1, \dots, q_n)$

Fig. 4. XQuery path extraction

Note that as previously the notations  $S_1/S_2$  where  $S_1$  and  $S_2$  are sets of paths stands for:

$$\bigcup_{P \in S_1} \bigcup_{P' \in S_2} \{P/P'\}$$

The case of a simple step composed with a path expression is handled similarly by Rule 10 and we recall that the notation  $\text{Step}/S$  where  $S$  is a set of path is syntactic sugar for the set:

$$\bigcup_{P \in S} \{\text{Step}/P\}$$

Rule 11 is more intricate, but its complexity is only bureaucratic. This rule allows the extraction process to retrieve Full XPath expressions from a FLOWR expression. In the present work, path

$$\begin{array}{lll} \mathbf{M}(\text{count},1) = 0 & \mathbf{M}(\text{string},1) = 1 & \mathbf{M}(\text{position},1) = 0 \\ \mathbf{M}(\text{last},1) = 0 & \mathbf{M}(\text{number},1) = 1 & \mathbf{M}(\text{not},1) = 0 \end{array}$$

Fig. 5. Value of the parameter  $m$  for various built-in XPath functions

extraction and path approximation are two separate processes. Path extraction only occurs at the level of XQuery terms and returns sets of full XPath expressions (which reflects the exact path that may be evaluated during query execution). Approximation from XPath to XPath<sup>ℓ</sup> is handled at the XPath level. The issues solved by Rule 11 is to recursively traverse an XQuery expression, using a recursive call to the auxiliary function  $\mathbf{E}'(\_,\_,\_)$  which builds a set of XPath conditions into which all variable bindings have been resolved. Therefore what we obtain after  $\mathbf{E}'(\_,\_,\_)$  is a set of XPath conditions free of any XQuery construct (especially variables). We can now explain how  $\mathbf{E}'(\_,\_,\_)$  works. In Rules 1' to 3' use a recursive descent into the production of the XQuery grammar, starting at the condition levels and reconstruct Boolean XPath condition (1'), relational XPath expressions (2') or arithmetic XPath expressions (3'). More interesting is Rule 4' which traverses the arguments of a function call and uses the auxiliary function  $\mathbf{M}(\_)$  to determine a suitable value for  $m$ . Lastly, if the input matches any other constructs Rule 5' applies and recursively applies  $\mathbf{E}(\_,\_,\_)$  to construct a set of XPath paths.

We can resume our description of  $\mathbf{E}(\_,\_,\_)$  for the remaining cases, the high level constructs “if then else”, “let return” and “for return” handled by Rules 12, 13 and 14 respectively. Rule 12 recursively extracts paths on the Boolean test  $q$ , the “then” case  $q_1$  and the “else” case  $q_2$ . The only point of interests is that the Boolean test cannot generate a result and therefore can be called with parameter  $m = 0$ . The let binding handled by Rule 13 augments the environment  $\Gamma$  with the path extracted from  $q_1$  and extracts the paths of query  $q_2$  in this augmented environment. Note that the path bound to  $x$  are added to the final results by Rule 5 or 6 only if the variable is used. On the contrary in Rule 14, for loops will perform their iterations even if the bound variable is never used, as long as the paths extracted from  $q_1$  yield a non-empty result. It is therefore mandatory to add the paths extracted from  $q_1$  to the final result.

These rules subsume and enhance the technique of Marian and Siméon [Marian and Siméon 2003]. In particular, (i) the technique we use to exclude useless intermediate paths is simpler and more compact, (ii) we do not need to distinguish between two kinds of extracted paths but, more simply, we always manage a unique set of path expressions, and last but not least, (iii) our path extractor can be used even if the user cannot access an XQuery to XQuery-Core compiler, which is necessary for [Marian and Siméon 2003].

Before applying the extraction function  $\mathbf{E}(\_,\_,\_)$  to some query  $q$  we apply some heuristics that rewrite  $q$  so as to improve the pruning capability of the inferred paths. Among these heuristics the most important is the one that rewrites

```

for y in Q/descendant-or-self::node
  return if C(y) then q else ()
into
for y in
  Q/descendant-or-self::node[C(self::node)]
  return q

```

whenever  $C(y)$  is a condition referring only to  $y$  and does not use external functions ( $C(\text{self}::\text{node})$  is obtained by replacing  $\text{self}::\text{node}$  for all occurrences of  $y$  free in  $C$ ). If we apply  $\mathbf{E}(\_,\_,\_)$  to the first query, then a path ending by the step

descendant-or-self::node

is extracted thus annulling further pruning: the entire forest selected by  $Q$  is loaded in main memory. This also happens with the approaches of Bressan *et al.* [Bressan et al. 2005] and of Marian and Siméon [Marian and Siméon 2003]. In ours and Marian and Siméon’s approach the query can be rewritten as above, while this is not possible with Bressan *et al.* formalisms since their subset of XQuery does not include predicates. However, Marian and Siméon’s path based pruning

degenerates (no further pruning is performed) also for the second query, since the step

`descendant-or-self::node`

ends up in the set of pruner paths, thus selecting all nodes. This is because their approach cannot manage predicates. In our approach, predicates are taken into account and therefore only nodes satisfying  $C(y)$  are kept by the projector, thus yielding a very precise pruning.

It is important to stress that despite their specific form the first kind of query is very common in practice since they are generated from XQuery to XQueryCore compilation of a non negligible class of queries or when rewriting upward axes into downward ones. This latter observation shows that the application of rewriting rules of [Olteanu et al. 2002] to extend Marian and Siméon’s approach to upward axes is not feasible since the rewriting may completely compromise pruning.

## 8. EXTENSION TO OTHER TYPING POLICIES

### 8.1. Handling un-typed documents

Although the usage of schema is being more and more wide-spread, it still is interesting to see how to perform type-based projection in an untyped world. A first, rather blunt, approach is to consider a fixed corpus of un-typed documents. For such sets of documents it is possible to *infer* a DTD. For instance, Bex *et al.* propose several automata-based methods to infer DTDs [Bex et al. 2006] and even XMLSchemas ([Bex et al. 2007; Bex et al. 2008]). Once a schema is inferred, our technique can be applied as-is.

More interestingly, this untyped problem can be reduced to a precise typing problem. Indeed, an un-typed document is nothing but a document of type  $(\{X\}, \{X \rightarrow Any\})$ . If we apply the type inference-algorithm of Section 5.1 to such an input type, then the result would be  $(\{X\}, \{X \rightarrow Any\})$  itself (meaning that the nodes selected by the query have type *Any*). Therefore in this case, since none of the intermediary steps of the query results in an empty-type, the type-projector inference algorithm of Section 5.2 cannot remove any rule from the input type which remains  $(\{X\}, \{X \rightarrow Any\})$ : the input document cannot be pruned. However, even though the input type does not contain any meaningful information, the query itself might. Imagine a query `//a/b`. It is easy to deduce, by a simple examination of the query a projector which keeps only “b” nodes occurring below “a”-nodes. While the solution in this case is straightforward, solving this problem in general is a tricky issue. The solution for a forward fragment of XPath can be found in [Nguyễn 2008] (cf. Chapter 7). Let us briefly outline it on our example. The first issue is the representation of types. For such precise algorithms, regular tree grammars are not well suited. Indeed instead of the type  $(\{X\}, \{X \rightarrow Any\})$ , it is more desirable to have a type  $(\{X\}, E)$  where  $E$  is the set of rules:

$$\{X \rightarrow String, X \rightarrow \_ [X*]\}$$

where  $\_$  denotes the set of all possible tree labels. Then the result of a query `//a/b` applied to a a tree of the above type (i.e., any XML tree) would be the type projector:

$$\{X \rightarrow \_ [X+], X \rightarrow a [B+], B \rightarrow b [Y*], Y \rightarrow String, Y \rightarrow \_ [Y*]\}$$

Note that this type-projector is non-deterministic top-down. It matches (and therefore keeps) any subtree  $t$  if and only if there is a subtree  $t'$  of  $t$  with tag “a” which itself has a non-empty sequence of children tagged “b”. Nodes that are children of an “a” node but whose tag is not “b” do not have any interpretation and therefore are discarded.

In [Nguyễn 2008], in order to achieve such a precise typing, the inference algorithm makes a heavy use of CDuce’s type algebra (see [Frisch et al. 2008] or, for an overview, [Castagna and Frisch 2005]) in particular of intersection and negation types. Also note that the projector above is not obtained by erasing some of the rules of the original type but by the mean of set theoretic operations. In fact, three new rules were created and intuitively they were obtained by intersecting the initial “ $X \rightarrow \_ [X*]$ ” rules with a type “ $a [B+]$ ” which is the constraint represented by the XPath query. Note also that contrary to our approach where new rules are only erasures of existing ones (of which there only exists a finite number), special care must be taken to not introduce infinitely many refined rules and ensuring the termination of the algorithm becomes a very delicate issue.

## 8.2. Using regular tree languages as schemas

While our formal development remains in the very general case of regular tree grammars, our implementation only focuses on DTDs. The main reason is that for DTDs pruning is efficient memory-wise. For regular tree languages instead, validation (and pruning) may need to visit the whole tree before deciding which node to prune. At first, it seems that this completely defeats the purpose of pruning, but we argue that pruning can still be of practical use in these cases.

Indeed, a way of addressing this problem is to temporarily store the document in memory in the form of a succinct tree data-structure (based for instance on balanced parenthesis: a survey of the most popular succinct tree representations can be found in [Arroyuelo et al. 2010]). The final data-structure of the document can then be built from the temporary one, by replaying a sequence of SAX events while traversing the temporary data-structure and by not synthesizing events for pruned subtrees. An alternative solution is to store on disk the sequence of SAX events and process it backward, thus simulating a bottom-up evaluation (validation of regular tree grammars and therefore projection can be done in a deterministic bottom-up fashion). Such a technique was used in [Koch 2003] to efficiently evaluate node selecting queries bottom-up on documents up to 1 GB of size.

## 9. EXPERIMENTS

### 9.1. Prototype

To gauge the benefits of type-based projection, we have implemented our pruning algorithm into a prototype (using the OCaml language, and the PXP library for XML and DTD parsing). Our prototype takes as input an XQuery query, a DTD, and a document. It then performs the path extraction described in Section 7 and computes for each extracted path its XPath<sup>ℓ</sup> approximation, applying the rewriting rules given in Section 6. Based on this set of paths, our program performs the static analysis described in Section 5.2 and computes a type projector. Once this is done, the prototype parses the input document, prunes it according to the inferred type-projector and serializes the result.

Besides what is included in the formal description of our algorithm, our prototype is extended also to support the full set of XPath axes as well as attributes. If we call  $D$  the input document and  $S$  the input DTD, then, assuming that  $D$  is well-formed, the pruning process is performed in  $O(|D|)$  time and  $O(|S|)$  memory, where  $|D|$  and  $|S|$  denote the size of  $D$  and  $S$  respectively. Indeed, the type projector associated with a DTD is at most as big as the original schema (when no pruning is performed) and  $O(|S|)$  space is required to store it in memory. Our prototype can also perform well-formedness check and validation while pruning, in which case time complexity remains  $O(|D|)$  and memory complexity becomes  $O(|S| + \log(|D|))$  (it is well-known that checking well-formedness during validation requires to keep a stack whose size is at most the height of the document, see for instance [Segoufin and Vianu 2002]).

### 9.2. Benchmark suite

We used the XPathMark ([Franceschet 2005; XPathMark ]) and XMark ([Schmidt et al. 2002]) benchmark suites. The former consists of a large set of XPath queries while the latter provides XQuery queries to test against.

*9.2.1. Data-set.* Both XPathMark and XMark use the XMark document generators. These documents comply with the “auction” DTD representing an auction web-site. It defines 77 element types and 15 attributes. This size and complexity is comparable to “real-life” type definitions (for instance the XHTML transitional DTD also features 77 element definitions). Because the “auction” DTD falls outside the conditions of our completeness theorem (it features recursion and unguarded union), it is a very good test-case to illustrate the precision we achieve in practice even when completeness does not hold. The scalability of our approach was tested by using documents of varying size, ranging from 12MB to 3GB. An important aspect of the XMark generator is that the proportion of textual data versus tree structure stays the same, for all size of documents. We report here some statistics of interest which we use later-on to gauge the precision of our pruning algorithm. An XMark generated file consists of:



- 74% of text content (as PCDATA elements or attribute value)
- 65% of all the text content (that is, 49% of the total file size) resides in a description element or one of its descendants.

*9.2.2. XPathMark queries.* Since its original publication (in [Franceschet 2005]), the XPathMark benchmark suite has evolved to provide a very complete set of XPath queries. It provides in particular *performance test* queries which provide computationally difficult queries. We highlight some of the main design goals of this test suite (the complete rationale can be found in [XPathMark ]):

- (1) queries simulate *realistic* query needs of a potential user of the the auction site;
- (2) queries are divided into groups according to the intrinsic computational complexity of the corresponding evaluation problem. The XPath language can be stratified in a number of fragments for which different complexity bounds are known [Benedikt and Koch 2008]. Comparing the theoretical computational complexity of the query evaluation problem with the actual amount of resources consumed during query evaluation might be, at least, a stimulating and instructive exercise;
- (3) queries are defined to challenge data scalability of the XML processing system, that is the performance of the system as the data complexity (document size) grows. In particular, the queries talk about document sections (like open and closed auctions, items, people, descriptions) that become bigger when the XMark document scaling factor increases. Moreover, the results of the queries are small compared to the size, to prevent serialization to obfuscate obfuscates the pure query processing time.

These three points are exactly those we aim to address with the present work. Indeed our approach drastically increases the data scalability (3.) of XML processing systems for realistic queries (1.) and potentially complex ones (2.). XPathMark queries are divided in 5 groups, labeled from *A* to *E* that we briefly describe now.

Group *A* contains unary tree pattern queries. These queries use only child and descendant axes, node tests equal to `*` or to a tag name, and filters (predicates). Conjunctive and disjunctive Boolean operators are allowed, but negation, relational and arithmetic operators and functions are disallowed. These queries fall therefore in the category of queries that we handle without any approximation.

Group *B* contains the so-called core or navigational XPath queries. This fragment extends Group *A* by admitting all XPath axes and negation. It mostly corresponds to queries for which our algorithm introduces a very lightweight approximation (we only need to approximate negations and those axes we did not treat formally such as `preceding`).

Group *C* extends Group *B* with comparisons (`=`, `!=`, `<`, `>`, `>=`, `<=`) and the `id()` function.

Group *D* extends Group *C* by allowing all arithmetic operators and the aggregate functions `sum()` and `count()`.

Group *E* contains all XPath 1.0 queries. In particular, it extends Group *D* by allowing all functions (like `position()` and `contains()`).

*9.2.3. XMark test suite.* To validate the extension of our approach to XQuery and in particular the path extraction algorithm, we use queries from the XMark benchmark suite ([Schmidt et al. 2002]). These queries feature “`for`” expressions guarded by “`where`” conditions and make use of element constructor to format their results. The corresponding code for the queries under consideration is given in Figure 6.

### 9.3. Protocol

We have designed two experiments, based on two different XQuery engines to validate our approach. For each engine and each query we described in the previous section we applied the following protocol. First, we tested the engine against original documents of increasing size and stopped when the query engine could not handle the input document anymore. Then we repeated the experiment a second time but used a document pruned by our prototype as input for the query engine.

```

A 1 /site/closed_auctions/closed_auction/annotation/description/text/keyword
A 6 /site/people/person[profile/gender and profile/age]/name
B 1 /site/regions/*/item[parent::america or parent::samerica]/name
B 2 //keyword/ancestor::listitem/text/keyword
C 3 /site/people/person[profile/@income =
      /site/open_auctions/open_auction/current]/name
C 4 /site/people/person[watches/watch/id(@open_auction)/seller/@person = @id]/name
D 1 /site/open_auctions/open_auction[(count(bidder) mod 2) = 0]/interval
D 2 count(//text) + count(//bold) + count(//emph) + count(//keyword)
E 5 /site/regions/*/item[preceding::item[100] and following::item[100]]/name
E 7 /site/regions/*/item[contains(substring-before(description,'eros'),'passion')
      and contains(substring-after(description, 'eros'), 'dangerous')]/name
M 3 for $b in $doc/site/open_auctions/open_auction
      where zero-or-one($b/bidder[1]/increase/text()) * 2
      <= $b/bidder[last()]/increase/text()
      return
      <increase
      first="$b/bidder[1]/increase/text()"
      last="$b/bidder[last()]/increase/text()"/>
M 6 for $b in $doc/site/regions return count($b//item)
M 7 for $p in $doc/site
      return
      count($p//description) + count($p//annotation) + count($p//emailaddress)
M 14 for $i in $doc/site//item
      where contains(string(exactly-one($i/description)), "gold")
      return $i/name/text()
M 15 for $a in
      $doc/site/closed_auctions/closed_auction/annotation/description/parlist/
      listitem/parlist/listitem/text/emph/keyword/text()
      return <text>$a</text>

```

Fig. 6. XPathMark (A-E) and XMark (M) queries

We detail now our experimental settings for the two engine we considered: Saxon-b/XQuery and MonetDB/XQuery.

*9.3.1. Test machine.* The experiments were performed on a desktop PC, with an Intel Core 2 Xeon 3Ghz CPU, 3.5 GB of RAM and a S-ATA hard-drive. We used Ubuntu Linux 9.10 64bits (featuring a 2.6.31 kernel) as operating system. The file-system used was ext3, with default settings. The OS was allocated 6 GB of swap space and tests were done in a reduced environment, where only essential services were running concurrently with our experiments. In what follows, when timings are reported they are obtained by removing the best and worst timing of 5 runs and averaging the remaining three. Also, all the parameters we measured (running time, memory consumption, I/O operations, ...) were measured in independent runs in order to have as little impact as possible on the experiments. The memory consumption of a running program was measured by monitoring the so-called "resident set size" field of the `/proc/pid/statm` pseudo file (this fields summarizes the amount of private data mapped in physical memory by the process, excluding shared segments such as shared libraries or shared mmaped files). I/O operations were monitored using the `iostat` utility.

*9.3.2. Saxon-b/XQuery.* Saxon [Kay 2008] is a popular XML library which implements various W3C standards (XPath, XQuery and XSLT) and has full schema support. We used version 9.0 of the Saxon-b XQuery engine (which is the Open Source one). Saxon being a main-memory query engine we focused on the following measurement both for pruned and unpruned documents: query answering time (excluding the parsing time of the document and serialization of the results, as reported by Saxon's debugging flags) and memory consumption. Saxon being written in Java, we

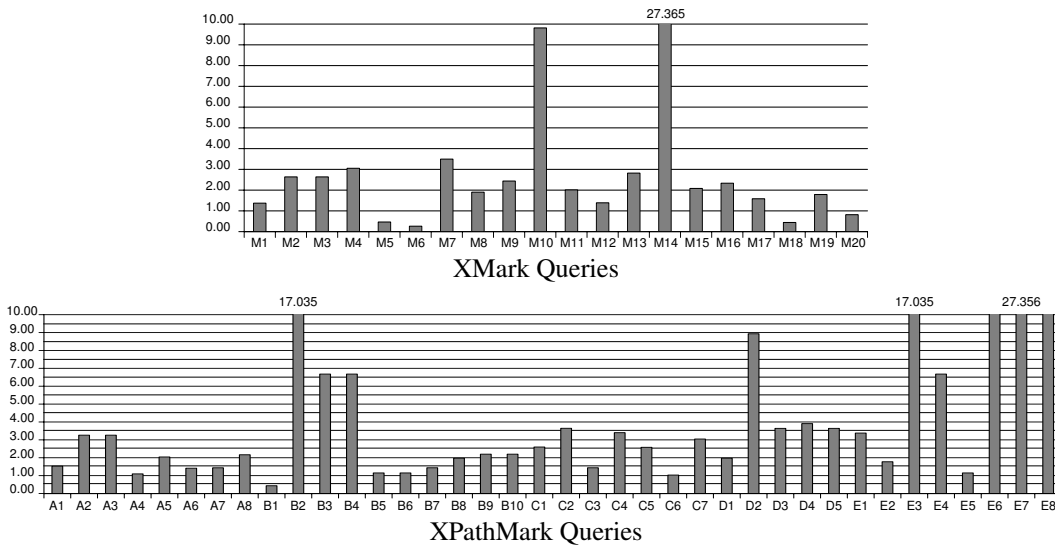


Fig. 7. Pruning ratio (% of the original file size) for XMark and XPathMark queries

used the latest version of the Sun's JVM available (1.6.0, 64 bit version) and set the amount of memory available to the JVM to the total physical memory of the machine.

**9.3.3. MonetDB/XQuery.** MonetDB/XQuery [Boncz et al. 2006] is a well established native XML database with full XQuery support. Contrary to Saxon, MonetDB stores *on disk* an index allowing fast navigation and query answering. In particular since it uses the disc as secondary storage, MonetDB is not limited by the amount of physical memory (it uses as much memory as possible to answer a query efficiently and performs its own page management by mapping memory pages to the disk and reading them back when needed). Therefore for such a query engine, speed is directly proportional to memory: the more memory is available, the less swapping occurs between pages on disk and pages in main memory. The three parameters we measured for MonetDB were the query answering time (again we did not consider document parsing nor serialization time), the size of the generated index on disk for a given document and the amount of I/O performed to answer the query. Indeed since MonetDB tries to max out its memory use to favour query answering time, measuring memory consumption does not reflect the actual scalability improvement one could expect when pruning documents. Disk access on the contrary are the bottleneck for such an engine and their frequency directly impacts query answering time.

## 9.4. Experimental results

**9.4.1. Pruning precision.** We gauged the precision of our pruning algorithm for the full set of XPathMark and XMark queries by comparing the size of the pruned document (serialized on disk) with the size of the corresponding original document. We report in Figure 7 the pruning ratio in percent of the original file size for XPathMark (labelled A1 to E8) and XMark (labelled M1 to M20) queries.

**9.4.2. Saxon-b/XQuery.** In our testing environment the biggest un-pruned document that the Saxon engine could handle was 671 MB large. We report in Figure 8 the original size of the largest *pruned* document Saxon could handle and the size of its projection (both in MB). Also, for a document of size 671 MB, we report the running time and memory consumption for the original and pruned version (as well as the size of the pruning). Lastly, we report the speed-up factor obtained thanks to pruning and the memory improvement we achieved (in percent of the original memory consumption) for the projected document. Due to the lack of space, we do not detail all of the XPathMark

	A 1	A 6	B 1	B 2	C 3	C 4	D 1	D 2	E 5	E 7	M 3	M 6	M 7	M 14	M 15
(i)	1.9	4.3	12.0	8.2	43.7	14.3	3.4	2.3	33.5	30.1	21.8	13.3	9.8	6.5	6.2
(ii)	3363	3363	3363	3363	3363	3363	3363	3363	3363	2242	3363	3363	3363	2242	3363
(iii)	447	67.6	50.5	571	68.3	137	65.5	297.25	60.8	605	92.8	9	121	605	67.6
(iv)	10	9	3	113	9	23	13	59	12	190	18	2	24	190	65
(v)	20	17	22	2.9	17	5.8	12.1	5.6	12.3	2.1	9	15.8	14.3	2.19	7.5
(vi)	3.7	5.5	2.2	22.2	3.7	8	5.4	9.7	7.44	22	9	1.8	4.9	29	15.2

(i): Computation of the type projector (ms)

(ii): Largest queryable document (MB).  
We stopped our testing at 3363 MB

(iii): Pruned size (MB)

(iv): Pruned size for 671 MB (MB)

(v): Speed up ( $\times$  faster)

(vi): Memory use in % of original

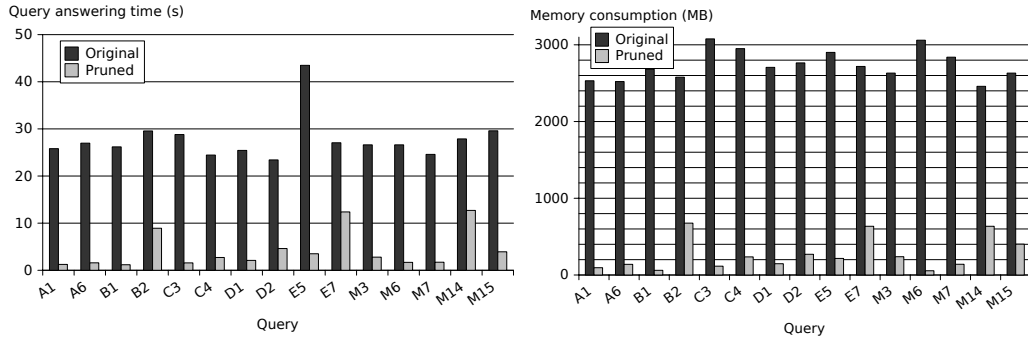


Fig. 8. Experimental results for the Saxon-b/XQuery engine

and XMark queries but rather for each category we give the “best performing query” –that is, the one for which we could achieve biggest speed-up– and the “worst performing query”, the one for which the speed-up was the smallest.

**9.4.3. MonetDB/XQuery.** Since MonetDB makes use of the secondary storage (disk) to query arbitrarily large documents, we chose a different approach to validate our pruning algorithm. We fixed the size of the input document to 3363 MB and then indexed it into the MonetDB document repository, yielding an index (on disk) of 4644 MB (as reported by the MonetDB administrative interface). Then for each query, we pruned the 3363 MB document with respect to the input query and indexed it. We summarize the results in Figure 9 (note that we did not report the time spent computing the type projector, since it is the same as in Figure 8). The first line in the table reports the size in MB of the index corresponding to the pruned document. The second line reflects the ratio between the amount of I/O operation performed by the MonetDB server for the pruned file and the amount of I/O performed on the original file. We only take into account the amount of data *read* from disk which helps us gauge the amount of data fetched from the index on disk into main-memory. In this same figure, the graphics represent the absolute query answering time in seconds for the original and pruned document.

Finally, the third line gives the speed-up in query answering achieved through pruning. We were not able to run the query E5 on our version of MonetDB (the server segfaulted at some point during the query computation).

## 9.5. Interpretation

**9.5.1. Type-projector inference.** As already stated, our type-projector inference algorithm is very efficient for practical queries. For all the queries we tested and the XMark DTD, computing the type projector took at most 49.4ms and about 15ms on average. We see that in practice, computing a type projector adds little overhead (despite the theoretical complexity). The analysis can therefore

	A1	A6	B1	B2	C3	C4	D1	D2	E5	E7	M3	M6	M7	M14	M15
(i)	675	119	95	843	117	232	93	443	116	1073	138	14	111	1073	436
(ii)	0.5	0.5	0.3	24.8	0.5	45.4	10.0	100.0	-	59	1.1	84.7	50.8	54.0	94.0
(iii)	6.8	10.5	8.5	3.3	18.1	2.1	21.6	1.0	-	1.7	5.0	2.8	1.7	2.06	1.0

(i): Index size (MB)

(ii): Amount of I/O (% of the original)

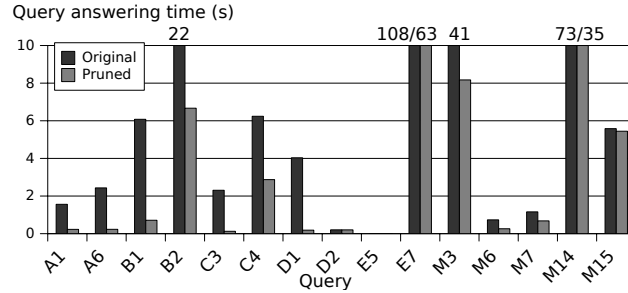
(iii): Speed up ( $\times$  faster)

Fig. 9. Experimental results for the MonetDB engine

be performed also at run-time. We include the detail of the timing of the type-projector computation for the selected set of queries that we study in the following subsections.

**9.5.2. Pruning precision.** The results from Figure 7 shows that, for the vast majority of the queries we considered, the document can be pruned to less than 10% of its original size. More precisely on the 58 queries we considered (20 XMark and 38 XPathMark queries):

- 47 queries yielded a projected document whose size was less than 5% of the original
- 5 queries (M10, B3, B4, D2, E4) had a pruning ratio between 5% and 10%
- 2 queries (B2, E3) had a pruning ratio of 17.035%
- 4 queries (M14, E6, E7, E8) had a pruning ratio of 27.35%

It should be noted that queries such as M14 return the content of a `description` element, consisting of almost all the textual data contained in the original XMark document. Since in these queries the value of the whole element is needed at runtime to perform string searching operations, there is little that can be done from the point of view of static pruning.

**9.5.3. Saxon-b/XQuery.** As we can see from the results in Figure 8 pruning the document before querying it always yields a speed-up and a reduction in memory use for a main memory engine such as Saxon. Furthermore, query answering time always dominates projector-inference time (several seconds for the former and less than 50ms for the latter).

On the one hand, for queries whose main bottleneck are string operations (such as calls to the `contains` function in M14 and E7), document projection gives very little speed-up. On the other hand, query A1 or C3 see a dramatic speed-up (20 and 17 times faster than the original respectively). This shows that despite the various optimizations built into the engine, a significant amount of time is often spent by iterating over “non-relevant” nodes which are discarded by the pruning process.

On the side of memory consumption, pruning the document unsurprisingly reduces memory usage drastically. Indeed, document projection reduces the number of elements and therefore simplifies the tree-structure of the XML document. This aspect is critical for main-memory engines which often (as in the case of Saxon) represent the document as a pointer-based data-structure (e.g., following the DOM model [DOM 2004] where each element is represented as a node which contains a pointer to its first child to the next sibling, and to its parent). Indeed, we experienced for Saxon (but we observed similar behaviour in other main-memory query engines) a 112MB XMark document would occupy 430MB of RAM while the same document stripped of its data—amounting to only 36MB on disk—would occupy 340MB of memory. As Figure 8 illustrates, our pruning technique

precisely addresses this issue, reducing in most cases the memory consumption to a few percents of what is needed to handle an un-pruned file.

*9.5.4. MonetDB/XQuery.* MonetDB is known to be one of the fastest XML databases available. The efficiency of the MonetDB/XQuery engine is essentially due to the stair-case join operation ([Grust et al. 2003]) which minimizes the amount of intermediate sets constructed to answer an XPath query. Even so, the use of type-based document projection often improves query answering time. In particular, as shown in Figure 9, a smaller index often yields less I/O operations which in turn increases the speed of the query engine. On the contrary for queries such as D2, M14 and M15, the document is already optimally indexed and reducing the size of the index does not reduce the amount of I/O which explains why for these queries the gain in speed is null. Yet for some queries the speed-up can be up to twenty-folds (D1).

*9.5.5. Comparison with related work.* These results are a clear-cut improvement over current technology. While we cannot directly compare processing performances since no implementation of the other pruning approaches is publicly available, we want to stress two points. First, for XMark queries the pruning precision we achieve is equal or better than what is obtained with other approaches (with the exception of query M10 for which [Marian and Siméon 2003] achieves a pruning ratio of 4.5% where we could only prune the document down to 9.2% of its original size). Second, performing pruning never is a bottleneck in our case thanks to the fact that our solution consists of a single buffer-less traversal of the input document (on our test machine we were able to efficiently prune arbitrary large documents, while in case of [Marian and Siméon 2003] pruning can end up using as much memory as the execution of the query).

The experiments also illustrate that our approach retains a very high precision even in the presence of complex XPath features (like backward axes and external functions). While it is true that the technique of [Olteanu et al. 2002] could be used to allow Marian and Simeon's work to handle backward axes, it would still not be, to our sense, a satisfactory solution. The first reason is that the rewrite rules given in [Olteanu et al. 2002] do not support the use of data-value or negation in the filters of the original query (see [Benedikt and Koch 2008]). For instance the query

`descendant::keyword [not(ancestor::item)]`

cannot be written without backward axis. Second the query generated may be exponentially bigger than the original one (and its computation takes exponential time in the size of the original) and may introduce several predicates as well as `descendant-or-self` axes. Both features degrades the pruning precision of [Marian and Siméon 2003].

## 10. CONCLUSION AND FUTURE WORK

Our experiments show the clear advantages of applying our optimisation technique to query XML documents, and the characteristics of our solution make it profitable in all application scenarios. We discussed several aspects for which our approach improves the state of the art: for performances (better pruning, greater speedup, smaller memory footprint), for the analysis techniques (linear pruning time, negligible memory and time consumption), for its generality (handling of full XPath), and, last but not least, for the formal foundation it provides (correctness formally proved, limits of the approach formally stated).

The present work extends and improves the current state of the art in several aspects. From a formal point of view, the use of regular tree grammars as schema model makes the technique applicable to the various kind of schemas currently in use. Furthermore, the closure properties that we proved ensure that type-based projection is at most as expensive as validation for a given class of schema language. We also handle a richer set of queries formally (in particular we handle nested predicates in XPath<sup>ℓ</sup>) and took special care to document how to encode or approximate several important XPath idioms that were lacking from the formal presentation. On a practical level, we have validated our approach against state of the art query engines, using realistic queries and data sets. In particular, not only did we test against an efficient main memory query engine (Saxon) but also

demonstrated that our approach can be used to improve, sometimes by a double digit factor, the performances of an already very optimized disk-based XML database such as MonetDB/XQuery.

Future work will be pursued both at a formal and practical level. At a formal level, one of the main shortcomings of our approach is its reliance on XPath syntax. Indeed, even though we managed to isolate a fragment of XPath that we could formally reason with, it still leaves us with a syntax-directed approach. The problem with this is twofold. First, it makes the proofs and the specification of algorithms quite tedious and unnecessarily intricate. Second and more importantly, our pruning inference algorithm might yield different type projectors depending on the syntax of the original query. For future work, we would like to tackle a semantic based approach. In particular it seems worthwhile to consider more theoretically sound formalisms for tree queries such as, for instance, MSO formula or tree automata. The latter in particular would allow us to reuse our pruning algorithm for pattern-matching based languages (such as the CDuce language [CDuce 2011] and its query language CQL [Miachon 2006; Benzaken et al. 2005; Castagna 2005]). It is also known that tree-automata have better closure properties than XPath expressions and support fine-grained set-theoretic operations (intersection, union, complement) that have been used with success to devise very precise type-systems for XML [Frisch et al. 2008].

At a practical level we would like to see a tighter integration between document-projection and query engines. Firstly, although quite crude, our experiments show that even a carefully designed indexed system such as MonetDB can benefit from document pruning. It seems interesting to develop further such preliminary results and design a projection aware XML index. In other words we would like to be able to equip any native XML query engine optimizer with a type-projector component. In particular, one could think of an index consisting of the original document together with its projected versions. Textual data could be shared between the main document and the projected ones which would merely become a projected view of the tree structure of the document. We make the hypothesis that the overhead of such pruned tree structures would be quite small compared to the size of an XML index while providing significant speed-up in query answering time.

Secondly, coupling our projection algorithm with early query answering techniques would allow us to achieve further pruning, especially when runtime conditions are involved. For instance we could use our type-inference algorithm to determine on what type of elements a given built-in function is applied to; for instance in an expression such as

```
contains(./.*, "foo").
```

This information could then be used to discard elements that do not match the predicate.

## ELECTRONIC APPENDIX

The electronic appendix for this article can be accessed in the ACM Digital Library.

## ACKNOWLEDGMENTS

???

## References

- CDuce 2004-2011. The CDuce language. <http://www.cduce.org>.
- ARROYUELO, D., CÁNOVAS, R., NAVARRO, G., AND SADAKANE, K. 2010. Succinct trees in practice. In *Proc. 11th Workshop on Algorithm Engineering and Experiments (ALENEX)*. SIAM Press. To appear.
- BENEDIKT, M., FAN, W., AND GEERTS, F. 2008. Xpath satisfiability in the presence of dtds. *J. ACM* 55, 2.
- BENEDIKT, M. AND KOCH, C. 2008. XPath leashed. *ACM Comput. Surv.* 41, 1, 1–54.
- BENZAKEN, V. 1990. An Evaluation Model for Clustering Strategies in the O<sub>2</sub> Object-Oriented Database System. In *ICDT*. 126–140.
- BENZAKEN, V., CASTAGNA, G., COLAZZO, D., AND NGUYỄN, K. 2006. Type-based XML projection. In *VLDB '06*. 271–282.
- BENZAKEN, V., CASTAGNA, G., AND MIACHON, C. 2005. A full pattern-based paradigm for XML query processing. In *PADL '05*. Number 3350 in LNCS. 235–252.

- BENZAKEN, V. AND DELOBEL, C. 1990. Enhancing Performance in a Persistent Object Store: Clustering Strategies in O<sub>2</sub>. In *Proceedings of Persistent Object Systems 4*. Morgan Kaufmann.
- BENZAKEN, V., DELOBEL, C., AND HARRUS, G. 1992. Clustering strategies in o<sub>2</sub>: an overview. In *Building an Object-Oriented Database System: the Story of O<sub>2</sub>*, F. Bancilhon, C. Delobel, and P. Kanellakis, Eds. Morgan Kaufman, 385–410.
- BEX, G. J., GELADE, W., NEVEN, F., AND VANSUMMEREN, S. 2008. Learning deterministic regular expressions for the inference of schemas from XML data. In *WWW*. 825–834.
- BEX, G. J., NEVEN, F., SCHWENTICK, T., AND TUYLS, K. 2006. Inference of concise dtds from XML data. In *VLDB*. 115–126.
- BEX, G. J., NEVEN, F., AND VANSUMMEREN, S. 2007. Inferring XML schema definitions from XML data. In *VLDB*. 998–1009.
- BONCZ, P. A., GRUST, T., VAN KEULEN, M., MANEGOLD, S., RITTINGER, J., AND TEUBNER, J. 2006. MonetDB/XQuery: a fast XQuery processor powered by a relational engine. In *SIGMOD*. 479–490.
- BRESSAN, S., CATANIA, B., LACROIX, Z., LI, Y.-G., AND MADDALENA, A. 2005. Accelerating queries by pruning XML documents. *Data Knowl. Eng.* 54, 2, 211–240.
- CASTAGNA, G. 2005. Patterns and types for querying XML. In *DBPL '05*. Number 3774 in LNCS.
- CASTAGNA, G. AND FRISCH, A. 2005. A gentle introduction to semantic subtyping. In Proc. of *PPDP '05* (full version) and *ICALP '05*, LNCS n. 3580, (summary). Joint ICALP-PPDP keynote talk.
- CHAMBERLIN, D., FANKHAUSER, P., FLORESCU, D., MARCHIORI, M., AND ROBIE, J. 2003. XML Query Use Cases. Tech. Rep. 20030822, World Wide Web Consortium.
- CHOI, B. 2002. What are real dtds like? In *WebDB*. 43–48.
- COMON, H., DAUCHET, M., GILLERON, R., JACQUEMARD, F., LUGIEZ, D., TISON, S., AND TOMMASI, M. 1997. Tree automata techniques and applications. <http://www.grappa.univ-lille3.fr/tata>.
- DOM 2004. W3C: DOM specifications. <http://www.w3.org/TR/DOMTR>.
- DRAPER, D., FANKHAUSER, P., FERNANDEZ, M., MALHOTRA, A., ROSE, K., RYS, M., SIMÉON, J., AND WADLER, P. 2007. XQuery 1.0 and XPath 2.0 Formal Semantics. Tech. rep., World Wide Web Consortium.
- FRANCESCHET, M. 2005. XPathMark - An XPath benchmark for XMark generated data. In *XSym 2005, 3rd Int. XML Database Symposium*. LNCS n. 3671.
- FRISCH, A., CASTAGNA, G., AND BENZAKEN, V. 2008. Semantic subtyping: dealing set-theoretically with function, union, intersection, and negation types. *Journal of the ACM* 55, 4, 1–64.
- GREEN, T. J., GUPTA, A., MIKLAU, G., ONIZUKA, M., AND SUCIU, D. 2004. Processing XML streams with deterministic automata and stream indexes. *ACM Trans. Database Syst.* 29, 4, 752–788.
- GRUST, T., VAN KEULEN, M., AND TEUBNER, J. 2003. Staircase join: Teach a relational DBMS to watch its (axis) steps. In *VLDB '03*. 524–535.
- KAY, M. 2008. Ten reasons why Saxon XQuery is fast. *IEEE Data Eng. Bull.* 31, 4, 65–74.
- KOCH, C. 2003. Efficient processing of expressive node-selecting queries on XML data in secondary storage: a tree automata-based approach. In *VLDB '03*.
- KOCH, C. 2005. On the complexity of nonrecursive xquery and functional query languages on complex values. In *ACM Transactions on Database Systems*. 2006.
- MARIAN, A. AND SIMÉON, J. 2003. Projecting XML documents. In *VLDB '03*. 213–224.
- MARTENS, W., NEVEN, F., SCHWENTICK, T., AND CENTRUM, L. U. 2005. Which xml schemas admit 1-pass preorder typing. In *In ICDT*. Springer, 68–82.
- MIACHON, C. 2006. Langages de requêtes pour XML à base de patterns : conception, optimisation et implantation. Ph.D. thesis, Université Paris-Sud 11.
- MURATA, M., LEE, D., MANI, M., AND KAWAGUCHI, K. 2005. Taxonomy of XML schema languages using formal language theory. *ACM Trans. Internet Techn.* 5, 4, 660–704.
- NGUYỄN, K. 2008. Combinator language for XML: design, typing, and implementation. Ph.D. thesis, Université Paris-Sud 11.
- OLTEANU, D., MEUSS, H., FURCHE, T., AND BRY, F. 2002. XPath: Looking forward. In *Proc. EDBT Workshop (XMLDM)*. LNCS Series, vol. 2490. Springer, 109–127.
- SCHMIDT, A., WAAS, F., KERSTEN, M. L., CAREY, M. J., MANOLESCU, I., AND BUSSE, R. 2002. XMark: A benchmark for XML data management. In *VLDB '02*. 974–985.
- SEGOUFIN, L. AND VIANU, V. 2002. Validating streaming XML documents. In *PODS '02*. 53–64.
- XML 2008. W3C: XML Version 1.0 (Fifth Edition). <http://www.w3.org/TR/REC-xml/>.
- XPath 1999. W3C: XML Path Language (XPath) Version 1.0. <http://www.w3.org/TR/xpath>.



XPATHMARK. XPathMark web site (M. Franceschet). <http://sole.dimi.uniud.it/~massimo.franceschet/xpathmark/>.

XQuery 2010. W3C: XML Query (XQuery). <http://www.w3.org/TR/xquery>.

Received June 2011; revised July 2012; accepted October 2012

## Online Appendix to: Optimizing XML querying using type-based document projection

VÉRONIQUE BENZAKEN, GIUSEPPE CASTAGNA, DARIO COLAZZO, and KIM NGUYỄN

---

### A. DETAILED PROOFS

#### A.1. Type projectors

LEMMA ((3.3)) *Let  $\pi$  be a type projector for  $(\mathcal{S}, E)$ . Then for every tree  $t \in_{\mathcal{T}}(\mathcal{S}, E)$  it holds  $(t \setminus_{\mathcal{T}} \pi) \preceq t$ .*

PROOF. The proof is a straightforward induction on  $t$ .  $\square$

PROPOSITION (ERASURE PRESERVES LOCALITY (3.4)) *Let  $(\mathcal{S}, E)$  be a local tree grammar and  $(\mathcal{S}', E')$  a regular tree grammar. If  $(\mathcal{S}', E') <: (\mathcal{S}, E)$  then  $(\mathcal{S}', E')$  is a local tree grammar.*

PROOF. By contradiction, suppose that  $(\mathcal{S}', E')$  is not a local tree grammar. By Definition 2.12,  $\mathcal{S}' \subseteq \mathcal{S}$  therefore,  $|\mathcal{S}'| \leq |\mathcal{S}| \leq 1$ .

- Then, either there exist two competing rules  $A \rightarrow l[r'_a]$  and  $B \rightarrow l[r'_b]$  in  $E'$ . Then by definition of erasure, there exist two rules  $A \rightarrow l[r_a]$  and  $B \rightarrow l[r_b]$  in  $E$  such that  $r'_a = r_a|_{N_a}$  and  $r'_b = r_b|_{N_b}$  for some  $N_a \subseteq \mathbf{Names}(r_a)$  and  $N_b \subseteq \mathbf{Names}(r_b)$ . But then these two rules share the same label  $l$ , and therefore are competing one with the other, which contradicts the fact that  $(\mathcal{S}, E)$  is a local tree grammar.
- Or there exists two rules  $C \rightarrow l[r_1]$  and  $C \rightarrow l'[r_2]$  in  $E'$  with the same left hand-side (and distinct labels). But then, by definition of erasure, there exists two corresponding rules in  $E$ ,  $C \rightarrow l[r'_1]$  and  $C \rightarrow l'[r'_2]$  such that  $r_i = r'_i|_{N_i}, i \in \{1, 2\}$  for some names  $N_i$ . Therefore there are two rules in  $E$  with the same left-hand side, which contradicts the fact that  $(\mathcal{S}, E)$  is a local-tree grammar.

$\square$

PROPOSITION (ERASURE PRESERVES SINGLE-TYPEDNESS (3.5)) *Let  $(\mathcal{S}, E)$  be a single-type tree grammar and  $(\mathcal{S}', E')$  a regular tree grammar. If  $(\mathcal{S}', E') <: (\mathcal{S}, E)$  then  $(\mathcal{S}', E')$  is a single-type tree grammar.*

PROOF. By contradiction suppose that  $(\mathcal{S}', E')$  is not a single-type tree grammar and proceed by case analysis:

---

© YYYY ACM 0362-5915/YYYY/01-ARTA \$15.00  
DOI 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

Véronique Benzaken, Giuseppe Castagna, Dario Colazzo, and Kim Nguyễn. Optimizing XML querying using type-based document projection.  
*ACM Transactions on Database Systems, TO APPEAR.*

- either there exist two competing non terminals  $A$  and  $B$  in  $\mathcal{S}'$ . But by definition of erasure,  $\mathcal{S}' \subseteq \mathcal{S}$  and  $(\mathcal{S}, E)$  has two competing start symbols, which contradicts the hypothesis that  $(\mathcal{S}, E)$  enjoys the single-type property.
- or there exists a rule  $X \rightarrow l[r']$  and there exist two competing non terminals  $A$  and  $B$  in  $\mathbf{Names}(r')$ . Since  $(\mathcal{S}', E') <: (\mathcal{S}, E)$ , then there exists  $X \rightarrow l[r]$  such that  $r' = r|_N$  for some  $N \subseteq \mathbf{Names}(r)$ . But that means that  $A$  and  $B$  are in  $\mathbf{Names}(r)$ , which implies that  $(\mathcal{S}, E)$  is not a single-type tree grammar, thus contradicting our hypothesis.  $\square$

**PROPOSITION (UNION CLOSURE OF LOCAL TYPE PROJECTORS (3.6))** *Let  $(\mathcal{S}, E)$  be a local tree grammar. Let  $(\mathcal{S}_1, E_1)$  and  $(\mathcal{S}_2, E_2)$  be two tree grammars such that  $(\mathcal{S}_1, E_1) <: (\mathcal{S}, E)$  and  $(\mathcal{S}_2, E_2) <: (\mathcal{S}, E)$ . Then  $(\mathcal{S}_1 \cup \mathcal{S}_2, E_1 \cup E_2)$  is a local tree grammar.*

**PROOF.** Consider  $(\mathcal{S}_1 \cup \mathcal{S}_2, E_1 \cup E_2)$  and suppose, by contradiction, that it is not local. First, remark that by definition of erasure,  $\mathcal{S}_1 \subseteq \mathcal{S}$  and  $\mathcal{S}_2 \subseteq \mathcal{S}$ , therefore  $\mathcal{S}_1 \cup \mathcal{S}_2 \subseteq \mathcal{S}$  and consequently  $|\mathcal{S}_1 \cup \mathcal{S}_2| \leq |\mathcal{S}| \leq 1$ . Second:

- either we have two rules  $A \rightarrow l[r_a]$  and  $B \rightarrow l[r_b]$  (with  $A$  and  $B$  distinct). By Lemma 3.4 we know that  $(\mathcal{S}_1, E_1)$  and  $(\mathcal{S}_2, E_2)$  are local tree grammars. Then it must be that one of two rules at issue is in  $(\mathcal{S}_1, E_1)$  and the other in  $(\mathcal{S}_2, E_2)$ , otherwise one of the two grammars would not be local (it would have the competing pair in its rules). Without loss of generality we can suppose that  $A \rightarrow l[r_a] \in E_1$  and  $B \rightarrow l[r_b] \in E_2$ . Since  $(\mathcal{S}_1, E_1) <: (\mathcal{S}, E)$ , then by definition of erasure there exists a rule  $A \rightarrow l[r'_a] \in E$  with  $r_a = r'_a|_N$  for some  $N \subseteq \mathbf{Names}(r'_a)$ . Similarly, there exists  $B \rightarrow l[r'_b] \in E$  with  $r_b = r'_b|_{N'}$  for some  $N' \subseteq \mathbf{Names}(r'_b)$ . But then, this means that we have two competing rules in  $E$ , which contradicts the hypothesis that  $(\mathcal{S}, E)$  is a local tree grammar.
- or, there are two rules  $C \rightarrow l[r_1]$  and  $C \rightarrow l'[r_2]$  with the same left-hand side and distinct labels in  $E_1 \cup E_2$ . Similarly to the previous case, we must have  $C \rightarrow l[r_1] \in E_1$  and  $C \rightarrow l'[r_2] \in E_2$ , otherwise  $(\mathcal{S}_1, E_1)$  or  $(\mathcal{S}_2, E_2)$  would not be local. But then by definition of erasure, it means that there exists two rules,  $C \rightarrow l[r'_1] \in E$  and  $C \rightarrow l'[r'_2] \in E$  such that  $r_i = r'_i|_{N_i}, i \in \{1, 2\}$  for some names  $N_i$ . This means that there are two rules in  $E$  with distinct labels and the same left-hand side, which contradicts the assumption that  $(\mathcal{S}, E)$  is a local tree grammar.

$\square$

**PROPOSITION (UNION CLOSURE OF SINGLE-TYPE TYPE PROJECTORS (3.7))**

*Let  $(\mathcal{S}, E)$  be a single-type tree grammar. Let  $(\mathcal{S}_1, E_1)$  and  $(\mathcal{S}_2, E_2)$  be two tree grammars such that  $(\mathcal{S}_1, E_1) <: (\mathcal{S}, E)$  and  $(\mathcal{S}_2, E_2) <: (\mathcal{S}, E)$ . Then  $(\mathcal{S}_1 \cup \mathcal{S}_2, E_1 \cup E_2)$  is a single-type tree grammar.*

**PROOF.** Consider  $(\mathcal{S}_1 \cup \mathcal{S}_2, E_1 \cup E_2)$  and suppose by contradiction that it does not enjoy the single-type property. This implies that there exists a rule  $X \rightarrow l[r]$ , such that  $\mathbf{Names}(r)$  contains two competing non-terminals  $A$  and  $B$ . Moreover, by Lemma 3.5 we know that  $(\mathcal{S}_1, E_1)$  and  $(\mathcal{S}_2, E_2)$  are single-type tree grammars. Therefore  $X \rightarrow l[r]$  cannot be in  $E_1$  nor in  $E_2$  because they have the single-type property. The only solution is that there exists a rule  $X \rightarrow l[r_1]$  in  $E_1$  with  $A \in \mathbf{Names}(r_1)$  and  $X \rightarrow l[r_2]$  in  $E_2$  with  $B \in \mathbf{Names}(r_2)$ , and that  $r = r_1|r_2$  (remember that we identify two rules with the same left-hand side and the same label by merging them into a single rule). Therefore, by the definition of erasure, there exists a rule  $X \rightarrow l[r'_1]$  in  $E$

such that  $A \in \text{Names}(r'_1)$ . Similarly, there exists a rule  $X \rightarrow l[r'_2]$  in  $E$  such that  $B \in \text{Names}(r'_2)$ . Since we identify such rules, there is a rule  $X \rightarrow l[r'_1|r'_2]$  in  $E$ . But then, this rule contains both  $A$  and  $B$  which are competing. This contradicts the hypothesis that  $(\mathcal{S}, E)$  is a single-type tree grammar.  $\square$

## A.2. Static analysis

LEMMA (5.2) *Let  $t$  be a tree  $\mathcal{J}$ -valid with respect to the schema  $(\mathcal{S}, E)$ . For every  $S \subseteq \text{Ids}(t)$  and type  $\tau$ , if  $\mathcal{J}(S) \subseteq \tau$ , then*

- (1)  $\mathcal{J}(\llbracket \text{Axis} \rrbracket_t(S)) \subseteq \mathbf{A}_E(\tau, \text{Axis})$
- (2)  $\mathcal{J}(S : :_t \text{Test}) \subseteq \mathbf{T}_E(\tau, \text{Test})$

PROOF. The proof is done by case analysis on the possible axes (for (1.)) and tests (for (2.)).

- (1) **self.** By Definition 4.3, we have that  $\llbracket \text{self} \rrbracket_t(S) = S$ , therefore  $\mathcal{J}(\llbracket \text{self} \rrbracket_t(S)) = \mathcal{J}(S)$ . By Definition 5.1,  $\mathbf{A}_E(\tau, \text{self}) = \tau$ . Since  $\mathcal{J}(S) \subseteq \tau$  by hypothesis, we can conclude that

$$\mathcal{J}(\llbracket \text{self} \rrbracket_t(S)) \subseteq \mathbf{A}_E(\tau, \text{self})$$

**child.** we suppose  $\mathbf{i}' \in \llbracket \text{child} \rrbracket_t(S)$ . Let us show that  $\mathcal{J}(\mathbf{i}') \in \mathbf{A}_E(\tau, \text{child})$ . If  $\mathbf{i}' \in \llbracket \text{child} \rrbracket_t(S)$  then by Definition 4.3:

$$\exists \mathbf{i} \in \text{Ids}(t) \text{ such that } (\mathbf{i}, \mathbf{i}') \in \text{Edg}(t)$$

By definition of **Edg**, this means that  $t@i = l[\dots t' \dots]$  and  $\text{RootId}(t') = \mathbf{i}'$ . Let us now call  $X \rightarrow l[r] = \mathcal{J}(\mathbf{i})$ . Since  $t$  is  $\mathcal{J}$ -valid with respect to  $E$  and by Definition 2.9, we have that  $\text{Dn}(\mathcal{J}(\mathbf{i}')) \in \mathcal{L}(r)$ , or equivalently that  $\mathcal{J}(\mathbf{i}') = Y \rightarrow R'$  and  $Y \in \text{Names}(r)$ , which implies that  $\mathcal{J}(\mathbf{i}') \in \mathbf{A}_E(\tau, \text{child})$  and therefore that

$$\mathcal{J}(\llbracket \text{child} \rrbracket_t(S)) \subseteq \mathbf{A}_E(\tau, \text{child})$$

**descendant.** we suppose  $\mathbf{i}' \in \llbracket \text{descendant} \rrbracket_t(S)$ . Let us show that  $\mathcal{J}(\mathbf{i}') \in \mathbf{A}_E(\tau, \text{descendant})$ . If  $\mathbf{i}' \in \llbracket \text{descendant} \rrbracket_t(S)$  then by Definition 4.3:

$$\exists \mathbf{i} \in \text{Ids}(t) \text{ such that } (\mathbf{i}, \mathbf{i}') \in \text{Edg}(t)^+$$

By definition of **Edg** this means that there exists a sequence  $\mathbf{i}_0, \mathbf{i}_1, \dots, \mathbf{i}_n$ . The property holds by induction on  $n$  with the base case  $n = 1$  is the one of **child**.

**parent.** dual of **child**

**ancestor.** dual of **descendant**

- (2) By case on the test:

**node.** By Definition 4.2 we have:

$$\text{node} : :_t S = S$$

By Definition 5.1 we have  $\mathbf{T}_E(\tau, \text{node}) = \tau$ . Since  $\mathcal{J}(S) \subseteq \tau$  by hypothesis, we have that:

$$\mathcal{J}(S : :_t \text{node}) \subseteq \mathbf{T}_E(\tau, \text{node})$$

*a (for some element name a).* suppose  $\mathbf{i} \in a : :_t S$ . Let us show that  $\mathcal{J}(\mathbf{i}) \in \mathbf{T}_E(\tau, a)$ . By Definition 4.2, we know that  $t@i = a[f]$  for some forest  $f$ . Since  $t$  is  $\mathcal{J}$ -valid with respect to  $E$ , then  $\mathcal{J}(\mathbf{i}) = Y$  and there exists a rule  $Y \rightarrow a[R]$  in  $E$ . Since  $\mathbf{i} \in S$ , we also have that  $\mathcal{J}(\mathbf{i}) \in \tau$  (by hypothesis). By Definition 5.1, since  $Y \rightarrow a[R] \in \tau$ , then  $Y \rightarrow a[R] \in \mathbf{T}_E(\tau, a)$  and therefore  $\mathcal{J}(\mathbf{i}) = Y \in \mathbf{T}_E(\tau, a)$ , hence

$$\mathcal{J}(S : :_t a) \subseteq \mathbf{T}_E(\tau, a)$$

**text.** similar to the previous case.

$\square$

LEMMA (TERMINATION OF TYPE INFERENCE (5.4)) *Let  $(\mathcal{S}, E)$  be a type,  $P$  a path, and  $\Sigma$  and  $\Sigma'$  two environments. If there is a derivation for the judgment  $\Sigma \vdash_E P : \Sigma'$ , then this derivation is unique and finite. Furthermore, the algorithm run in time  $O(|E|^d \times |P|)$ , where  $d$  is the maximum number of nested predicates.*

PROOF. Uniqueness of the derivation is immediate, since the rules are syntax-directed: at each step, at most one of the rules applies (if no rule applies, there is no derivation and the output type is  $\emptyset$ ). Finiteness can be shown by a simple induction on the length of the path, noted  $l(P)$ , that we define as follows:

$$\begin{aligned} l(\text{Axis}::\text{Test}) &= 1 \\ l(\text{Axis}::\text{Test}[C]) &= 1 + \sum_{P \in C} l(P) \\ l(P/P') &= l(P) + l(P') \\ l(P|P') &= l(P) + l(P') \end{aligned}$$

**Basic case:** The query has length 1, meaning it is a single step without predicate. Then the only rules we can apply are **(down-axis)**, **(up-axis)** or **(test)**. These rules have no premise, therefore the derivation is finite and has length 1.

**Inductive case:** If the query has several steps, then the rule **(seq)** applies. The lengths of the queries in the premises of the rule are strictly less than the length of the query in the goal, by definition of  $l(\_)$ . By induction hypothesis both premises have a finite derivation, therefore the goal can be derived with a finite derivation.

Similarly, if the query is a top-level union, the typing rule **(union)** applies.

If the query is a single step with a predicate, then rule **(predicate)** applies. We should first remark that there is a finite set of rules in  $\Sigma_{\text{typ}}$ . Thus, there are exactly  $|\Sigma_{\text{typ}}|$  premises for this rule. Consider the condition  $C$  in the filter. It contains  $n$  paths  $P_i$ , for  $i$  in  $1..n$ . Each one of these path  $P_i$  is such that  $l(P_i) < l(P)$ , by definition of  $l(\_)$ . By induction hypothesis, every premise has finite derivation, therefore the judgment in the goal of the rule has a finite derivation.

As for the complexity result, consider at first the case where  $P$  does not contain any predicate. The algorithm performs  $l(P)$  recursive calls. At each call, it performs set-theoretic operations between subsets of the input type, which can be implemented in  $O(|E|)$ , thus yielding a time complexity of  $O(|E| \times |P|)$ . However, if the path contains a predicate, (that is if rule **(predicate)** is used) then type-checking is recursively called for subpaths whose nesting decreases by one, for each rule in  $\Sigma_{\text{typ}}$  (which is at most  $E$ ). This gives an overall complexity of  $O(|E|^d \times |P|)$ , where the depth of a path without predicate is 1.  $\square$

To prove soundness of type inference, we require an auxiliary definition:

*Definition A.1 (Environment well-formedness).* Let  $(\tau, \kappa)$  be an environment and  $E$  a set of rules. If  $\tau \subseteq E$  and  $\kappa = \tau \cup \mathbf{A}_E(\tau, \text{ancestor})$ , then we say that  $(\tau, \kappa)$  is *well formed* with respect to  $E$ .

In other words, a context is well-formed if it contains only rules from which the names in  $\mathbf{Dn}(\tau)$  are reachable. We say that a judgment  $\Sigma \vdash_E P : \Sigma'$  is well formed if both  $\Sigma$  and  $\Sigma'$  are well formed with respect to  $E$ . We can remark that the rules in Figure 1 are syntax directed—at most one rule apply for a given judgment—and they preserve context well-formedness.

**THEOREM (SOUNDNESS OF TYPE INFERENCE (5.5))** *Let  $(\mathcal{S}, E)$  be a type and  $P$  a path. Let  $E_0 = \{X \rightarrow R \mid X \rightarrow R \in E, X \in \mathcal{S}\}$ . If  $(E_0, E_0) \vdash_E P : (\tau, \kappa)$  then:*

$$\tau \supseteq \bigcup_{t \in \mathcal{J}(\mathcal{S}, E)} \mathcal{J}(\llbracket P \rrbracket_t(\mathbf{RootId}(t)))$$

**PROOF.** We consider the following, more general judgment:

$$(\tau, \kappa) \vdash_E P : (\tau', \kappa')$$

We show simultaneously the following properties:

- (1) **Soundness** : for all tree  $t$   $\mathcal{J}$ -valid with respect to  $(\mathcal{S}, E)$  and all set  $S \subseteq \mathbf{Ids}(t)$ , if  $\mathcal{J}(S) \subseteq \tau$  then:

$$\mathcal{J}(\llbracket P \rrbracket_t(S)) \subseteq \tau'$$

- (2) **Context well-formedness**, if

$$\kappa = \tau \cup \mathbf{A}_E(\tau, \mathbf{ancestor})$$

then

$$\kappa' = \tau' \cup \mathbf{A}_E(\tau', \mathbf{ancestor})$$

Property 1 is a generalization of the soundness property we are proving. Given a set of context nodes  $S$  whose types are in  $\tau$ , the type of  $\llbracket P \rrbracket_t(S)$  is in  $\tau'$ . Property 2 states that the algorithm preserves the well-formedness of contexts. We prove both properties by induction on the depth of the typing derivation, which is finite by Lemma 5.4:

**Base case:**

**(down-axis):** Property 1 is true by a direct application of Lemma 5.2. Property 2 holds by definition of  $\mathbf{A}_E(-, -)$

**(up-axis):** By Lemma 5.2:

$$\mathcal{J}(\llbracket \mathbf{Axis}::\mathbf{node} \rrbracket_t(S)) \subseteq \mathbf{A}_E(\tau, \mathbf{Axis})$$

We must now show that  $\mathcal{J}(\llbracket \mathbf{Axis}::\mathbf{node} \rrbracket_t(S))$  is in  $\Sigma_{\text{ctx}} = \kappa$ , for it to be in the intersection of both. Since  $\kappa$  is a well-formed context:

$$\kappa = \tau \cup \mathbf{A}_E(\tau, \mathbf{ancestor})$$

Let us first consider the case  $\mathbf{Axis} = \mathbf{ancestor}$ . We see that

$$\mathcal{J}(\llbracket \mathbf{ancestor}::\mathbf{node} \rrbracket_t(S)) \subseteq \mathbf{A}_E(\tau, \mathbf{ancestor})$$

by Lemma 5.2, therefore

$$\mathcal{J}(\llbracket \mathbf{ancestor}::\mathbf{node} \rrbracket_t(S)) \subseteq \kappa$$

from which we can conclude

$$\mathcal{J}(\llbracket \mathbf{ancestor}::\mathbf{node} \rrbracket_t(S)) \subseteq (\mathbf{A}_E(\tau, \mathbf{Axis})) \cap \kappa$$

For the case  $\mathbf{Axis} = \mathbf{parent}$ , it is sufficient to remark that since

$$\kappa = \tau \cup \mathbf{A}_E(\tau, \mathbf{ancestor})$$

then

$$\kappa = \tau \cup \mathbf{A}_E(\tau, \mathbf{parent}) \cup \mathbf{A}_E(\mathbf{A}_E(\tau, \mathbf{parent}), \mathbf{ancestor})$$

which proves Property 1. As for Property 2,  $\kappa$  is the set of rules used to derive the context node type  $\tau$ .  $\mathbf{A}_E(\kappa, \mathbf{Axis})$  is the set of all the parent rules (or ancestor rules) of the rules in  $\kappa$ . Consequently, the intersection is still a well formed context.

(*test*). : Similarly to the case of Rule (**down-axis**), Property 1 is a direct application of Lemma 5.2. For Property 2, we can remark that

$$\kappa' = \kappa \cap \mathbf{A}_E(T_E(\tau, Test), \text{ancestor})$$

contains all the rules leading to a node in  $\tau$  for which *Test* succeeds (including the ones of the selected node), therefore it is a well-formed context.

**Inductive case.**

(*predicate*). : Let us consider:

$$\llbracket \text{self} :: \text{node}[C] \rrbracket_t(S)$$

By Definition 4.4, we have

$$\llbracket \text{self} :: \text{node}[C] \rrbracket_t(S) = \bigcup_{\mathbf{i} \in T} \mathbf{i}$$

where  $T$  is the set of ids satisfying the predicate:

$$T = \{\mathbf{i} \mid \mathbf{i} \in S \wedge \mathbf{Check}_t[C](\mathbf{i})\}$$

Let us consider  $\mathbf{i} \in T$ . We have that  $\mathbf{i} \in S$ , and since  $\mathbf{i}$  is part of an  $\mathcal{J}$ -valid tree  $t$ , there exists  $X_i \rightarrow R_i = \mathcal{J}(\mathbf{i})$ . Now, by induction on the derivation of

$$\{X_i \rightarrow R_i\}, \kappa \vdash_E^{Cond} C : b$$

it is clear that if  $\mathbf{Check}_t[C](\mathbf{i})$  and  $\mathcal{J}(\mathbf{i}) = X_i \rightarrow R_i$  then  $b = \text{true}$  and therefore,  $X_i \rightarrow R_i \in \tau'$ , hence

$$\mathcal{J}(\llbracket \text{self} :: \text{node}[C] \rrbracket_t(S)) \subseteq \tau'$$

which proves Property 1. Property 2 holds for the same argument as in rule (**test**).

(*sequence*). : Property 1 is true by induction hypothesis on both premises. Property 2 is true for the first premise, by induction hypothesis. In particular,  $\Sigma''_{\text{ctx}}$  is a well-formed context. We can then apply the induction hypothesis on  $\Sigma''$  and we have that  $\Sigma'_{\text{ctx}}$  is a well-formed context too.

(*union*). : is similar to the previous case.

□

LEMMA (WITNESS OF A GRAMMAR (5.7)) *Let  $(\mathcal{S}, E)$  be a non-recursive, \*-guarded, parent-unambiguous local tree grammar. There exists a document  $t$ ,  $\mathcal{J}$ -valid with respect to  $(\mathcal{S}, E)$  such that:*

$$\forall X \rightarrow R \in E, \exists \mathbf{i} \in \mathbf{Ids}(t) \text{ such that } \mathcal{J}(\mathbf{i}) = X \rightarrow R$$

*we call such a document a witness of the schema  $(\mathcal{S}, E)$ .*

PROOF. Since the tree grammar is non recursive and parent unambiguous, we can prove the lemma by induction (I) on the height of the grammar, seen as a DAG.

**Basic case.**: the grammar has height 1. It consists therefore of a single rule. The rule is either  $X \rightarrow \text{String}$  and a document  $s_i$  is a suitable witness; or the rule is  $X \rightarrow a[\ ]$  for some label  $a$  and the document  $a_i[\ ]$  is a witness of the grammar.

**Inductive case.**: Consider  $(\{X\}, E)$  (the start symbol is unique since we consider a local tree grammar). The rule for the start symbol  $X$  is  $X \rightarrow a[r_1 \cdots r_n]$  for some label  $a$  (since  $E$  is

\*-guarded, the rule must have this shape). We show by induction **(II)** on the structure of the regular expression  $r_i$  that there is a witness for this regular expression.

**Basic case:**. Either  $r_i = \varepsilon$ , and therefore the empty forest  $()$  is a suitable witness. Or  $r_i = Z$ . Then consider the grammar  $(\{Z\}, E')$  where  $E' = \{Y \rightarrow R \mid Y \rightarrow R \in \mathbf{A}_E(\text{descendant}, Z)\}$ , that is the restriction of  $E$  to  $Z$ . Then,  $\text{height}(E') < \text{height}(E)$  since at least the rule associated with  $X$  is not in  $E'$  (and because  $E$  is not recursive and parent unambiguous). Therefore, by induction hypothesis **(I)** there exists a witness  $t_z$  for  $Z$ .

**Inductive case:**. Either  $r_i = (r'_i | r''_i)^*$  and by induction hypothesis **(II)**, there is a witness  $t'_i$  for  $r'_i$  and  $t''_i$  for  $r''_i$ . Then, the forest  $t'_i, t''_i$  is a witness for  $r_i$  (the first iteration of  $*$  matches  $t_i$  and the second one matches  $t''_i$ ).

Or  $r_i = (r'_i)^*$  and  $r'_i$  is not a union. Then by induction hypothesis **(II)**, there is a witness  $t'_i$  and  $t'_i$  is also a witness for  $r_i$ . Or  $r_i = r'_i r''_i$ . By induction hypothesis **(II)**, there is a witness  $t'_i$  for  $r'_i$  and  $t''_i$  for  $r''_i$ . Then, the forest  $t'_i, t''_i$  is a witness for  $r_i$ .

Therefore, for each  $r_i$  there is a witness  $t_i$ . Then the tree  $a[t_1 \dots t_n]$  is a witness of the rule  $X \rightarrow a[r_1 \dots r_n]$ .

□

The following corollary states that in a witness tree, all possible derivations of the grammar are represented:

**COROLLARY A.2**  $()$ . *Let  $(\{X\}, E)$  be a non-recursive, \*-guarded, parent-unambiguous local tree grammar and  $t$  be its witness. Let  $\{Y_1 \rightarrow R_1, \dots, Y_n \rightarrow R_n\} \subseteq E$ . If  $Y_{i+1} \rightarrow R_{i+1} \in \mathbf{A}_E(\text{child}, \{Y_i \rightarrow R_i\})$  for  $i \in 1..n-1$ , then there exists  $\{\mathbf{i}_1, \dots, \mathbf{i}_n\} \subseteq \mathbf{Ids}(t)$  such that*

$$\forall i \in \{2 \dots n\}, ((\mathbf{i}_{i-1}, \mathbf{i}_i) \in \mathbf{Edg}(t)) \wedge \mathcal{J}(id_{i-1}) = Y_{i-1} \wedge \mathcal{J}(\mathbf{i}_i) = Y_i$$

**PROOF.** This is a direct application of Lemma 5.7. We know that for all  $Y \rightarrow R \in E, \exists \mathbf{i} \in \mathbf{Ids}(t)$  such that  $\mathcal{J}(\mathbf{i}) = Y \rightarrow R$ . This is true in particular for  $\{Y_1, \dots, Y_n\}$ . Consider  $Y_i$  and  $Y_{i+1}$ . We have  $Y_{i+1} \rightarrow R_{i+1} \in \mathbf{A}_E(\text{child}, \{Y_i \rightarrow R_i\})$  which means that in  $E$ , there is a rule  $Y_i \rightarrow a[r_i]$  for some label  $a$  and with  $Y_{i+1} \in \mathcal{L}(r_i)$ . Consequently,  $t @ \mathbf{i}_i = a[\dots, id_{i+1}, \dots]$ . Therefore,  $(\mathbf{i}_i, \mathbf{i}_{i+1}) \in \mathbf{Edg}(t)$ . □

**THEOREM (COMPLETENESS OF TYPE INFERENCE (5.8))** *Let  $(\mathcal{S}, E)$  be a \*-guarded non-recursive and parent unambiguous local tree grammar, and  $P$  a path. Let*

$$E_0 = \{X \rightarrow R \mid X \rightarrow R \in E, X \in \mathcal{S}\}.$$

*If  $(E_0, E_0) \vdash_E P : (\tau, \kappa)$  then:*

$$\tau \subseteq \bigcup_{t \in \mathcal{J}E} \mathcal{J}(\llbracket P \rrbracket_t(\mathbf{RootId}(t)))$$

**PROOF.** Like for the proof of Theorem 5.5, we consider the following, more general judgment:

$$(\tau, \kappa) \vdash_E P : (\tau', \kappa')$$

let  $t$  be the witness of  $E$ . We show that if  $\tau \subseteq \mathcal{J}(S)$  then,  $\tau' \subseteq \mathcal{J}(\llbracket P \rrbracket_t(S))$ . If this holds for the witness  $t$  then it holds for the union of all trees  $\mathcal{J}$ -valid w.r.t to  $E$  (which contains  $t$ ). Informally,



this means that if the type  $\tau$  “describes precisely” the nodes in  $S$ , that is, if there are no unneeded rules in  $\tau$ , then the type  $\tau'$  describes exactly the result of the query: for each rule in  $\tau'$ , there is a node in the result of the query typed by that rule. We proceed by induction on the depth of the typing derivation:

**Basic case:**

**(down-axis):** self axis: We supposed  $\tau \subseteq \mathcal{J}(S)$ . We have  $\tau' = \mathbf{A}_E(\tau, \text{self}) = \tau$ . We also have:

$$\llbracket \text{self} :: \text{node} \rrbracket_t(S) = S$$

by Definition 4.1. Therefore,  $\tau' \subseteq \mathcal{J}(S)$  and so

$$\tau' \subseteq \mathcal{J}(\llbracket \text{self} :: \text{node} \rrbracket_t(S))$$

descendant axis: Let us consider rules,  $X \rightarrow R \in \tau$  and  $Y \rightarrow R' \in \mathbf{A}_E(\{X \rightarrow R\}, \text{descendant})$ . By using Corollary A.2, we have that there exists a sequence:  $\mathbf{i}_1, \dots, \mathbf{i}_n$  in  $t$  such that  $X \rightarrow R = \mathcal{J}(\mathbf{i}_1)$  and  $Y \rightarrow R' = \mathcal{J}(\mathbf{i}_n)$ . We also have that

$$\forall i \in \{1 \dots n-1\}, (\mathbf{i}_i, \mathbf{i}_{i+1}) \in \mathbf{Edg}(t)$$

thus  $(\mathbf{i}_1, \mathbf{i}_n) \in \mathbf{Edg}^+(t)$  and therefore that

$$\mathbf{i}_n \in \llbracket \text{descendant} :: \text{node} \rrbracket_t(\{\mathbf{i}_1\})$$

Subsequently:

$$\mathbf{A}_E(\{X\}, \text{descendant}) \subseteq \mathcal{J}(\llbracket \text{descendant} :: \text{node} \rrbracket_t(\{\mathbf{i}_1\}))$$

child axis: is a particular instance of the previous case.

**(up-axis)**. We only treat the case of the ancestor axis, of which the parent axis is a particular instance. This case is the symmetric of the descendant axis. Let  $X \rightarrow R \in \tau$ . Let  $Y \rightarrow R' \in \mathbf{A}_E(\{X\}, \text{ancestor}) \cap \kappa$ . By using Corollary A.2, we have that there exists a sequence:  $\mathbf{i}_1, \dots, \mathbf{i}_n$  in  $t$  such that  $Y \rightarrow R' = \mathcal{J}(\mathbf{i}_1)$  and  $X \rightarrow R = \mathcal{J}(\mathbf{i}_n)$ . We also have

$$\forall i \in \{1 \dots n-1\}, (\mathbf{i}_i, \mathbf{i}_{i+1}) \in \mathbf{Edg}(t)$$

thus  $(\mathbf{i}_1, \mathbf{i}_n) \in \mathbf{Edg}^+(t)$  and therefore

$$\mathbf{i}_n \in \llbracket \text{ancestor} :: \text{node} \rrbracket_t(\{\mathbf{i}_1\})$$

Thus, we have

$$\mathbf{A}_E(\{X\}, \text{ancestor}) \subseteq \mathcal{J}(\llbracket \text{ancestor} :: \text{node} \rrbracket_t(\{\mathbf{i}_1\}))$$

We must also show that if  $Y \rightarrow R' \in \kappa$  then

$$Y \rightarrow R' \in \mathcal{J}(\llbracket \text{ancestor} :: \text{node} \rrbracket_t(\{\mathbf{i}_1\}))$$

(because the output type is intersected with the context for this rule). This is an immediate consequence of the well-formedness of contexts.  $\kappa$  is well-formed only if  $\kappa \in \tau \cup \mathbf{A}_E(\tau, \text{ancestor})$ .

**(test)**: is similar to the case self of Rule **(down-axis)**.

**Inductive case:**

**(predicate)**. First, since  $t$  is a witness, for each  $\mathbf{i} \in S$ , there exists  $X_i \rightarrow R_i \in \tau$ . Second, consider We consider the premise:

$$(\{X_i \rightarrow R_i\}, \kappa) \vdash_E^{\text{Cond}} C : b$$

By induction on  $C$  and the derivation of  $(\{X_i \rightarrow R_i\}, \kappa) \vdash_E^{\text{Cond}} C : b$ , it is clear than  $b \Rightarrow \mathbf{Check}_t[C](\{\mathbf{i}_i\})$  (for the basic case (rule **(path)**) we can apply the main induction hypothesis on the path which is an atom of  $C$ , and the inductive case (rule **(boolean)**) is straightforward). Therefore, whenever  $X_i \rightarrow R_i \in \tau'$ , we have that  $b$  is true, thus  $\mathbf{Check}_t[C](\{\mathbf{i}_i\})$  and therefore  $\mathbf{i} \in \llbracket \text{self} :: \text{node}[C] \rrbracket_t(\{\mathbf{i}\})$ , which proves  $\tau' \subseteq \mathcal{J}(\llbracket \text{self} :: \text{node}[C] \rrbracket_t(S))$

(*seq.*). : By applying straightforwardly the induction hypothesis on the premises.  
 (*union*). : By applying straightforwardly the induction hypothesis on the premises.

□

LEMMA (TERMINATION OF TYPE-PROJECTOR INFERENCE (5.10)) *Let  $(\mathcal{S}, E)$  be a type,  $P$  a path, and  $\Sigma$  and  $\Sigma'$  environments. The judgment  $\Sigma \Vdash_E P : \Sigma'$  has a unique and finite derivation. Furthermore the algorithm runs in time  $O(|E|^{o(d)} \times |P|^2)$  where  $d$  is the maximum number of nested predicates*

PROOF. The uniqueness of the derivation follows from the fact that all the rules are mutually exclusive (although not strictly syntax directed) thanks to their side conditions.

To prove termination, we need some more care than for the type inference algorithm. For the judgment:

$$\Sigma \Vdash_E P : \Sigma'$$

we give it as weight the triple  $(l(P), r(P), |\Sigma_{\text{typ}}|)$  ordered lexicographically, where:

$l(P)$ . is the length of the path  $P$ , as defined previously  
 $r(P)$ . is the number of occurrences of a recursive step, that is the number of occurrences of `descendant::node` or `ancestor::node` in the  $P$   
 $|\Sigma_{\text{typ}}|$ . is the number of rules in the input

The proof is straightforward and consists that for every rule the weight strictly decreases in the premises:

**Basic case.**: the base of induction is an application of (**p-step**) or (**p-erase**) does not have any premises.

**Inductive case.**: For the rules (**p-union**), (**p-test**) and (**p-predicate**), the weight strictly decreases in  $l(P)$  in the premises. For the rule (**p-iterate**),  $|\Sigma_{\text{typ}}|$  strictly decreases in the premises, since in the conclusion the weight has at least two for this component and exactly one in each of the premises. Also,  $P$  is unchanged in the premises therefore  $l(P)$  and  $r(P)$  do not increase. For the rule (**p-many**), the  $l(P)$  part is unchanged in the premises since  $l(\text{descendant::node}/P) = l(\text{child::node}/P) = 1 + l(P)$  and  $r(P)$  decreases strictly.

The complexity result is a straightforward consequence of the one given by Lemma 5.4. If one ignores the recursive calls to the typing algorithm, it is clear that during projection inference, one perform at most  $O(|E|)$  operations and performs at most one recursive call to a strict sub-path of the argument path. □

LEMMA (WELL-FORMEDNESS OF TYPE-PROJECTOR INFERENCE (5.11)) *Let  $(\mathcal{S}, E)$  be a type,  $\tau$ ,  $\tau'$ , and  $\kappa$  sets of rules, and  $P$  a path. If  $(\tau, \kappa) \Vdash_E P : \tau'$ , then  $(\tau, \kappa) \vdash_E P : (\tau'', \kappa'')$  implies  $\kappa'' \subseteq \tau'$ .*

PROOF. We use a structural induction on the derivation of  $(\tau, \kappa) \Vdash_E P : \tau'$  which is finite by Lemma 5.10.

**Basic case:** The property is trivially true for the rule (**p-step**) since the result is the union of the output type and its associated context.

Rule (**p-erase**) can only be applied if the side conditions of the other rules fail, which means in the case where the judgment  $(\tau, \kappa) \vdash_E P : (\tau'', \kappa'')$  does not hold. Therefore the lemma is true too in that case.

**Inductive case:**

(**p-union**). We suppose  $(\tau, \kappa) \vdash_E P_1 P_2 : (\tau'', \kappa'')$ . This means that the typing rule (**union**) (cf. Figure 1) holds and that:

$$(\tau, \kappa) \vdash_E P_1 : (\tau'_1, \kappa'_1)$$

and

$$(\tau, \kappa) \vdash_E P_2 : (\tau'_2, \kappa'_2)$$

let  $P_1$  produce a type-projector  $\tau'_1$  and  $P_2$  produce  $\tau'_2$ ; by induction hypothesis  $\kappa'_1 \subseteq \tau'_1$  and  $\kappa'_2 \subseteq \tau'_2$ . But since  $\kappa'' = \kappa'_1 \cup \kappa'_2$ , we have  $\kappa'' \subseteq \tau'_1 \cup \tau'_2 = \tau'$

(**p-iterate**). similar to the previous case.

(**p-test**). we suppose

$$(\{Y \rightarrow R\}, \kappa) \vdash_E \text{self} :: \text{Test}/P : (\tau'', \kappa'')$$

According to the (**test**) typing rule, this means that:

$$(\{Y \rightarrow R\}, \kappa) \vdash_E \text{self} :: \text{Test} : (\tau_1, \kappa_1)$$

and

$$(\tau_1, \kappa_1) \vdash_E P : (\tau'', \kappa'')$$

the induction hypothesis can be applied on the second premise of the rule (**p-test**) and we have  $(\tau_1, \kappa_1) \vdash_E P : \tau'$  with  $\kappa'' \subseteq \tau'$ . Since  $\tau' \subseteq \{Y \rightarrow R\} \cup \tau'$ , we have  $\kappa'' \subseteq \{Y \rightarrow R\} \cup \tau'$  which proves this case.

(**p-predicate**). similar to the previous case, we can observe that the context resulting of the typing of the first step is passed as argument for the inference of the projector of the remainder of the path.

(**p-single**). similar to the previous case.

(**p-many**). we only treat the case for  $\text{Axis} = \text{descendant}$ , the case for ancestor being similar. We suppose:

$$(\{Y \rightarrow R\}, \kappa) \vdash_E \text{descendant} :: \text{node}/P : (\tau_0, \kappa_0) (*)$$

and we want to show that  $\kappa_0 \subseteq \{Y \rightarrow R\} \cup \tau' \cup \tau''$ . Let us write:

$$(\{Y \rightarrow R\}, \kappa) \vdash_E \text{descendant} :: \text{node} : (\tau_1, \kappa_1) \quad \mathbf{(1)}$$

$$(\{Y \rightarrow R\} \cup \tau_1, \kappa) \vdash_E \text{child} :: \text{node} : (\tau_2, \kappa_2) \quad \mathbf{(2)}$$

then  $\tau_1 = \tau_2$  and  $\kappa_1 = \kappa_2$ . Indeed:

$$\mathbf{(1)} \begin{cases} \tau_1 = \mathbf{A}_E(\{Y \rightarrow R\}, \text{descendant}) \\ \kappa_1 = \tau_1 \cup \kappa \end{cases}$$

$$\mathbf{(2)} \begin{cases} \tau_2 = \mathbf{A}_E(\{Y \rightarrow R\} \cup \tau_1, \text{child}) \\ \kappa_2 = \tau_2 \cup \kappa \end{cases}$$

by definition of  $\mathbf{A}_E(\_, \_)$ :

$$\mathbf{(1)} \begin{cases} \tau_1 = \mathbf{A}_E(\{Y \rightarrow R\}, \text{child}) \cup \mathbf{A}_E(\mathbf{A}_E(\{Y \rightarrow R\}, \text{child}), \text{descendant}) \\ \kappa_1 = \tau_1 \cup \kappa \end{cases}$$

$$\mathbf{(2)} \begin{cases} \tau_2 = \mathbf{A}_E(\{Y \rightarrow R\} \cup \tau_1, \text{child}) \\ \kappa_2 = \tau_2 \cup \kappa \end{cases}$$

For (2), we have:

$$\begin{aligned}\tau_2 &= \mathbf{A}_E(\{Y \rightarrow R\} \cup \tau_1, \text{child}) \\ \tau_2 &= \mathbf{A}_E(\{Y \rightarrow R\}, \text{child}) \cup \mathbf{A}_E(\tau_1, \text{child}) \\ \tau_2 &= \mathbf{A}_E(\{Y \rightarrow R\}, \text{child}) \cup \mathbf{A}_E(\mathbf{A}_E(\{Y \rightarrow R\}, \text{descendant}), \text{child}) \\ \tau_2 &= \mathbf{A}_E(\{Y \rightarrow R\}, \text{child}) \cup \mathbf{A}_E(\mathbf{A}_E(\{Y \rightarrow R\}, \text{child}), \text{descendant}) \\ \tau_2 &= \tau_1 \\ \kappa_2 &= \tau_2 \cup \kappa = \tau_1 \cup \kappa = \kappa_1\end{aligned}$$

therefore, for the path  $P$

$$(\{Y \rightarrow R\} \cup \tau_1, \kappa) \vdash_E \text{child}::\text{node}/P : (\tau_0, \kappa_0)$$

by application of the (**seq**) typing rule. We can finally remark that

$$((\{Y \rightarrow R\} \cup \tau_1) \cap \tau', \kappa) \vdash_E \text{child}::\text{node}/P : (\tau_0, \kappa_0)$$

all the rules ins  $\tau_1 \setminus \tau'$  yield an empty projector. Therefore

$$(\tau', \kappa) \vdash_E \text{child}::\text{node}/P : (\tau_0, \kappa_0)$$

which allows us to apply the induction hypothesis on the third premise of (**p-many**), which gives us  $\kappa_0 \subseteq \tau''$ , and therefore:  $\kappa_0 \subseteq \{Y \rightarrow R\} \cup \tau' \cup \tau''$

□

**THEOREM (SOUNDNESS OF TYPE-PROJECTOR INFERENCE (5.12))** *Let  $(\mathcal{S}, E)$  be a type and  $P$  an  $XPath^l$  query. Let  $S$  be the set of rules:  $S = \{X \rightarrow R \mid X \in \mathcal{S}\}$ . If*

$$(S, S) \Vdash_E P : \tau$$

*then  $\tau$  is a type-projector for  $(\mathcal{S}, E)$  and for every  $t \in_{\mathcal{J}}(\mathcal{S}, E)$  we have:*

$$\llbracket P \rrbracket_{t \setminus_{\mathcal{J}} \tau}(\mathbf{RootId}(t)) = \llbracket P \rrbracket_t(\mathbf{RootId}(t))$$

**PROOF.** By simple structural induction on the path. □

**THEOREM (COMPLETENESS OF PROJECTOR INFERENCE (5.14))** *Let  $(\mathcal{S}, E)$  be a  $*$ -guarded, non-recursive, and parent-unambiguous local tree grammar, and  $P$  a strongly-specified  $XPath^l$  path. Let  $S$  be the set of rules:  $S = \{X \rightarrow R \mid X \in \mathcal{S}\}$ . If*

$$(S, S) \Vdash_E P : \tau$$

*then there exists  $t \in_{\mathcal{J}}(\mathcal{S}, E)$  such that for each  $Y \rightarrow R \in \tau$ , if  $\pi = \tau \setminus (\{Y \rightarrow R\} \cup \mathbf{A}_E(\{Y \rightarrow R\}, \text{descendant}))$ , then:*

$$\llbracket P \rrbracket_{t \setminus_{\mathcal{J}} \pi}(\mathbf{RootId}(t)) \neq \llbracket P \rrbracket_t(\mathbf{RootId}(t))$$

**PROOF.** By induction on the length of the derivation which is finite. We use Theorem 5.8 to show that if we remove a name  $Y$  inferred by the type inference algorithm, then we remove

nodes from the result of the query applied to the projected document. The fact that  $P$  is strongly specified is used for the treatment of predicates. Indeed, it forces any path in a predicate to be matched exactly by one node. If a path in a predicate could be matched by two (or more) nodes, then removing one of the nodes would not change the semantics of the query, since there would still be a node present to make the predicate succeed.  $\square$

### Detailed answers to the second round of reviewer comments

We thank both reviewers for their insightful comments and their in-depth review of our work. Since Reviewer 1 did not have any major comments (we have silently fixed the typos (s)he mentioned), we address the questions raised by Reviewer 2 in detail. Related modifications are in highlighted in orange.

#### S.2. Reviewer Two

We detail our answers to the three main points raised by the reviewer as well as the technical comments or minor points whose answer is not a straightforward fix (typos, more explanations when required and so on). Note that this section is labeled S.2 (Second round of reviews, reviewer two) to disambiguate its references from references to our previous answers (R.2.x).

COMMENT S.2.1. *On the reply to Comment R.2.1:*

*I was indeed very confused by the combination of the sentence "Indeed, production rules such as  $X \rightarrow l_1[R_1], X \rightarrow l_2[R_2]$  are \*necessary\* if one wants to encode complex schemas" and the motivating example on page 16. The issue with the motivating example is now clear. However, I still think that the authors need to reword the above sentence, because I believe that it is false as stated.*

*The high-level reason why I do not believe it is that E-DTDs define the set of unranked regular tree languages and so do regular tree grammars.*

*On a more detailed level, E-DTDs actually \*are\* regular tree grammars in which productions such as  $X \rightarrow l_1[R_1], X \rightarrow l_2[R_2]$  do not occur. In fact, it is possible to translate a regular tree grammar to an equivalent E-DTD in polynomial time by the following transformation:*

*As long as there is a non-terminal  $X$  with rules  $X \rightarrow l_1[R_1], \dots, X \rightarrow l_n[R_n]$  and  $n > 1$ , do the following:*

- Create  $n$  fresh non-terminals  $X_1, \dots, X_n$
- Make  $n$  new rules  $X_1 \rightarrow l_1[R_1], \dots, X_n \rightarrow l_n[R_n]$
- Replace every occurrence of  $X$  in the grammar by the disjunction  $(X_1 + \dots + X_n)$
- Remove non-terminal  $X$

*Therefore, I still do not understand the above sentence, claiming the necessity of production rules of the form  $X \rightarrow l_1[R_1], X \rightarrow l_2[R_2]$ .*

*I do not want to say that the authors should change their definition of grammars, but the motivation of these rules should be changed. The main reason why I think that rules of the kind  $X \rightarrow l_1[R_1], X \rightarrow l_2[R_2]$  could be convenient is that they may simplify Definition 5.6. Since the above transformation introduces new disjunctions, I assume that defining a similar class of grammars than you have now in Definition 5.6 would require a more intricate condition for \*-guardedness. Therefore, the authors could perhaps say that they allow production rules  $X \rightarrow R_1, X \rightarrow R_2$  in order to simplify the definition of \*-guarded grammars later in the paper.*

**Answer:** Indeed, this remark is completely justified and it is clear that by using E-DTDs, we can normalize our types so as to avoid having two rules with the same left-hand side. We thank the reviewer for his/her insight on the matter. It seems we presented and mixed in our presentation the two following facts:

- (1) The need for the types to support top-down non-determinism (in order to express Regular Tree languages). Note that this is also the case in the construction described by the reviewer where the  $X_1 + \dots + X_n$  may yield a non deterministic (top-down) choice.
- (2) The fact that having rules of the form  $X \rightarrow R_1, X \rightarrow R_2$  provides non-determinism but is not necessary per se.

We have reworded completely the offending statement offering another justification for the usefulness (but not necessity) of ambiguous rules:

- It indeed simplifies the specification of \*-guarded grammars, as suggested by the reviewer

— It also allows us to freely take the union of two type projectors without introducing the normalization step described by the reviewer

COMMENT S.2.2. *On the reply to Comment R.2.2:*

*My comment was indeed motivated by the question why the authors did not consider using existing algorithms for XPath satisfiability or validity as a black box. However, I am not sure whether I understand the authors' answer regarding the problems one may run into wrt the typing context.*

*I had the impression that a relatively simple algorithm that uses XPath satisfiability tests as a black box could process the DTD and query using Definition 5.1 and could remove types from the DTD (or tree grammar) for which the remainder of the XPath query is unsatisfiable. Such an algorithm would be, just like the algorithm in the paper, a refinement of the very naive algorithm that directly builds upon Definition 5.1.*

*More precisely, in case that the current XPath axis is downward navigational ("child" or "descendant"), the algorithm could consider the rules  $X \rightarrow a[R]$  in the current context and delete all symbols from  $R$  for which the remaining XPath query is unsatisfiable. (As such, the obtained rule  $X \rightarrow a[R']$  would be the type projection of the rule  $X \rightarrow a[R]$ .) If, for a certain rule  $X \rightarrow a[R]$ , we have computed a projection in this fashion, we mark it, and we will not compute an (even smaller) projection for it later on.*

*It was my belief that such an approach could also be used in combination with Def.5.1 in order to avoid the problems stated in, e.g., the new Example on p.15.*

*Indeed; consider the grammar*

$$A \rightarrow a[BC], B \rightarrow b[D], C \rightarrow c[D], D \rightarrow d[]$$

*and the query `self::a / child::b / child::d / parent::node` as on p.15.*

*This algorithm would obtain in the first step  $\{A \rightarrow a[BC]\}$  and remaining XPath query `child::b / child::d / parent::node`*

*This XPath query is unsatisfiable for the C-child of this rule and therefore we would change the DTD to  $A \rightarrow a[B], B \rightarrow b[D], C \rightarrow c[D], D \rightarrow d[]$*

*Furthermore, since we know that unreachable and useless rules in context-free grammars can be easily removed (in polynomial time), we can simplify the grammar to  $A \rightarrow a[B], B \rightarrow b[D], D \rightarrow d[]$*

*The rest of the example would go analogous to the one in the paper.*

*Analogously, for the example on p.17, the algorithm would immediately simplify the DTD to  $A \rightarrow a[B], B \rightarrow b[D], D \rightarrow d[]$  in the first step.*

*For the example on p.23, the algorithm would simplify the DTD to  $A \rightarrow a[(B)*], B \rightarrow b[D], D \rightarrow d[E], E \rightarrow e[]$  in the first step. After the first step, no more rules from the DTD would be removed.*

*I will provide a somewhat more elaborate example to illustrate how filters in XPath could be handled:*

*Consider the DTD  $A \rightarrow a[(BCD)*], B \rightarrow b[D], C \rightarrow c[E], D \rightarrow d[E], E \rightarrow e[]$  and XPath query `self::a [child::d] / child::b / child::d / parent::node`*

*Here, the algorithm would, in the first step, consider the context  $A \rightarrow a[(BCD)*]$  and consider \*both\* subqueries `child::d` and `child::b / child::d / parent::node` for satisfiability tests.*

*(That is, we consider all subqueries below the `self::a` root to test for satisfiability. We keep all non-terminals in  $(BCD)^*$  for which at least one subquery is satisfiable.)*

*As such, `child::d` is found to be satisfiable at  $D$  and `child::b / child::d / parent::node` is found to be satisfiable at  $B$ . Therefore, the rule for  $A$  is projected onto  $A \rightarrow a[(BD)*]$  and the rule  $C \rightarrow C[E]$  is removed from the DTD because it became unreachable.*

*The projected type would be  $A \rightarrow a[(BD)*], B \rightarrow b[D], D \rightarrow d[E], E \rightarrow e[]$  (Or in other words, the type projector would be  $\{A, B, D, E\}$ .)*

*Lastly, one could, of course, remove all types from the grammar that have never been visited in the procedure; since they seem to be irrelevant for the query.*

*I understand that satisfiability of XPath in the presence of DTDs is NP-hard rather quickly but, as I understand it, the above procedure would run in PTIME, given an oracle for solving XPath satisfiability in the presence of DTDs. As such, the complexity of the algorithm seems to be similar to the complexity of XPath satisfiability for all fragments considered. Furthermore, it would compute a type projection for the schema at the same time. I may have missed it, but as far as I could see, the present paper does not claim a time upper bound for the complexity of computing a projector. (Lemma 5.10 would, e.g. be a place to state such a time complexity.) As such, I was confused about why the authors did not claim upper bounds on computational complexity; whereas it seemed to me that upper bounds could actually be obtained from XPath satisfiability. Also, it seemed to me that the above procedure could compute a type projector very similar to the type projector presented in this paper.*

*Would it be a worse type projector?*

*I understand that the authors already have implemented their current type projector and I am not asking to deeply consider the type projector I outlined above. I simply wanted to better understand the reasons of the authors to not consider an "XPath satisfiability" algorithm for obtaining type projectors.*

**Answer:** We believe that indeed, in the case where one only considers downward axes (child and descendant), then an XPath Satisfiability algorithm could replace our own typing algorithm during the computation of type projectors.

However, when considering also backward axes, it seems that this satisfiability algorithm would have to be tailored a bit to preserve the soundness and the precision of our algorithm. We consider the DTD given by the reviewer and the slightly modified path

$$A \rightarrow a[(BCD)*], B \rightarrow b[D], C \rightarrow c[E], D \rightarrow d[E], E \rightarrow e[ ]$$

and the path

$$/self::a[child::d]/child::b/child::d[ancestor::a/child::c]$$

Applying the same procedure as the reviewer we find that also in that case, the subquery `child::d` is satisfiable at  $D$  and `child::b/child::d[ancestor::a/child::c]` is satisfiable at  $B$ . But here, if we remove  $C$  we find later on that `ancestor::a/child::c` is not satisfiable anymore. So we have two choices for this case: either we tailor the satisfiability algorithm to return the information of the types used to compute its answer (meaning we are not using it as a blackbox) or we need to make the projection algorithm more involved (that is when we arrive at the `child::d[ancestor::a/child::c]` part of the query, evaluate it on the *original* DTD and add back in the projector the  $C \rightarrow c[E]$  rule that was removed.

As for the possible loss in precision (in the case we use an XPath Sat. Algorithm) again the problem seem to stem from upward axes within predicates. Again thanks to our use of *context*, the typing of upward axes is precise. In a sense, if the information in the context is enough to satisfy an upward axis, then we do not add other types to the type-projector. It is unclear whether an XPath Sat. Algorithm used as a blackbox allows us to reproduce this optimization.

Again we do not wish to say that previous work on XPath satisfiability cannot be used to compute type projectors, but we feel that they would need to be modified to ensure soundness and precision (in the presence of upward axes within filters where the *context* computed by our type inference plays a crucial role) of projector inference.

Our approach with this paper as been a different one: give a self-contained (albeit ad-hoc) method for performing document projection based on schema information, *that handle both upward and downward axes natively* (this is one of the main distinction with the work in the literature) and show that it performs well in practical cases (by providing extensive experimental results).

An interesting but maybe slightly off-topic question for this paper is to see how our typing algorithm compares to existing XPath Satisfiability procedures, in the restricted case of downward axes.



COMMENT S.2.3. *I have another misunderstanding in the author's response. I do not understand "The algorithm we give is fixed-parameter tractable (if the depth of nested predicates is bounded then our algorithm runs in polynomial time)." An algorithm that runs in time  $O(E^d \times P)$  (Lemma 5.4) is not necessarily fixed-parameter tractable. For fixed-parameter tractability one would need complexity of, e.g.,  $O(F^d \cdot \text{poly}(E, P))$ , where  $F$  is a constant that does not depend on the input and  $\text{poly}(E, P)$  is a polynomial in  $E$  and  $P$ . However,  $d$  can depend on the input.*

*Would it be possible for the authors to give a time complexity for Lemma 5.10?*

**Answer:** Again the reviewer is correct, we were too hasty in our answer and the complexity we give in the paper does not imply fixed-parameter tractability (which, fortunately we did not mention in the paper).

Also, following the same reasoning as the reviewer in the previous point, we give now a complexity result for Lemma 5.10 (which is polynomial time when the type-inference algorithm is considered as an Oracle). See also our answer to point S.2.16 below for the running time in practice of that algorithm.

COMMENT S.2.4. *Technical Comment:*

*I disagree with the statement in the last paragraph of Section 3 that, for restrained competition tree grammars, pruning might require a full buffering of the input document. [...] Notice that, since restrained-competition EDTDs are the one-pass preorder typeable schemas, we know the type of an element at the point where its opening tag is met. Furthermore, note that, in order to compute the projection of the XML document, we do \*not\* need to compute the projection of the grammar. It suffices to have the original (restrained competition) grammar, and the sets of types to prune in each rule.*

**Answer:** Indeed, in the case of restrained competition EDTDs, since the type of an element is known on the opening tag, we can choose the correct pruning rule accordingly and decide whether to keep this element or not. We have corrected the offending paragraph accordingly (with a citation to "Which XML schemas admit 1-pass preorder typing" (Martens et al, ICDT 2005).

COMMENT S.2.5. *p7:10 we have  $\text{Names}(E) = \text{Dn}(E)$  I missed why this is implied by the definitions. Don't we only have  $\text{Names}(E) \subseteq \text{Dn}(E)$ ?*

**Answer:** Correct. We have modified the sentence accordingly.

COMMENT S.2.6. *p7:36-43 This definition confuses me. What do you require the mapping to do for IDs that have an ancestor mapped to  $X \rightarrow \text{Any}$ ? Do they have to map to anything?*

**Answer:** We fixed the definition so that descendants of nodes that map to  $X \rightarrow \text{Any}$  also map to  $X \rightarrow \text{Any}$ .

COMMENT S.2.7. *p9:55 "to each node of  $t$  a non terminal of  $E$ " Is this true? What about nodes that have an ancestor mapped to  $X \rightarrow \text{Any}$ ?*

**Answer:** It is now true with the fix mentioned above.

COMMENT S.2.8. *p11:53 Notice that the name "path" can be confusing to readers familiar with XPath tree pattern queries and path queries. Traditionally, "path query" is used for queries without filter operations, i.e., without the production  $\text{Axis}::\text{Test}[\text{Cond}]$  Queries with this production are usually called "tree pattern queries" or "twig queries". (There are very many variations of such queries in the literature.) Personally, I find that the name "tree pattern" or "twig" expresses more clearly the "branching" behavior of the query.*

**Answer:** Since it would have been difficult to use "tree pattern" or "twig queries" instead of path (and it would have been more confusing in our opinion), we added a statement that our notion of "path" is sometimes called tree pattern or twig query in the literature (just before the formal definition).

COMMENT S.2.9. *p14:42 I'm confused by your notion of completeness. Wouldn't it be stronger to require that the converse inclusion holds in the equation on l 32? That is, that tau is a subset of  $\mathcal{I}([\text{path}]_i, (\text{RootID}(t)))$  ?*

**Answer:** It is true that  $\tau \subseteq \mathcal{I}([\text{path}]_i, (\text{RootID}(t)))$  is a stronger result than the one we give as equation (2). What we wanted to stress is that the “empty-precision” property is the one that is needed to ensure the precision of type-projector inference (and of course, a general completeness statement implies this property).

COMMENT S.2.10. *p16:17 When you informally introduce contexts, it may also be good to give an intuitive explanation of what a context is. In my understanding, a context has the grammar productions that type ancestors of nodes that are (currently) selected by the query.*

**Answer:** Modified as suggested by the reviewer.

COMMENT S.2.11. *p19:17 Could you elaborate on "our personal experience is that most of the DTDs available on the web are \*-guarded"? Do you conclude this from papers investigating the structure of DTDs on the Web? E.g.: - Choi: What are real DTDs like? WebDB 2002: 43-48 - Bex, Neven, Van den Bussche: DTDs versus XML Schema: A Practical Study. WebDB 2004: 79-84 - Laender, Moro, Nascimento, Martins: An X-ray on web-available XML schemas. SIGMOD Record 38(1): 37-42 (2009) Or did you perform your own study? (Actually, the word from the sentence in the paper that confuses me the most is "most". How do you know that most DTDs on the web are \*-guarded?)*

**Answer:** Clearly “most” is too strong a statement. We have reworded the sentence and added relevant references (mainly Choi).

COMMENT S.2.12. *p20:52-p21:14 This is a very long sentence that is also very subjective, in my opinion. I am not at all sure if most readers would agree with the entire sentence. Maybe it should be reworded to "However, we feel that (i) ... (ii) ... parent-ambiguity \*may\* make the event rare in practice (iii) ... seems to be negligible in most cases (iv) ..." The reason for rewriting (ii) would be that you are not providing any hard evidence for the statement. Therefore, it may be better to weaken it.*

**Answer:** Done as suggested.

COMMENT S.2.13. *p21:11 "it is well-known that exact type-inference for XPath routinely escapes regular tree languages..." I am not sure if I understand this. This observation seems to be orthogonal to what you want to do in the paper. I understand that, given a regular tree language, there are very simple transducers that output a non-regular tree language. However, you want to compute the set of types of nodes selected by an XPath query, right? This is just a set of types. Therefore, I do not fully see what the relevance of part (iv) of the remark is.*

**Answer:** In some sense, our goal was to say that using something more precise than sets of rules for contexts (and therefore for types and types projectors) would still yield an approximation at some point, and that the added complexity of managing e.g. chains of rules may not be worth it.

COMMENT S.2.14. *p24:54 Does the algorithm have the same complexity as the algorithm implied by Lemma 5.4?*

**Answer:** We have put a complexity results for this Lemma 5.10 as well.

COMMENT S.2.15. *p26:35 What do you mean by "multiple paths". It seems that the query  $a / b [ \text{child}::c \text{ and } \text{child}::d ] / e$  has multiple paths in the predicate. Which of the queries below have multiple paths in a predicate:  $\text{self}::a / \text{child}::b [ \text{child}::d [ \text{child}::e ] ] / \text{child}::f$   $\text{self}::a / \text{child}::b [ \text{child}::d [ \text{child}::e ] / \text{child}::f ] / \text{child}::g$   $\text{self}::a / \text{child}::b [ \text{child}::d / \text{child}::e ] / \text{child}::f$*

**Answer:** None of the above have multiple paths in our sense. The intuition is that for completeness to not hold it is sufficient that, from the same context node, a predicate evaluates to “true” for two possible reasons, e.g. `child:b[ child:a or child:d]`. We have made the statement less ambiguous by requiring that predicates do not contain `or` or `and` operators.

COMMENT S.2.16. *p41: Fig 8: Maybe I missed it: do your time measurements (in particular speed-up) include the time needed to compute the type projector and to compute the projection of the document? I think that it is quite important to clearly discuss this issue, because it can potentially change the scope of application of the paper. If the time for computing the type projector (etc.) is included, then the approach can be used for performing type projections in run-time. But, if not, then it seems to be more aimed towards optimizing "long standing" queries, i.e., queries that need to be evaluated against a dataset regularly. (A similar remark holds for further time measurements in this section.)*

**Answer:** This is a huge mistake on our part, and we thank the reviewer for pointing out this problem. For the practical cases we consider, computing the projector is always very fast: tens of milliseconds at most, well below query answering time for both Saxon and MonetDB tests. We have updated the table and experimental result section to highlight that fact.