

Destabilizers and Independence of XML Updates

Michael Benedikt
Computing Laboratory, Oxford University

James Cheney
LFCS, University of Edinburgh

ABSTRACT

Independence analysis is the problem of determining whether an update affects the result of a query, e.g. a constraint or materialized view. We develop a new, modular framework for static independence analysis that decomposes the problem into two orthogonal subproblems: approximating the *destabilizer*, that is, a finite representation of the set of updates that can change the result of the query, and testing whether the update and destabilizer overlap via an intersection analysis. Focusing on XML queries as the view language and the XQuery Update Facility as the update language, we present a syntactic query rewriting algorithm for translating queries to destabilizers, and show that intersection checking can be reduced to satisfiability problems for which efficient checkers already exist. We present an implementation based on an expressive tree satisfiability checker and a Satisfiability Modulo Order package, and give experiments confirming that the resulting analysis is both fast and effective.

1. INTRODUCTION

View maintenance – that is, recomputing integrity constraints or derived data as base data is updated – is a classical and well-studied problem for relational databases. It is much less well-understood for other data models, particularly XML databases. View maintenance is just as important for XML as for relational databases and is an active area of study [2, 27, 26, 7, 12].

View recomputation is avoidable if the query and update are statically *independent*, that is, there is no possible source data for which the update can affect the result of the query. In this common case, full or incremental recomputation might require time proportional to the size of the data just to determine that no work needs to be done. Query and update expressions are usually fairly small, so even a relatively expensive static analysis may be competitive with full recomputation or incremental techniques if it offers the hope of avoiding dynamic recomputation cost proportional to data size.

Static independence analysis of queries and updates has been explored both in the context of relational databases [8, 23] and XML [2, 27, 26]. For SQL, a crude independence analysis often suffices, checking whether the relation and field names in the up-

Permission to make digital or hard copies of all or part 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 bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were presented at The 36th International Conference on Very Large Data Bases, September 13-17, 2010, Singapore.

Proceedings of the VLDB Endowment, Vol. 3, No. 1
Copyright 2010 VLDB Endowment 2150-8097/10/09... \$ 10.00.

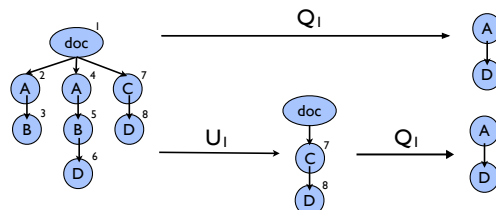


Figure 1: Document, view and an independent update

date are disjoint from those in the query. For XML there is no such “obvious analysis” – queries and updates generally run over the same document tree, and variables may not have a syntactically transparent scope. Prior techniques either require a schema [2] or apply only to downward fragments of XPath [27, 26].

We develop a new, modular framework for static independence analysis that decomposes the problem into two orthogonal subproblems: (1) identifying a *destabilizer*, that is, a finite representation of a set of run-time updates that includes all updates that could change the result of the query, and (2) statically determining whether the destabilizer and the update are disjoint. Destabilizers are a generalization of ideas such as statically approximating the set of *accessed nodes* found in some prior work [24, 6, 18]. However, the destabilizer framework provides increased flexibility and modularity, easing exploration of the tradeoff between precision and performance. Our approach, in principle, is applicable to independence analysis in any data model and query/update language; here, we apply it to XQuery/XQuery Updates.

To illustrate our approach, consider the following examples:

Example 1. Consider the situation depicted in Figure 1, where we have a view

$$Q_1 = \text{for } \$x \in \$doc/C \text{ return } \langle A \rangle \$x / D \langle /A \rangle .$$

Consider the set of all atomic update sequences that could affect the result of Q_1 . Any set containing this collection of sequences is called a *subtree value destabilizer* for Q_1 . In this example, the set of all update sequences that include an update to $\{1, 2, 4, 7, 8\}$ is a value destabilizer. We can be even more precise and say that for an update sequence to change the value of Q_1 it must contain either a deletion or replace on one of $\{1, 7, 8\}$, a rename on one of $\{2, 4, 7, 8\}$, or an insert into one of $\{1, 7, 8\}$, etc. Our analysis can determine that the update expression $U_1 = \text{delete } \$doc/A$ can never perform any of the updates above, so it is statically independent of Q_1 , as suggested in Figure 1. \diamond

Example 2. Consider the XPath query $Q_2 = \$doc/A/B$ where variable $\$doc$ again refers to the root of the tree shown in Figure 1.

We view Q_2 as a Boolean constraint asserting that there must always be a node matching path $/A/B$. To ensure the constraint is maintained, we would like to statically determine whether an update can change the result of Q_2 from nonempty (true) to empty (false). A *negative Boolean destabilizer* is a set including all update sequences that can invalidate Q_2 . In this example, the set of all sequences containing updates to one of $\{1, 2, 3, 4, 5\}$ is a negative Boolean destabilizer. If we take the update operation into account, then for an update sequence to change Q_2 to false, it must contain a delete or replace to one of $\{1, 2, 3, 4, 5\}$ or a rename to one of $\{2, 3, 4, 5\}$: insert operations can never play a role in making the result of Q_2 empty. Our analysis can determine that update $U_2 = \text{insert } \langle A \rangle \langle /A \rangle$ into $\$doc$ can never invalidate Q_2 , since it can only perform inserts. \diamond

Example 3. Consider a query and update

```

 $Q_3 = \text{if } \$doc/A \text{ then } \langle E \rangle \langle /E \rangle \text{ else } ()$ 
 $U_3 = \text{if } \$doc/A \text{ then } () \text{ else delete } \$doc/C$ 

```

Clearly, for the example document in Figure 1, Q_3 and U_3 are independent; in fact they are independent for any input, since U_3 cannot change the Boolean value of the test $\$doc/child :: A$. Our analysis can determine this by describing a destabilizer and then showing that the updates generated by U_3 can not intersect the destabilizer via reduction to a decidable tree logic. \diamond

The destabilizer framework thus works in several steps: we symbolically describe the sequences of updates that can change the query result, and also represent the sequences that can be generated by the update expression, and then use satisfiability testing to see that they cannot overlap. Though this is a clean and deceptively simple approach, there are nontrivial obstacles that must be overcome to turn it into a working, useful analysis. Calculating exact (static or dynamic) destabilizers is not feasible for our core XML query language; instead, we present a query-rewriting technique that provides a useful, sound approximation. Given an approximate destabilizer, it then suffices to establish that none of the updates that U can perform are in Q 's destabilizer. This can be reduced to solving intersection problems for XPath or XQuery expressions, but again, this is nontrivial in practice: for example, intersection testing for even simple path expressions is NP-hard [21].

Contributions. The main contributions of this work are:

- We introduce a framework for independence analysis based on the notion of destabilizer, which is, in principle, applicable to any data model and query/update language.
- We introduce two varieties of destabilizer for XQuery, called *generic* and *operation(op)-sensitive*. We provide algorithms that efficiently rewrite queries to approximate destabilizers.
- We develop novel techniques for efficiently solving the disjointness problems that arise through destabilizers, using a combination of heuristics and reduction to SMT solvers.
- We present experiments showing that independence analysis using destabilizers is fast and effective.

Organization: We begin in Section 2 by reviewing the XML query and update languages considered here. In Section 3 we give our destabilizer representations and algorithms. In Section 4 we apply these to independence analysis. Section 5 presents our experiments. Section 6 reviews related work, and Section 7 Proofs and many routine details are omitted or placed in appendices.

2. PRELIMINARIES

As is conventional in most work on XML databases (e.g. [24, 6, 26, 11, 16, 2, 18, 27]) we employ a simplified XML data model and query language, focusing on tree structure and ignoring attributes

and some other features. Generalizing to handle these features is straightforward. Throughout this paper, we use the symbol $=$ for mathematical equality, not one of XQuery's equality operations.

A *store* σ (sometimes called an *instance*, *document* or *database*) is an ordered labeled forest, whose nodes l, l', m (also referred to as *locations*) are either element nodes or text nodes. An element node has a label, while a text node has an associated string; text nodes have no children. In addition to its label or string, each node has an identifier, which is assumed to be unique within a store. We often write location sequences as L, L', L'' .

Queries. We will use a simple XQuery-like core language, denoted XQ :

```

 $q ::= s | () | q, q' | \langle a \rangle q \langle /a \rangle | \$x | \$x/step$ 
 $| \text{if } q \text{ then } q_1 \text{ else } q_2 | \text{for } \$x \in q \text{ return } q'$ 
 $step ::= ax :: \phi \quad \phi ::= a | * | \text{text}()$ 
 $ax ::= \text{self} | \text{child} | \text{parent} | \text{foll\_sib} | \text{prec\_sib}$ 
 $| \text{desc} | \text{desc\_or\_self} | \text{anc} | \text{anc\_or\_self}$ 

```

XQuery expressions include variables $\$x, \$y, \$z, \dots \in Var$. The expression s is a test literal (constructing a text node from a string), expression $()$ denotes the empty sequence; the expression $\langle a \rangle q \langle /a \rangle$ builds an XML tree with root a and content q ; and q, q' concatenates sequences of XML values. Conditionals $\text{if } q \text{ then } q_1 \text{ else } q_2$ branch depending on whether their first argument is nonempty. The iteration expression $\text{for } \$x \in q \text{ return } q'$ evaluates q , and for each node l in the result evaluates q' with $\$x$ bound to l , concatenating the results in order. Without loss of generality, we omit XQuery's *let*-binding construct; since we consider a recursion-free language, all occurrences of *let* can be inlined.

The expression $\$x/ax::\phi$ performs an XPath step starting from $\$x$, where ax is one of the standard XPath axes and ϕ is an XPath node test. The XPath axes allow navigation in the document tree. The *self* axis is the identity, and the *child* and *parent* axes link nodes to their parents and children. The *desc* and *anc* axes are the transitive closures of *child* and *parent*, and the *desc_or_self* and *anc_or_self* are their reflexive, transitive closures. The axes *foll_sib* and *prec_sib* select the following or preceding siblings of a node, respectively; other axes, such as *following*, can be built up from these using composition. The semantics of XPath axes is standard; see for example [17, 13]. The expressions $\$x/child :: \phi$ and $\$x/desc_or_self :: */child :: \phi$ are often abbreviated $\$x/\phi$ and $\$x//\phi$, respectively.

Updates. The W3C XQuery Update proposal [9] defines *atomic updates* that describe concrete changes to XML trees::

```

 $\iota ::= \text{ins}(L, d, l) | \text{del}(l) | \text{repl}(l, L) | \text{ren}(l, a)$ 
 $d ::= \leftarrow | \rightarrow | \downarrow | \swarrow | \searrow$ 

```

Here, the symbols d stand for the qualifiers on inserts (before, after, into, as first into and as last into, respectively). Atomic updates describe concrete changes to the tree using node identifiers; for example $\text{delete}(n)$ or $\text{ins}(L, \downarrow, n)$ where n is a target node and L is a node sequence addressing data that will be inserted into the child list of a target node. Atomic updates can be classified by the following *operations*:

```

 $op ::= \text{insert}(d) | \text{delete} | \text{replace} | \text{rename}$ 

```

where the insert, delete, replace and rename operations match the corresponding atomic updates in the obvious way.

The W3C proposal also defines high-level *updating expressions* u that declaratively describe updates without referring to node identifiers. Update expressions include forms such as $\text{delete } q_0$ and

insert q into q_0 that evaluate to sequences of atomic updates. In our implementation, we handle a core language for XQuery Updates similar to that formalized in [3]. Due to space limits, in the body of the paper we discuss only updates targeting XPath queries. More details are presented in Appendix A.

Semantics. Due to space limits, we cannot give the full formal semantics of queries and updates. We outline a simplified semantics that illustrates the basic ideas. Our results hold for the full query and update language described in Appendix A. The semantics of a query is given by a relation $\sigma, \gamma \models q \Rightarrow \sigma', L$ where σ is an input store, γ is an environment mapping variables to nodes in σ , σ' is an output store, and L is a list of nodes. Our semantics explicitly models nondeterministic creation of element nodes. Node construction queries allocate new trees in the store, whereas all other query constructs leave the contents of the store alone; hence, σ is always a subforest of σ' .

The semantics of atomic updates to XML trees is formalized in several places, e.g. [3], and we include a formalization in Appendix A. One subtlety relevant to our work is that we model a delete operation as detaching a node from its parent, creating a disconnected store; downward axis steps from the detached node can be evaluated as usual after the deletion occurs. We can easily modify our algorithms to account for alternative semantics for deletion.

We assume (purely for convenience) that each store has a distinguished root $root(\sigma)$, and we define a *document query* or *document update* to be a query or update with a single distinguished free variable $\$doc$ which is bound to the root of the store.

We also define a sublanguage of XQ called $SelXQ$, consisting of *selection queries* that exclude the element node construction operation $\langle a \rangle q \langle /a \rangle$. This restriction implies that queries always return nodes already present in the input and never construct new nodes. We can interpret document selection queries as maps from stores to sets of nodes as follows:

$$Sel[q]\sigma = \{l \in L \mid \sigma, [\$doc := root(\sigma)] \models q \Rightarrow \sigma, L\}$$

Equivalence. As discussed in the introduction, we are interested in several different interpretations of the results of a query, including Boolean, node, and subtree equivalence. We capture these notions precisely in the following definitions.

DEFINITION 1 (QUERY OUTPUT EQUIVALENCES). *Given a document query q and stores $\sigma, \hat{\sigma}$, let σ_1, L_1 be the result of evaluating q on σ and σ_2, L_2 be the result of evaluating q on $\hat{\sigma}$. Then we say that σ_1, L_1 and σ_2, L_2 are Boolean equivalent ($(\sigma_1, L_1) \cong_{bool} (\sigma_2, L_2)$) if either both L_1 and L_2 are empty or both are nonempty. Similarly, we say that σ_1, L_1 and σ_2, L_2 are subtree value equivalent ($(\sigma_1, L_1) \cong_{ns} (\sigma_2, L_2)$) if L_1 and L_2 are isomorphic (as hedges, that is, sequences of trees).*

Finally, we say that σ_1, L_1 and σ_2, L_2 are node equivalent if $L_1 = l_1, \dots, l_n$ $L_2 = \hat{n}_1, \dots, \hat{n}_n$ such that

- *for every $i \leq n$ if l_i is in σ or \hat{n}_i is in $\hat{\sigma}$ then $l_i = \hat{n}_i$*
- *for every $i \leq n$ if l_i is not in σ then the isomorphism type of l_i within its connected component in σ_1 is the same as the isomorphism type of \hat{n}_i within its connected component in σ_2 .*

We write $\sigma_1, L_1 \cong_{node} \sigma_2, L_2$ in this case, although the equivalence depends on the original stores $\sigma, \hat{\sigma}$ as well.

In the above definition, node equivalence says that the resulting lists return exactly the same nodes in the original store, and the new nodes that are created via node construction are isomorphic. It is easy to see that if $\sigma_1, L_1 \cong_{node} \sigma_2, L_2$, then the same holds for any other L'_1 that can result from q on σ , since these will only vary on

the identities of nodes outside of the input store. The term ‘‘isomorphism type’’ is used above with its standard meaning in model theory, regarding the store as a finite structure.

Independence. Informally, independence means that an update cannot change the result of a query. However, since both queries and updates can be nondeterministic, we need to define this carefully. We formalize this as follows.

DEFINITION 2. *Let \equiv be one of $\cong_{ns}, \cong_{node}, \cong_{bool}$. Given an update u and a query q , we say q is independent of u (modulo \equiv) on input σ_1 provided that for all σ_2 that can result from applying u to σ_1 , the possible behaviors of q on σ_1 are the same as those on σ_2 (modulo \equiv). If q and u are independent on all input stores, then we say they are statically independent.*

3. DESTABILIZERS FOR XML QUERIES

As discussed in the previous section, we can view the result of a query in a number of ways. We may only care about the Boolean value of a query, for example, if the query is only used as a Boolean test in a conditional, or as a constraint. We may care about the sequence of nodes produced (up to isomorphism). Or we may care about the complete XML tree values produced.

DEFINITION 3. *A destabilizer is a collection of update sequences containing all those sequences that change the result of a query in one of the above senses. More specifically, a Boolean destabilizer for q on σ is a collection of atomic update sequences that includes every update sequence that changes the result of q on σ modulo \cong_{bool} . A positive (negative) destabilizer is a collection containing every update sequence changing the result from false to true (respectively from true to false). A node destabilizer is a collection of update sequences containing every sequence changing the result of q on σ modulo \cong_{node} , while a subtree value destabilizer contains every sequence changing the result modulo \cong_{ns} .*

In this paper we will be most concerned with value destabilizers, but Boolean and node destabilizers are also of practical interest. Moreover, to calculate useful value destabilizers we will need to calculate positive, negative and node destabilizers as well.

3.1 Statically approximating destabilizers

For a given update sequence ω and input σ it is possible to determine dynamically whether ω destabilizes q on σ , by applying ω and rerunning q . However, this cost is exactly what we wish to avoid by a static independence analysis. It is not feasible to deal with destabilizers statically as concrete sets of sequences. Instead, we want to work with destabilizers indirectly, at a symbolic level.

Selection queries are a natural way to represent sets of nodes statically; a set of nodes S in turn can be considered as a representation of the collection of update sequences that modify a node in S . For a set of nodes S , let $SeqTargets(S)$ be the set of update sequences that contain some update with target in S . A *static value destabilizer* for q is a selection query δ such that for each σ , $SeqTargets(Sel[\delta]\sigma)$ is a value destabilizer, and analogously define positive (respectively negative, node, Boolean) static destabilizers. Unwinding the definition, this says that the query returns a set of nodes S such that every update sequence that destabilizes q on σ modifies a node in S .

The generic static destabilizer defined above is frequently a significant over-approximation of the runtime destabilizer. For example, a query may be sensitive to update sequences that can only affect a node via a particular update operation. Precision can be increased substantially using finer-grained representations of collections of update sequences. We exhibit one such refinement, using

$$\begin{aligned}
\Delta(()) &= () \\
\Delta(\$x) &= \Delta(\$x/\text{self} :: *) \\
\Delta(q, q') &= \Delta(q), \Delta(q') \\
\Delta^{b+}(\langle A \rangle q / \langle A \rangle) &= () \\
\Delta^{b-}(\langle A \rangle q / \langle A \rangle) &= () \\
\Delta^n(\langle A \rangle q / \langle A \rangle) &= \Delta^v(q) \\
\Delta^v(\langle A \rangle q / \langle A \rangle) &= \Delta^v(q) \\
\Delta(\text{if } q \text{ then } q' \text{ else } q'') &= \text{if } q \text{ then } (\Delta^{b-}(q), \Delta(q')) \\
&\quad \text{else } (\Delta^{b+}(q), \Delta(q'')) \\
\Delta^{b+}(\text{for } \$x \in q \text{ return } q') &= \Delta^{\text{for}}(q, \$x, \text{if } \neg q' \text{ then } \Delta^{b+}(q')) \\
\Delta^{b-}(\text{for } \$x \in q \text{ return } q') &= \Delta^{\text{for}}(q, \$x, \text{if } q' \text{ then } \Delta^{b-}(q')) \\
\Delta^n(\text{for } \$x \in q \text{ return } q') &= \Delta^{\text{for}}(q, \$x, \Delta^n(q')) \\
\Delta^v(\text{for } \$x \in q \text{ return } q') &= \Delta^{\text{for}}(q, \$x, \Delta^v(q')) \\
\Delta^{\text{for}}(q, \$x, q') &= \Delta^n(q), \text{for } \$x \in q \text{ return } q'
\end{aligned}$$

Figure 2: Query rewriting for destabilizers (non-axis steps)

operations to classify updates (as defined in Section 2). We then statically represent the update sequences that include an update with a particular target and a particular operation. For a set of nodes S , let $\text{SeqTargets}_{op}(S)$ be the set of update sequences that contain some update with operation op and a target in S . An indexed family of queries $(\delta_{op})_{op \in \text{Ops}}$ is a *static value op-destabilizer* if for each σ , $\bigcup_{op} \text{SeqTargets}_{op}(\text{Sel}[\delta]\sigma)$ is a value destabilizer, and we define positive, negative, node, and Boolean op-destabilizers analogously. Equivalently, every update sequence that destabilizes q on σ must contain an update of operation op that modifies a node returned by δ_{op} .

In developing a static destabilizer, we have a trade-off between the precision (how close it is to the dynamic analog) and the complexity of analyzing the resulting query. Ideally, one would calculate a query that always returns the best possible approximation of the runtime destabilizer. Formally, a *pointwise minimal* static destabilizer (or respectively op-destabilizer) is a query that returns on any document a minimal collection of nodes having the property that every update sequence that changes the query result contains an update (respectively an update of operation op) with the given node as target. However, the following result shows that minimal static boolean destabilizers can not be calculated efficiently. The proof is in Appendix B.

THEOREM 1. *There is no elementary time algorithm for constructing a pointwise minimal static destabilizer.*

On the other extreme, there is always a trivial, non-minimal destabilizer including *all* update sequences to the store σ . Of course, this is useless for static independence analysis.

Our approach is to compute a static destabilizer by a simple inductive rewriting algorithm. This is algorithmically feasible, while providing high precision in practice. We define four functions Δ^{b+} , Δ^{b-} , Δ^n and Δ^v simultaneously by mutual recursion on the structure of expressions, which will return static positive, negative, node, and value destabilizers. We also define functions Δ_{op}^{b+} , Δ_{op}^{b-} , Δ_{op}^n and Δ_{op}^v for each update operation op ; the vector of queries over each op will give a op-destabilizer.

Figure 2 shows the definitions of Δ^{b+} , Δ^{b-} , Δ^n and Δ^v for all query constructs except for $\$x/\text{step}$. We use Δ to denote any one of Δ^{b+} , Δ^{b-} , Δ^n and Δ^v . In the case of these query constructs, the very same rules can be used for each element of the static op-destabilizer (e.g. for every $op \in \text{Ops}$, $\Delta_{op}^{b+}(q, q') = \Delta^{b+}(q)$, $\Delta_{op}^{b+}(q')$, and similarly for the other complex constructs).

We now discuss the rules in Figure 2 starting from the top. The case for $()$ is obvious: no update sequence can destabilize the query $()$. Variables $\$x$ are treated the same as $\$x/\text{self} :: *$ steps. For

$$\begin{aligned}
\Delta_*^{b+}(\$x/\text{self} :: *) &= () \\
\Delta_*^{b+}(\$x/\text{self} :: A) &= \$x \\
\Delta_*^{b+}(\$x/\text{child} :: \phi) &= (\$x, \$x/\text{child} :: *) \\
\Delta_*^{b+}(\$x/\text{parent} :: \phi) &= (\$x, \$x/\text{parent} :: *) \\
\Delta_*^{b+}(\$x/\text{desc} :: \phi) &= (\$x/\text{desc_or_self} :: *) \\
\Delta_*^{b+}(\$x/\text{anc} :: *) &= (\$x/\text{anc} :: *) \\
\Delta_*^{b+}(\$x/\text{foll_sib} :: \phi) &= (\$x, \$x/\text{parent} :: *, \$x/\text{foll_sib} :: *) \\
\Delta_*^{b-}(\$x/\text{self} :: *) &= () \\
\Delta_*^{b-}(\$x/\text{self} :: A) &= \$x \\
\Delta_*^{b-}(\$x/\text{child} :: \phi) &= (\$x, \$x/\text{child} :: \phi) \\
\Delta_*^{b-}(\$x/\text{parent} :: \phi) &= (\$x, \$x/\text{parent} :: \phi) \\
\Delta_*^{b-}(\$x/\text{desc} :: \phi) &= (\$x/\text{desc_or_self} :: *) \\
\Delta_*^{b-}(\$x/\text{anc} :: *) &= (\$x/\text{anc_or_self} :: *) \\
\Delta_*^{b-}(\$x/\text{foll_sib} :: \phi) &= (\$x/\text{parent} :: *, \$x/\text{foll_sib} :: \phi) \\
\Delta_*^n(\$x/\text{self} :: *) &= () \\
\Delta_*^n(\$x/\text{self} :: A) &= \$x \\
\Delta_*^n(\$x/\text{child} :: \phi) &= (\$x, \$x/\text{child} :: *) \\
\Delta_*^n(\$x/\text{parent} :: \phi) &= (\$x, \$x/\text{parent} :: *) \\
\Delta_*^n(\$x/\text{desc} :: \phi) &= (\$x/\text{desc_or_self} :: *) \\
\Delta_*^n(\$x/\text{anc} :: *) &= (\$x/\text{anc_or_self} :: *) \\
\Delta_*^n(\$x/\text{foll_sib} :: \phi) &= (\$x, \$x/\text{parent} :: *, \$x/\text{foll_sib} :: *) \\
\Delta_*^v(\$x/\text{ax} :: \text{text}()) &= \Delta_*^n(\$x/\text{ax} :: \text{text}()) \\
\Delta_*^v(\$x/\text{step}) &= \Delta_*^n(\$x/\text{step})/\text{desc_or_self} :: *
\end{aligned}$$

Figure 3: Generic destabilizers for axis steps. We omit the cases for desc_or_self, anc_or_self, and prec_sib.

sequential composition, the translation is straightforward. In the output of the destabilizer, order is unimportant, since we will only be considering its node selection semantics. This step could thus be read as saying that we over-approximate the update sequences destabilizing q, q' as those that destabilize either q or q' .

Now consider the static destabilizers for node construction queries on lines 4–7 of Figure 2. The rules for the Boolean cases reflect that no updates can change the Boolean result of this query. For node and value equivalence, we conservatively assume that any change that could affect the subtree value of the subquery q could also affect the result. Note that it would not be sound to translate $\Delta^n(\langle A \rangle q / \langle A \rangle)$ to $()$ instead. For example, in a query such as $\text{for } \$y \in \langle A \rangle \$x \langle A \rangle \text{ return } \y , we would thereby lose track of the fact that changes to $\$x$ can affect the result (modulo $\cong_{n,s}$). For conditional queries $\text{if } q \text{ then } q' \text{ else } q''$, we mimic the structure of q . If q is nonempty, then an update sequence that destabilizes q must either change q to be empty, or destabilize q' . Dually, if q is empty, then to destabilize q , the update must either change q to be nonempty, or destabilize q'' .

For iteration queries $\text{for } \$x \in q_1 \text{ return } q_2$, we introduce a helper function $\Delta^{\text{for}}(q, \$x, q')$ that conservatively includes any updates that destabilize the node sequence returned by q , and additionally returns any updates satisfying $\text{for } \$x \in q \text{ return } q'$. For the positive destabilizer, the return sequence can also change from empty to nonempty if some update changes the value returned in some iteration from empty to nonempty. Thus, in this case the third argument to Δ^{for} returns all updates that can positively destabilize an empty result from some iteration. (In Figure 2, we use $\text{if } \neg q \text{ then } q'$ as shorthand for $\text{if } q \text{ then } () \text{ else } q'$.) The negative destabilizer is symmetric. For node and value destabilizers, we return all nodes that may destabilize the node sequence or the value returned by some iteration of the body of the loop, respectively.

We now turn to the static destabilizer functions for axis steps $\$x/\text{step}$. We give only the generic versions in Figure 3; we em-

$$\begin{aligned}
\Delta_{\text{delete}}^{\text{b}^+}(\$x/ax :: \phi) &= () \\
\Delta_{\text{insert}(\downarrow)}^{\text{b}^+}(\$x/child :: \phi) &= \$x \\
\Delta_{op}^{\text{b}^+}(\$x/step) &= \Delta_*^{\text{b}^+}(\$x/step) \\
\Delta_{\text{delete}}^{\text{b}^-}(\$x/child :: \phi) &= \$x/child :: \phi \\
\Delta_{\text{insert}(\downarrow)}^{\text{b}^-}(\$x/child :: \phi) &= () \\
\Delta_{op}^{\text{b}^-}(\$x/step) &= \Delta_*^{\text{b}^-}(\$x/step) \\
\Delta_{\text{delete}}^{\text{n}}(\$x/child :: \phi) &= \$x/child :: \phi \\
\Delta_{\text{insert}(\downarrow)}^{\text{n}}(\$x/child :: \phi) &= \$x \\
\Delta_{op}^{\text{n}}(\$x/step) &= \Delta_*^{\text{n}}(\$x/step) \\
\Delta_{op}^{\text{v}}(\$x/ax :: \text{text}()) &= \Delta_{op}^{\text{n}}(\$x/ax :: \text{text}()) \\
\Delta_{op}^{\text{v}}(\$x/step) &= \Delta_{op}^{\text{n}}(\$x/step)/\text{desc.or.self} :: *
\end{aligned}$$

Figure 4: Kind-sensitive destabilizers for axis steps

phasize this by using notation $\Delta_*^{\text{b}^+}$, $\Delta_*^{\text{b}^-}$, Δ_*^{n} and Δ_*^{v} , the $*$ meaning “generic operation”. We will not discuss these rules in detail, but each case is easy to verify. Note that for the negative destabilizer rules, we do often take the node label ϕ into account, while for other rules we cannot do this. For example, only updates to nodes matching $\$x/child :: A$ can negatively destabilize (change from true to false) a step $\$x/child :: A$, while replacing or renaming any child of $\$x$ could positively destabilize (change from false to true) this step. This is one reason it is helpful to distinguish between positive and negative Boolean destabilizers.

The rules for these functions are highly conservative because we must handle all update operations, as the previous example illustrates. Similarly, the destabilizers for irreflexive steps such as `desc` need to be made reflexive since an insertion targeting $\$x$ would affect the results. Using the `op`-destabilizer we can improve accuracy, since we are computing different queries for different classes of updates. Figure 4 shows how to do this for `insert(↓)` and `delete` operations (additional `op`-specific rules are in Appendix C). For `deletes`, we can take advantage of the fact that the result of a child step can only be affected by deleting one of the selected children, and that there is no way to positively destabilize an XPath step by deleting. Analogously, for `inserts` into the child sequence of a query, we take advantage of the fact that such updates can only destabilize a child step $\$x/child :: \phi$ by inserting into $\$x$, and that no insert can negatively destabilize such a query. We also give rules for handling `desc` steps. Similar reasoning can be applied to the other axes and other update operators, but we will restrict attention to these common cases.

Note that the algorithm can produce large expressions — in principle, there can be an exponential blowup resulting from the rewriting steps in Figure 2, due to duplication of expressions in conditionals and `for`-expressions. However, this is not a problem for the typical queries in our benchmarks.

We state correctness for `op`-destabilizers; the statement of correctness for generic destabilizers is similar.

THEOREM 2. *For any query q and any update operation $op \in \text{Ops}$:*

1. $(\Delta_{op}^{\text{b}^+}(q))_{op \in \text{Ops}}$ is a static positive `op`-destabilizer for q ;
2. $(\Delta_{op}^{\text{b}^-}(q))_{op \in \text{Ops}}$ is a static negative `op`-destabilizer for q ;
3. $(\Delta_{op}^{\text{n}}(q))_{op \in \text{Ops}}$ is a static node `op`-destabilizer for q ; and
4. $(\Delta_{op}^{\text{v}}(q))_{op \in \text{Ops}}$ is a static value `op`-destabilizer for q .

Finally, recalling examples 1–3 from the introduction, we note that the dynamic destabilizers discussed there are obtained by evaluating the static destabilizers in this section on the example data shown in Fig. 1. For example, the (simplified) subtree-value delete-

destabilizer of Q_1 is $(\$doc, \$doc/C, \$doc/C/D/\text{desc.or.self} :: *)$, which returns nodes $\{1, 7, 8\}$ when run on the document in Figure 1.

4. INDEPENDENCE ANALYSIS

In this section, we show how to use destabilizers for static independence analysis, by reducing the problem to disjointness analysis, or the problem of determining whether two selection queries can ever select a node in common.

DEFINITION 4. *Let q_1 and q_2 be selection queries. We say q_1 and q_2 overlap if there exists an input store σ such that q_1 and q_2 select a node in common when evaluated on σ . Conversely, q_1 and q_2 are disjoint if they do not overlap.*

We employ the notation $\text{Targ}_{op}(u)$ to denote the selection query that identifies all targets of atomic updates of operation op that can be produced by u . We also write $\text{Targ}_*(u)$ for the query selecting all targets of any kind of update performed by u . The (routine) definition of this function is in Appendix C. We state the main correctness result for our independence analysis for arbitrary operations; the correctness property of the generic analysis is similar:

THEOREM 3. *Query q and update expression u are statically independent*

1. *modulo \cong_{bool} if for each $op \in \text{Ops}$, the queries $\text{Targ}_{op}(u)$ and $\Delta_k^{\text{b}^+}(q), \Delta_k^{\text{b}^-}(q)$ are disjoint.*
2. *modulo \cong_{node} if for each $op \in \text{Ops}$, the queries $\text{Targ}_{op}(u)$ and $\Delta_{op}^{\text{n}}(q)$ are disjoint.*
3. *modulo \cong_{ns} if for each $op \in \text{Ops}$, the queries $\text{Targ}_{op}(u)$ and $\Delta_{op}^{\text{v}}(q)$ are disjoint.*

4.1 Disjointness Analysis

Selection queries can be translated to first-order formulas over the child, descendant and sibling relations in a standard way (see e.g. [5]). Thus, disjointness analysis for *SelXQ* reduces to satisfiability for first-order formulas over trees $FO(Tree)$, involving the child, descendant, sibling and node test predicates. However it follows from results of [29] that there is a non-elementary lower bound on the complexity of satisfiability.

We investigated a number of practical ways to solve the overlap problems arising in typical XPath or XQuery queries, including both exact algorithms and approximate heuristics that offer better performance on typical examples. Here, we present general techniques that can prove static disjointness for arbitrary pairs of selection queries. Additional heuristic techniques and further details are presented in Appendix C. We should stress here that in our implementation, we use the fast, approximate techniques first, and only if these fail try progressively stronger but more expensive methods.

Exact approaches. The most direct, and expensive, approaches we considered employ off-the-shelf techniques for checking satisfiability of tree logics. We experimented with solving these problems by translation to monadic second order logic over trees ($MSO(Tree)$), which can be solved by the `mona` tool [16, 22]. This approach can decide satisfiability for arbitrary first-order formulas over trees, thus can exactly determine whether two selection queries overlap. We also experimented with an XPath solver based on μ -calculus by Genevès et al. [17].

Reduction to SMT solving. We developed a novel approach that worked particularly well: a translation from existential first-order logic over trees ($EFO(Tree)$) to *satisfiability modulo the theory of linear order*. Concretely, we can take the linear order to be $(\mathbb{N}, <)$, a subtheory of linear arithmetic. Satisfiability modulo theories (SMT) solvers such as `yices` [14], `z3` [10], and `cvc3` [1]

have been developed that can solve such linear arithmetic satisfiability problems very quickly; we obtained the best results with the `yices` and `z3` solvers.

This approach requires that we first approximate selection queries by positive $EFO(Tree)$ formulas. We consider two such approximations: The *existential first-order abstraction* replaces conditionals `if q then q1 else q2` by the monotone overapproximation (`if q then q1`), `q2` and then applies the translation to $FO(Tree)$ from [5]. The *simple path abstraction* approximates a query by overapproximating it with a set of simple path expressions, that is, plain sequences of XPath axis steps. Both of these abstractions are defined and verified in Appendix C.

Once we have approximated the destabilizer and target queries by $EFO(Tree)$ formulas, we translate the satisfiability question to one concerning sentences about a linear order. The translation is based on the well-known interval encoding [11] of XML documents, widely-used for indexing. We represent a tree node by a triple $(pre, post, tag)$ of numbers, representing the starting position, ending position, and tag of x respectively. We constrain each node x to satisfy $0 < x.pre \leq x.post$ and each pair of nodes x, y must satisfy:

$$\begin{aligned} \text{(well nesting)} \quad & x.pre < x.post < y.pre < y.post \\ \text{(uniqueness)} \quad & x.pre = y.pre \iff x.post = y.post \\ \text{(tags unique)} \quad & x.pre = y.pre \implies x.tag = y.tag \end{aligned}$$

Then, given a formula ϕ over *child*, *desc* and *sibling*, whose free variables are drawn from a finite set X , we can define:

$$\begin{aligned} \text{follows}(x, y) &= x.post < y.pre \\ \text{desc}(x, y) &= x.pre < y.pre \wedge y.post < x.post \\ \text{btwn}(x, y, z) &= \text{desc}(x, y) \wedge \text{desc}(y, z) \\ \text{child}(x, y) &= \text{desc}(x, y) \wedge \bigwedge_{z \in X} \neg(\text{btwn}(x, z, y)) \\ \text{sibling}(x, y) &= \exists z. \text{child}(z, x) \wedge \text{child}(z, y) \wedge \text{follows}(x, y) \end{aligned}$$

Node tests are represented using constraints $x.tag = i$, where i is an integer index surrogate for the actual element tag. We constrain text nodes to be leaves, i.e. $x.pre = x.post$, and to have a special tag not used for any element tag, such as -1 . Note that in the definition of $\text{child}(x, y)$, we check that none of the nodes described by the variables of the formula ϕ are between x and y . This is correct because it suffices to restrict attention to these nodes when considering satisfiability (we omit a full proof but the details are similar to the satisfiability proofs in [21]).

Additional heuristics. Calling an external SMT solver is expensive, and often we can avoid this cost using heuristics to solve easy common cases of disjointness problems. For downward-only paths, overlap is exactly decidable in quadratic time [20], and we use this algorithm as well as other heuristics, described further in Appendix C.

5. EXPERIMENTS

5.1 Experimental Setup

We have implemented a static independence analyzer in OCaml based on destabilizers. The analysis takes a query and update expression and attempts to determine their independence, first using heuristics, and then optionally by calling an external solver.

We evaluated several analyses of varying degrees of complexity:

1. L_1 : Solving via simple path approximations, using only heuristics and an exact algorithm for downward-only paths [20].
2. L_2 : Solving via simple path approximations using `yices`.
3. L_3 : Solving XQuery disjointness directly using `mona`. We permit `mona` at most one second to find a solution.

4. SCH : For comparison, using the schema-based analysis from our earlier work [2].

Each of the above levels L_i is cumulative, and stops as soon as a definite answer is obtained: e.g. L_2 involves first trying the heuristics and Hammerschmidt et al.’s algorithm [20], then trying `yices`. We also evaluated several other translations and external solvers in place of `yices`, in step L_2 .

All external solvers were called as separate processes, and each such call is expensive; this represents an opportunity for improvement. Some of the solvers (`yices`, `mona`, `cvc3`) are available as libraries or in source code form, but others (`z3`, Genevès et al.’s XPath solver [17]) are not. We therefore chose not to try to integrate our implementation more closely with any external solvers in order to permit fair comparisons among the different solvers.

We used queries and updates derived from two standard sets of benchmark queries: twenty XQuery queries comprising the XMark benchmark [28], and sixteen queries from the XPathMark benchmark [15]. Our simplified forms of these queries exercise many aspects of XQuery and were used to evaluate previous work on schema-based independence analysis (e.g. [2]). The XMark queries were all translated to the more-restrictive core language accepted by our system, and we automatically generated rename, delete and insert updates based on each XPathMark query. We refer to these as UA_1 – UA_8 and UB_1 – UB_8 . See Appendix E for further details of the queries and updates. Experiments were conducted on an Intel Pentium D (3.0 Ghz) running Ubuntu Linux 8.10, with 2GB of main memory.

5.2 Experimental Results

Running Time. For each update, we measured the total time taken to analyze its independence with respect to all queries. For each pair, we recorded system time elapsed between the start of execution and the time when the independence result is produced, including all time taken by external solvers and system overhead. We report times for the op-insensitive analysis only; the running times for op-sensitive analysis were similar. Figure 5 shows the total time taken by each analysis for each update to analyze all 36 benchmark queries. We also show the time for the schema-based analysis SCH from [2].

For comparison, we also evaluated the performance of our approach using Genevès et al.’s XPath solver, and using `z3` and `cvc3` (both via the interval encoding). The running times obtained using `z3` were similar to those using `yices`, while `cvc3` and Genevès et al.’s solver were not competitive. (We used a preliminary, unoptimized version of Genevès et al.’s solver; a more efficient version is under development but has not yet been released.)

Analysis Results. We summarize the results of the analysis by giving the percentages of queries that each level of the analysis was able to prove independent of each update, in Figure 6. These results are for the generic version of the analysis, so the operation of the update (insert, delete or rename) was irrelevant; the results are grouped by update numbers UA_i and UB_i . We also report the results of the schema-based analysis SCH and the combination of L_3 and SCH (labeled $L_3 + SCH$).

We also compare the generic analysis with the op-sensitive analysis. These results are summarized in Figure 7, where we show the total number of problems proved independent by the schema-based analysis SCH , destabilizer-based analyses L_1 – L_3 , and combined analyses $L_3 + SCH$.

Maintenance time. Finally, to evaluate the overall effectiveness of the approach, we measured the savings in view recomputation time for analyses SCH , L_1 – L_3 , and $L_2 + SCH$. We summarize the results in Table 1 by presenting the total analysis overhead time and

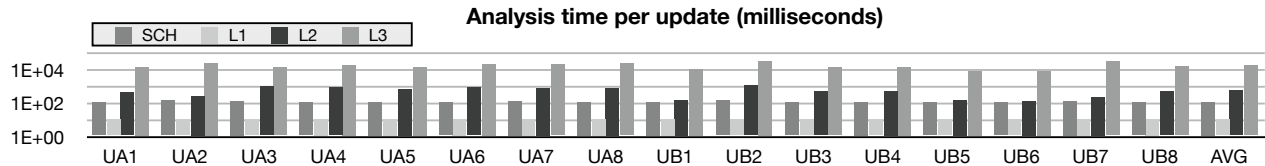


Figure 5: Running times for the generic analysis, in milliseconds (logarithmic scale), broken down by update and analysis level.

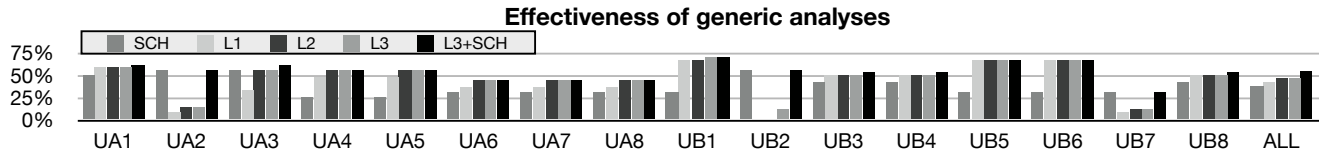


Figure 6: Effectiveness of the generic analysis, expressed as a percentage of query-update pairs determined independent, broken down by update and by analysis level.

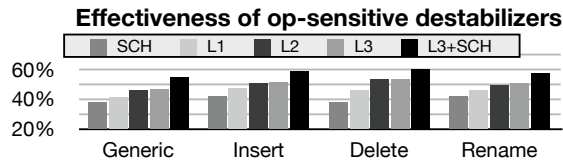


Figure 7: Effectiveness of the generic and op-sensitive analysis, expressed as a percentage of query-update pairs determined independent, by analysis level and update operation.

	<i>SCH</i>	<i>L</i> ₁	<i>L</i> ₂	<i>L</i> ₃	<i>L</i> ₂ + <i>SCH</i>
Static overhead	1.8	0.16	8.4	265	10
1.1MB	52	53	55	55	66
2.3MB	95	80	86	85	130
5.7MB	375	312	311	311	555

Table 1: Static analysis overheads and maintenance savings across whole benchmark (seconds).

total time saved for each technique on documents of increasing size. The documents are 1.1MB, 2.3MB and 5.7MB XMark data [28]. We used BaseX 6.1, an efficient XQuery engine that implements the full W3C update recommendation.

5.3 Discussion

Efficiency. The running times reported in Figure 5 show that the simple techniques used by *L*₁ are around ten times faster than *SCH*; the difference is likely due to the need for processing and analyzing the schema. The more sophisticated analysis level *L*₂ is generally slower than *SCH*, and *L*₃ is generally much slower, even with a 1 second time limit.

Effectiveness. Using the op-insensitive analysis, *L*₁ is able to determine independence in approximately 41% of cases. *L*₂ improves this to around 46%. However, the most sophisticated analysis, *L*₃, is only able to find a handful of additional independent pairs. Using op-sensitive destabilizers led to improvements of 4–6% for each level of the analysis. This demonstrates the value of taking the update operators into account.

We cannot easily measure the completeness of this approach, and

do not have a theoretical characterization of its (relative) completeness. Informal inspection indicates that op-sensitive analysis is reasonably precise. For example, we manually constructed counterexamples to independence for all but five of the deletion-based problems that are not proved independent by op-sensitive *L*₃. Automating this heuristic counterexample-generation process could be an interesting topic for future work.

Figure 6 also gives an idea of how independence of the benchmark update/query pairs depends on a schema. Although our approach is good at detecting schema-free independence, there are certain pairs that can interfere in the absence of a schema but are independent for the XMark schema — e.g. *UA*₂, *UB*₂, *UB*₇ involve the descendant or ancestor axes and can modify data almost anywhere in the tree, so actually do interfere with most queries in the absence of a schema – thus schema-based analysis does better on these updates. Note that our approach can certainly accommodate schemas, simply by using schema-aware satisfiability tests – this is another opportunity for future work.

Impact on View Maintenance. Table 1 shows that all of the techniques save time overall on our benchmark compared to full recomputation, and that the new techniques in this paper find savings that previous approach *SCH* did not. Moreover, as document size increases, the savings-to-overhead ratio increases; already for a 5.7MB document the cost of analysis is minuscule compared to the time savings obtained. This is promising, but some important caveats apply to these results: we are comparing with full recomputation, not incremental view maintenance, and our benchmark is synthetic. However, dynamic view maintenance techniques do not yet appear to be available in standard XML databases. More work needs to be done to develop realistic view maintenance benchmarks that can be used to evaluate the techniques.

6. RELATED WORK

Due to space limits, we discuss only directly related work. Additional related work (including work on incremental view maintenance for XML) is reviewed in Appendix D.

This work is inspired partly by work on XML projection, where the goal is to identify nodes that can safely be deleted without affecting the result of a query [24, 6]; this is similar to independence problems involving deletion only.

Ghelli et al. [18] developed a static commutativity analysis, using paths to represent the sets of nodes accessed and updated by an ex-

pression. However, their analysis is based on a different XML update language proposal and does not address independence. Raghavachari and Shmueli [26] investigated the complexity of the independence problem for different fragments of downward XPath, while Sawires et al. [27] implemented an independence checker for downward XPath queries and updates, motivated by difficulties with incremental view maintenance in loosely-coupled systems. Benedikt and Cheney [2] developed a schema-based independence analysis for core XQuery over arbitrary XPath axes.

Our approach improves on these techniques in several ways. We handle arbitrary XPath axes and all of the essential features of XQuery. By analyzing operations we gain accuracy over any purely node-based approach, by using a variety of XQuery overlap tests we can gain accuracy on path-based abstraction (as used in [18, 26, 27, 24, 6]) and we do not rely on the presence of a schema (unlike [2]). Moreover, we have evaluated our approach on benchmark queries and updates exercising many features of XPath/XQuery.

Beyond these advantages, the notion of destabilizer itself is a key conceptual contribution of our work, since it factors the analysis problem cleanly into the approximation of the runtime updates by a query and an overlap test. This cleanly generalizes ideas in previous work. It is conceptually straightforward to extend our approach to handle constraints or schemas, simply by testing disjointness with respect to such constraints. Destabilizers may be useful in other data models, and may also be applicable to the dynamic, incremental maintenance problem.

7. CONCLUSIONS

We have introduced the notion of destabilizer, or a representation of the set of updates that may change the result of a given query. The problem of determining whether a query and update are (statically) independent reduces to calculating a destabilizer and testing whether the destabilizer and update are (statically) disjoint. This yields a general, modular and extensible framework for reasoning about independence.

We defined generic and kind-sensitive destabilizers via query rewriting for XQuery queries and updates. We implemented an independence analysis based on these destabilizers, and demonstrated its effectiveness and efficiency. Our approach offers a range of tradeoffs between precision and analysis cost, and can solve independence problems that no prior approach could solve. Even for fairly small documents, the overhead of our approach is much smaller than the savings it obtains. Combining this approach with schema-based analysis yields further savings. Moreover, these savings increase with document size.

There are many interesting directions for future work, such as schema-conscious disjointness analysis, counterexample generation, and combining static and dynamic view maintenance techniques.

Acknowledgments. Benedikt is supported in part by EPSRC EP/G004021/1 (the Engineering and Physical Sciences Research Council, UK) and by FET FP7-ICT-233599 (European Research Consortium). Cheney is supported by a Royal Society University Research Fellowship.

8. REFERENCES

- [1] C. Barrett and C. Tinelli. CVC3. In *CAV 2007*, July 2007.
- [2] M. Benedikt and J. Cheney. Schema-based independence analysis for XML updates. In *VLDB*, 2009.
- [3] M. Benedikt and J. Cheney. Semantics, types and effects for XML Updates. In *DBPL*, 2009.
- [4] M. Benedikt, W. Fan, and F. Geerts. XPath satisfiability in the presence of DTDs. In *PODS*, 2005.
- [5] Michael Benedikt and Christoph Koch. From XQuery to relational logics. *ACM TODS*, 34(4):A25, 2009.
- [6] V. Benzaken, G. Castagna, D. Colazzo, and K. Nguyễn. Type-based XML projection. In *VLDB*, 2006.
- [7] H. Björklund, W. Gelade, M. Marquardt, and W. Martens. Incremental XPath evaluation. In *ICDT*, 2009.
- [8] J. Blakeley, N. Coburn, and P.-Å. Larson. Updating Derived Relations: Detecting Irrelevant and Autonomously Computable Updates. *ACM TODS*, 14(3):369–400, 1989.
- [9] D. Chamberlin, M. Dyck, D. Florescu, J. Melton, J. Robie, and J. Siméon. XQuery update facility 1.0. W3C Candidate Recommendation, August 2008.
- [10] L. Mendonça de Moura and N. Bjørner. Z3: An efficient SMT solver. In *TACAS*, 2008.
- [11] D. DeHaan, D. Toman, M. Consens, and M. T. Özsu. A comprehensive XQuery to SQL translation using dynamic interval encoding. In *SIGMOD*, 2003.
- [12] Katica Dimitrova, Maged El-Sayed, and Elke A. Rundensteiner. Order-Sensitive View Maintenance of Materialized XQuery Views. In *ER*, 2003.
- [13] D. Draper, P. Fankhauser, M. Fernández, A. Malhotra, K. Rose, M. Rys, J. Siméon, and P. Wadler. XQuery 1.0 and XPath 2.0 formal semantics. W3C Recommendation, January 2007.
- [14] B. Dutertre and L. de Moura. The Yices SMT solver. Tool paper at <http://yices.csl.sri.com/tool-paper.pdf>, August 2006.
- [15] M. Franceschet. XPathMark: an XPath benchmark for XMark generated data. In *XXSYM*, 2005.
- [16] P. Genevès and N. Layaïda. Deciding XPath containment with MSO. *Data Knowl. Eng.*, 63(1):108–136, 2007.
- [17] P. Genevès, N. Layaïda, and A. Schmitt. Efficient static analysis of XML paths and types. In *PLDI*, 2007.
- [18] G. Ghelli, K. Rose, and J. Siméon. Commutativity analysis for XML updates. *ACM TODS*, 33(4):1–47, 2008.
- [19] T. Griffin, L. Libkin, and H. Trickey. An improved algorithm for the incremental recomputation of active relational expressions. *IEEE TKDE*, 9(3):508–511, 1997.
- [20] B. Hammerschmidt, M. Kempa, and V. Linnemann. On the Intersection of XPath Expressions. In *IDEAS*, 2005.
- [21] J. Hidders. Satisfiability of XPath expressions. In *DBPL*, 2003.
- [22] N. Klarlund and A. Møller. Mona v. 1.4 user manual. Technical Report BRICS NS-01-1, U. Aarhus, 2001.
- [23] A. Levy and Y. Sagiv. Queries Independent of Updates. In *VLDB*, 1993.
- [24] A. Marian and J. Siméon. Projecting XML documents. In *VLDB*, 2003.
- [25] X. Qian and G. Wiederhold. Incremental Recomputation of Active Relational Expressions. *IEEE TKDE*, 3(3):337–341, 1991.
- [26] M. Raghavachari and O. Shmueli. Conflicting XML updates. In *EDBT*, 2006.
- [27] A. Sawires, J. Tatemura, O. Po, D. Agrawal, A. El Abbadi, and K. Candan. Maintaining XPath Views In Loosely Coupled Systems. In *VLDB*, 2006.
- [28] A. Schmidt, F. Waas, M. Kersten, M. Carey, I. Manolescu, and R. Busse. XMark: A Benchmark for XML Data Management. In *VLDB*, 2002.
- [29] Sergei G. Vorobyov. An improved lower bound for the elementary theories of trees. In *CADE*, 1996.

$$\begin{array}{c}
\frac{\sigma(1) = \text{text}[s] \quad 1' \notin \text{dom}(\sigma)}{\sigma, 1 \xrightarrow{\text{copy}} \sigma[1' := \text{text}[s]], 1'} \\
\frac{\sigma(1) = a[L] \quad \sigma, L \xrightarrow{\text{copy}} \sigma', L' \quad 1' \notin \text{dom}(\sigma')}{\sigma, 1 \xrightarrow{\text{copy}} \sigma'[1' := a[L']], 1'} \\
\frac{\sigma, () \xrightarrow{\text{copy}} \sigma, ()}{\sigma, L_1 \xrightarrow{\text{copy}} \sigma', L'_1 \quad \sigma', L_2 \xrightarrow{\text{copy}} \sigma'', L'_2} \\
\frac{\sigma, L_1 \cdot L_2 \xrightarrow{\text{copy}} \sigma'', L'_1 \cdot L'_2}{}
\end{array}$$

Figure 8: Copying rules

APPENDIX

A. SEMANTICS OF QUERIES AND UPDATES

A.1 Query semantics

A store is modeled as a mapping $\sigma : \text{Loc} \rightarrow \Sigma \times \text{Loc}^* \uplus \text{String}$, where Σ is the set of element tags; we often write $\sigma(1) = a[L]$ or $\text{text}[s]$, instead of $\sigma(1) = (a, L)$ or s , respectively.

An environment is modeled as a mapping $\gamma : \text{Var} \rightarrow \text{Loc}^*$ from variables to lists of nodes.

We define the semantics of queries using the following judgments:

- $\sigma, L \xrightarrow{\text{copy}} \sigma', L'$ — takes a store and a list of nodes in the store and copies the subtrees below the nodes, yielding a new store and a new list of nodes pointing to the roots of the new trees.
- $\sigma, \gamma \models q \Rightarrow \sigma', L$ — takes a store and environment and evaluates a query to a new store and a list of nodes.
- $\sigma, \gamma \models q \xrightarrow{\text{copy}} \sigma', L$ — an abbreviation for evaluating q and then making a copy of the result.
- $\sigma, \gamma, x \in L \models^* q \Rightarrow \sigma', L$ — takes a store, environment, variable x , and list of nodes L , and evaluates the query with x bound to each node in L , concatenating the results in order.
- $\sigma \models L/ax::\phi \xrightarrow{\text{step}} L'$ — given a store and list of nodes, evaluates XPath axis step $ax :: \phi$ to produce a new list of nodes. This judgment is standard and we omit its definition.

A.2 Update semantics

Atomic updates. We consider atomic updates of the form:

$$\begin{array}{l}
\iota ::= \text{ins}(L, d, 1) \mid \text{del}(1) \mid \text{repl}(1, L) \mid \text{ren}(1, a) \\
d ::= \leftarrow \mid \rightarrow \mid \downarrow \mid \sphericalangle \mid \searrow
\end{array}$$

Here, the direction d indicates whether to insert before (\leftarrow), after (\rightarrow), or into the child list in first (\sphericalangle), last (\searrow) or arbitrary position (\downarrow). We define the *target* of an atomic update ι to be the distinguished node L ; for example $\text{Targ}(\text{ins}(L, d, 1)) = L$.

Our basic update model will be *sequences of atomic updates*. We will let ω range over sequences with the empty sequence written ϵ and concatenation written $\omega; \omega'$. Formally, we can define the semantics of atomic updates as a relation $\sigma \models \iota \rightsquigarrow \sigma'$; for the precise definition, see Figure 10. We will refer to this relation as EvalUpd . Our semantics is a formalization of that given by the W3C's candidate recommendation for updates, see [9].

Updating expressions. The W3C XQuery Update proposal defines high-level *updating expressions* that can be used to generate atomic updates that are to be performed [9]. We use the following

$$\begin{array}{c}
\frac{}{\sigma, \gamma \models \$x \Rightarrow \sigma, \gamma(x)} \\
\frac{}{\sigma, \gamma \models () \Rightarrow \sigma, ()} \\
\frac{\sigma, \gamma \models q_1 \Rightarrow \sigma_2, L_1 \quad \sigma_2, \gamma \models q_2 \Rightarrow \sigma_3, L_2}{\sigma, \gamma \models q_1, q_2 \Rightarrow \sigma_3, L_1 \cdot L_2} \\
\frac{\sigma, \gamma \models q \xrightarrow{\text{copy}} \sigma_2, L \quad 1 \notin \text{dom}(\sigma_2)}{\sigma, \gamma \models \langle a \rangle q \langle /a \rangle \Rightarrow \sigma_2[1 := a[L]], 1} \\
\frac{\sigma, \gamma \models q \Rightarrow \sigma_2, 1 \cdot L \quad \sigma_2, \gamma \models q_1 \Rightarrow \sigma_3, L_1}{\sigma, \gamma \models \text{if } q \text{ then } q_1 \text{ else } q_2 \Rightarrow \sigma_3, L_1} \\
\frac{\sigma, \gamma \models q \Rightarrow \sigma_2, () \quad \sigma_2, \gamma \models q_2 \Rightarrow \sigma_3, L_2}{\sigma, \gamma \models \text{if } q \text{ then } q_1 \text{ else } q_2 \Rightarrow \sigma_3, L_2} \\
\frac{\sigma, \gamma \models q_1 \Rightarrow \sigma_2, L \quad \sigma_2, \gamma, x \in L \models^* q_2 \Rightarrow \sigma_3, L'}{\sigma, \gamma \models \text{for } \$x \in q_1 \text{ return } q_2 \Rightarrow \sigma_3, L'} \\
\frac{\sigma \models \gamma(x)/ax::\phi \xrightarrow{\text{step}} L}{\sigma, \gamma \models x/ax :: \phi \Rightarrow \sigma, L} \\
\frac{\sigma, \gamma \models q \Rightarrow \sigma_0, L_0 \quad \sigma_0, L_0 \xrightarrow{\text{copy}} \sigma', L}{\sigma, \gamma \models q \xrightarrow{\text{copy}} \sigma', L} \\
\frac{}{\sigma, \gamma, x \in () \models^* q \Rightarrow \sigma, ()} \\
\frac{\sigma, \gamma[x := 1] \models q \Rightarrow \sigma_2, L_1 \quad \sigma, \gamma, x \in L \models^* q \Rightarrow \sigma_3, L_2}{\sigma, \gamma, x \in 1 \cdot L \models^* q \Rightarrow \sigma_3, L_1 \cdot L_2}
\end{array}$$

Figure 9: Query evaluation rules

core language for update expressions:

$$\begin{array}{l}
u ::= () \mid u, u' \mid \text{insert } q \text{ d } q_0 \mid \text{delete } q_0 \\
\mid \text{rename } q_0 \text{ as } a \mid \text{replace } q_0 \text{ with } q \\
\mid \text{if } q \text{ then } u_1 \text{ else } u_2 \mid \text{for } \$x \in q \text{ return } u
\end{array}$$

The XQuery Update proposal re-uses existing query syntax for updates. The $()$ expression does nothing, while u, u' is sequential composition, and XQuery conditionals and for-loops can also be used. There are four atomic update expressions: insertion $\text{insert } q \text{ d } q_0$, which says to insert a copy of the value of q in position d relative to node q_0 ; deletion $\text{delete } q_0$, which says to delete the node q_0 , disconnecting it from its parent; renaming $\text{rename } q_0 \text{ as } a$, which says to rename the node q_0 to a and replacement $\text{replace } q_0 \text{ with } q$, which says to replace the value of node q_0 with a copy of q . In each case, the target expression q_0 is required to be a selection query and expected to evaluate to a single node; if not, evaluation fails. Again, we omit let -binding without loss of generality.

Semantics of updates. Updates have a multi-phase semantics. First, the updating expression is *evaluated*, resulting in a *pending update list* ω . We model this phase using an update evaluation judgement $\sigma, \gamma \models u \Rightarrow \sigma', \omega$, defined in Figure 11 (and based on a semantics developed elsewhere [3]). Note that, as for queries, the store may grow as a result of allocation, but trees already in σ do *not* change while the update expression is being evaluated. Next, ω is *validated* to avoid pathological problems such as renaming a node in two different ways. We do not model this validation phase explicitly here. Finally, the pending updates are *applied* to the store in some order. The W3C semantics mandates reordering the update list so that insertions and renamings are applied first, then replacements, and finally deletions. We will conservatively assume that atomic updates in the pending update list might be performed in *any* order. A static analysis that is sound with respect to this se-

$$\begin{array}{c}
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{ins}(L, \leftarrow, l) \rightsquigarrow \sigma[l' := a[L_1 \cdot L \cdot l \cdot L_2]]} \\
\frac{\sigma(l) = a[L']}{\sigma \models \text{ins}(L, \sphericalangle, l) \rightsquigarrow \sigma[l := a[L \cdot L']]} \\
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{ins}(L, \rightarrow, l) \rightsquigarrow \sigma[l' := a[L_1 \cdot l \cdot L \cdot L_2]]} \\
\frac{\sigma(l) = a[L']}{\sigma \models \text{ins}(L, \searrow, l) \rightsquigarrow \sigma[l := a[L' \cdot L]]} \\
\frac{\sigma(l) = a[L_1 \cdot L_2]}{\sigma \models \text{ins}(L, \downarrow, l) \rightsquigarrow \sigma[l := a[L_1 \cdot L \cdot L_2]]} \\
\frac{\sigma(l) = a[L]}{\sigma \models \text{ren}(l, b) \rightsquigarrow \sigma[l := b[L]]} \\
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{repl}(l, L) \rightsquigarrow \sigma[l' := a[L_1 \cdot L \cdot L_2]]} \\
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{del}(l) \rightsquigarrow \sigma[l' := a[L_1 \cdot L_2]]} \\
\frac{\sigma \models \epsilon \rightsquigarrow \sigma}{\sigma \models \omega_1 \rightsquigarrow \sigma' \quad \sigma' \models \omega_2 \rightsquigarrow \sigma''} \\
\frac{\sigma \models \omega_1, \omega_2 \rightsquigarrow \sigma''}{\sigma \models \omega_1, \omega_2 \rightsquigarrow \sigma''} \\
\frac{\sigma, \gamma \models u \Rightarrow \sigma', \omega \quad \omega' \text{ any permutation of } \omega \quad \sigma' \models \omega' \rightsquigarrow \sigma''}{\sigma, \gamma \models u \rightsquigarrow \sigma''}
\end{array}$$

Figure 10: Update application

mantics will also be sound with respect to any implementation of the W3C semantics.

To summarize, the semantics of updates is defined using the following judgments in Figure 10 and Figure 11.

- $\sigma \models \iota \rightsquigarrow \sigma'$ — given a store σ and an atomic update, produce a store resulting from applying ι to σ . (This, and all later judgments, can be nondeterministic because of node copying and because insert-into operations are nondeterministic.)
- $\sigma \models \omega \rightsquigarrow \sigma'$ — given a store σ and an atomic update sequence, perform the atomic updates in order.
- $\sigma, \gamma \models u \Rightarrow \sigma', \omega$ — given a store σ and variable environment γ , evaluate u to a new store and atomic update list ω .
- $\sigma, \gamma, x \in L \models^* u \Rightarrow \sigma', \omega$ — given a store σ , variable environment γ , variable x and node list L , evaluate update u with x bound to each element of L in turn, concatenating the resulting update sequences in order.
- $\sigma, \gamma \models u \rightsquigarrow \sigma'$ — given input store and variable environment γ , evaluate u and apply the resulting atomic updates in any order, yielding σ' .

B. PROOFS

B.1 Proof of Theorem 1

Recall the statement of the result: There is no elementary algorithm for constructing a pointwise minimal static destabilizer.

For the first part, recall that the satisfiability problem for first-order logic on labeled trees can not be resolved in elementary time [29]. From the results of [5] it follows that the satisfiability problem for *SelXQ* is non-elementary. Now consider the document Σ_0

$$\begin{array}{c}
\frac{\sigma, \gamma \models () \Rightarrow \sigma, \epsilon}{\sigma_1, \gamma \models u_1 \Rightarrow \sigma_2, \omega_1 \quad \sigma_2, \gamma \models u_2 \Rightarrow \sigma_3, \omega_2} \\
\frac{\sigma_1, \gamma \models u_1, u_2 \Rightarrow \sigma_3, \omega_1; \omega_2}{\sigma_1, \gamma \models q \Rightarrow \sigma_2, l \cdot L \quad \sigma_2, \gamma \models u_1 \Rightarrow \sigma_3, \omega_1} \\
\frac{\sigma_1, \gamma \models q \Rightarrow \sigma_2, () \quad \sigma_2, \gamma \models u_2 \Rightarrow \sigma_3, \omega_2}{\sigma_1, \gamma \models \text{if } q \text{ then } u_1 \text{ else } u_2 \Rightarrow \sigma_3, \omega_1} \\
\frac{\sigma_1, \gamma \models q \Rightarrow \sigma_2, () \quad \sigma_2, \gamma \models u_2 \Rightarrow \sigma_3, \omega_2}{\sigma_1, \gamma \models \text{if } q \text{ then } u_1 \text{ else } u_2 \Rightarrow \sigma_3, \omega_2} \\
\frac{\sigma_1, \gamma \models q \Rightarrow L, \sigma_2 \quad \sigma, \gamma, x \in L \models^* u \Rightarrow \sigma_3, \omega}{\sigma_1, \gamma \models \text{for } \$x \in q \text{ return } u \Rightarrow \sigma_3, \omega} \\
\frac{\sigma_1, \gamma \models q_1 \xrightarrow{\text{copy}} \sigma_2, L_1 \quad \sigma_2, \gamma \models q_2 \Rightarrow \sigma_3, l_2}{\sigma_1, \gamma \models \text{insert } q_1 \text{ d } q_2 \Rightarrow \sigma_3, \text{ins}(L_1, d, l_2)} \\
\frac{\sigma_1, \gamma \models q \Rightarrow \sigma_2, l}{\sigma_1, \gamma \models \text{delete } q \Rightarrow \sigma_2, \text{del}(l)} \\
\frac{\sigma_1, \gamma \models q_1 \Rightarrow \sigma_2, l_1 \quad \sigma_2, \gamma \models q_2 \xrightarrow{\text{copy}} \sigma_3, L_2}{\sigma_1, \gamma \models \text{replace } q_1 \text{ with } q_2 \Rightarrow \sigma_3, \text{repl}(l_1, L_2)} \\
\frac{\sigma_1, \gamma \models q \Rightarrow \sigma_2, l}{\sigma_1, \gamma \models \text{rename } q \text{ as } a \Rightarrow \sigma_2, \text{ren}(l, a)} \\
\frac{\sigma, \gamma, x \in () \models^* u \Rightarrow \sigma, \epsilon}{\sigma_1, \gamma[x := l] \models u \Rightarrow \sigma_2, \omega_1 \quad \sigma_2, \gamma, x \in L \models^* u \Rightarrow \sigma_3, \omega_2} \\
\frac{\sigma_1, \gamma, x \in l \cdot L \models^* u \Rightarrow \sigma_3, \omega_1; \omega_2}{\sigma_1, \gamma, x \in l \cdot L \models^* u \Rightarrow \sigma_3, \omega_1; \omega_2}
\end{array}$$

Figure 11: Update expression evaluation

consisting only of a single node, and a boolean query Q . If Q is unsatisfiable, then the empty set is a minimal set of nodes of Σ_0 having the property that every update sequence making Q turn from false to true on Σ_0 contains a node in the set. If Q is satisfiable but does not hold on Σ_0 , then the empty set is not such a set of nodes, since we could replace the root of Σ_0 by some document satisfying Q . Hence we can determine whether Q is satisfiable by first running Q on d_0 and then running an exact destabilizer for Q on d_0 and checking if the result is empty. Since evaluation of a *SelXQ* query is polynomial time (e.g. [5]) this would provide a contradiction. \square

C. DETAILS OF INDEPENDENCE ANALYSIS

C.1 Additional destabilizer rules

In Figure 12 we show additional rules for rename, insert before, and insert after destabilizers.

For renames, note that a rename operation can only affect the label of an element node, not the tree structure. Hence, for positive boolean and node destabilizers, rename operations that do not involve element label tests have empty destabilizers, no matter what axis is used. Furthermore, for negative boolean destabilizers, the only way to destabilize a step involving a node test $ax :: a$ is to rename a node already labeled a to something else, so in this case we can improve precision slightly.

The insert-before and insert-after axes are symmetric; we discuss only the former case. For positive boolean destabilizers, we can never make a $\text{child} :: *$ step nonempty by inserting before a node, since since if a node has no children we cannot add a child by inserting before a node, while if a node has children already then inserting before one of them does not change the Boolean value of

$$\begin{aligned}
\Delta_{\text{rename}}^{\text{b}+}(\$x/\text{ax} :: a) &= \$x/\text{ax} :: * \\
\Delta_{\text{rename}}^{\text{b}+}(\$x/\text{ax} :: \phi) &= () \quad (\phi \in \{*, \text{text}()\}) \\
\Delta_{\text{insert}(\leftarrow)}^{\text{b}+}(\$x/\text{child} :: *) &= () \\
\Delta_{\text{insert}(\leftarrow)}^{\text{b}+}(\$x/\text{child} :: \phi) &= \$x/\text{child} :: * \quad (\phi \neq *) \\
\Delta_{\text{insert}(\leftarrow)}^{\text{b}+}(\$x/\text{foll_sib} :: \phi) &= (\$x/\text{foll_sib} :: *) \\
\Delta_{\text{insert}(\leftarrow)}^{\text{b}+}(\$x/\text{prec_sib} :: \phi) &= (\$x, \$x/\text{prec_sib} :: *) \\
\Delta_{\text{insert}(\rightarrow)}^{\text{b}+}(\$x/\text{child} :: *) &= () \\
\Delta_{\text{insert}(\rightarrow)}^{\text{b}+}(\$x/\text{child} :: \phi) &= \$x/\text{child} :: * \quad (\phi \neq *) \\
\Delta_{\text{insert}(\rightarrow)}^{\text{b}+}(\$x/\text{foll_sib} :: \phi) &= (\$x, \$x/\text{foll_sib} :: *) \\
\Delta_{\text{insert}(\rightarrow)}^{\text{b}+}(\$x/\text{prec_sib} :: \phi) &= (\$x/\text{prec_sib} :: *) \\
\Delta_{\text{rename}}^{\text{b}-}(\$x/\text{ax} :: a) &= \$x/\text{ax} :: a \\
\Delta_{\text{rename}}^{\text{b}-}(\$x/\text{ax} :: \phi) &= () \quad (\phi \in \{*, \text{text}()\}) \\
\Delta_{\text{insert}(\leftarrow)}^{\text{b}-}(\$x/\text{ax} :: \phi) &= () \\
\Delta_{\text{insert}(\rightarrow)}^{\text{b}-}(\$x/\text{ax} :: \phi) &= () \\
\Delta_{\text{rename}}^{\text{n}}(\$x/\text{ax} :: a) &= \$x/\text{ax} :: * \\
\Delta_{\text{rename}}^{\text{n}}(\$x/\text{ax} :: \phi) &= () \quad (\phi \in \{*, \text{text}()\}) \\
\Delta_{\text{insert}(\leftarrow)}^{\text{n}}(\$x/\text{child} :: \phi) &= \$x/\text{child} :: * \\
\Delta_{\text{insert}(\leftarrow)}^{\text{n}}(\$x/\text{foll_sib} :: \phi) &= (\$x/\text{foll_sib} :: *) \\
\Delta_{\text{insert}(\leftarrow)}^{\text{n}}(\$x/\text{prec_sib} :: \phi) &= (\$x, \$x/\text{prec_sib} :: *) \\
\Delta_{\text{insert}(\rightarrow)}^{\text{n}}(\$x/\text{child} :: \phi) &= \$x/\text{child} :: * \\
\Delta_{\text{insert}(\rightarrow)}^{\text{n}}(\$x/\text{foll_sib} :: \phi) &= (\$x, \$x/\text{foll_sib} :: *) \\
\Delta_{\text{insert}(\rightarrow)}^{\text{n}}(\$x/\text{prec_sib} :: \phi) &= (\$x/\text{prec_sib} :: *)
\end{aligned}$$

Figure 12: Additional op-sensitive destabilizer rules

this step. Of course, for a node-test destabilizer $\text{child} :: \phi$ where ϕ is not $*$, we can make the result go from empty to nonempty by inserting a node matching ϕ before an already-present node. For negative destabilizers, there is no way an insert-before operation can make the result of a child step empty. Finally, for node destabilizers, we conservatively assume that any insert-before to a child list of a node can destabilize the node list returned by the child step. Other axes such as descendant can be handled similarly to the child axis. For the preceding and following axes, however, we can leverage the order behavior of the insert-before operation. Specifically, for the following sibling axis, only inserts before some following sibling can affect the result of this step, while inserts before both the node and its preceding siblings can affect the preceding-sibling step.

C.2 Update target queries

We define update target queries $\text{Targ}_{op}(u)$ as shown in Figure 13, where $op \in \text{Ops}$. Note that $\text{Targ}_{op}(u)$ is a selection query. Its key correctness property is:

LEMMA 1. *Suppose σ is a store and $(\sigma', \omega) \in \text{Upd}[\mathbb{U}]\sigma$. Then for any $\iota \in \omega$, if $\iota : op$ and $\text{Targ}(\iota) \in \sigma$ then*

$$\text{Targ}(\iota) \in \text{Sel}[\text{Targ}_{op}(u)]\sigma.$$

C.3 Query translations

We employ three translations from selection queries to other representations:

- The *first-order abstraction* translates a selection query to a first-order formula over the basic axis predicates.
- The *existential first-order abstraction* translates a selection query to a *positive* existential first-order formula (that is, one with no negation or universal quantifiers).

$$\begin{aligned}
\text{Targ}_{op}() &= () \\
\text{Targ}_{op}(u_1, u_2) &= \text{Targ}_{op}(u_1), \text{Targ}_{op}(u_2) \\
\text{Targ}_{op}(\text{if } q \text{ then } u_1 \text{ else } u_2) &= \text{if } q \text{ then } \text{Targ}_{op}(u_1) \\
&\quad \text{else } \text{Targ}_{op}(u_2) \\
\text{Targ}_{op}(\text{for } \$x \in q \text{ return } u) &= \text{for } \$x \in q \text{ return } \text{Targ}_{op}(u) \\
\text{Targ}_{op}(\text{insert } q \text{ d } q') &= q' \quad (\text{insert}(d) : op) \\
\text{Targ}_{op}(\text{delete } q') &= q' \quad (\text{delete} : op) \\
\text{Targ}_{op}(\text{rename } q' \text{ as } a) &= q' \quad (\text{rename} : op) \\
\text{Targ}_{op}(\text{replace } q \text{ with } q') &= q' \quad (\text{replace} : op) \\
\text{Targ}_{op}(u) &= () \quad (\text{otherwise})
\end{aligned}$$

Figure 13: Target query of an update (with respect to update operator op)

$$\begin{aligned}
\text{FO}_y() &= \perp \\
\text{FO}_y(\$x) &= x = y \\
\text{FO}_y(\$x/\text{axis} :: *) &= \text{Axis}(x, y) \\
\text{FO}_y(\$x/\text{axis} :: \text{text}()) &= \text{Axis}(x, y) \wedge \text{text}(y) \\
\text{FO}_y(\$x/\text{axis} :: A) &= \text{Axis}(x, y) \wedge A(y) \\
\text{FO}_y((q_1, q_2)) &= \text{FO}_y(q_1) \vee \text{FO}_y(q_2) \\
\text{FO}_y(\text{if } q \text{ then } q_1 \text{ else } q_2) &= (\exists z. \text{FO}_z(q) \wedge \text{FO}_y(q_1)) \\
&\quad \vee (\neg(\exists z. \text{FO}_z(q)) \wedge \text{FO}_y(q_2)) \\
\text{FO}_y(\text{for } \$x \in q_1 \text{ return } q_2(x)) &= \exists x. \text{FO}_x(q_1) \wedge \text{FO}_y(q_2(x))
\end{aligned}$$

Figure 14: Translating selection queries to first-order formulas. In each case we assume y is not free in the query.

- The *path abstraction* translates a selection query to a set of paths (each of which can be viewed individually as an EFO formula)

These translations from queries to first-order formulas or paths are well-understood and straightforward (see e.g. [5]); we provide them in this appendix for completeness.

First-order translation. Figure 14 defines a function $\text{FO}_x(q)$ that maps a selection query $q(y_1, \dots, y_n)$ to an equivalent first-order formula $\phi(x, y_1, \dots, y_n)$, in the following sense:

PROPOSITION 1. *For any document selection query $q(\$doc)$, we have $\text{Sel}[\mathbb{Q}]q = \{1 \mid \sigma \models \phi(1, \text{root}(\sigma))\}$.*

Existential first-order abstraction. For some analysis techniques, it is necessary to over-approximate a selection query by a positive existential formula, that is, avoiding the use of negation or universal quantification. In Figure 14, we can see that only one case can involve negation, namely the case for the else-branch of a conditional. We can safely over-approximate a selection query by

$$\begin{aligned}
\text{Pt}(\Gamma, ()) &= \emptyset \\
\text{Pt}(\Gamma, \$x) &= \Gamma(x) \\
\text{Pt}(\Gamma, \$x/\text{step}) &= \{p/\text{step} \mid p \in \Gamma(x)\} \\
\text{Pt}(\Gamma, (q_1, q_2)) &= \text{Pt}(\Gamma, q_1) \cup \text{Pt}(\Gamma, q_2) \\
\text{Pt}(\Gamma, \text{if } q \text{ then } q_1 \text{ else } q_2) &= \text{Pt}(\Gamma, q_1) \cup \text{Pt}(\Gamma, q_2) \\
\text{Pt}(\Gamma, \text{for } \$x \in q_1 \text{ return } q_2) &= \text{Pt}(\Gamma[x := \text{Pt}(\Gamma, q_1)], q_2) \\
\text{Pt}(q) &= \text{Pt}([\$doc = /], q)
\end{aligned}$$

Figure 15: Approximating selection queries by paths. Here, Γ is an environment mapping each variable to a set of paths.

rewriting each conditional in it as follows:

$$\text{if } q \text{ then } q_1 \text{ else } q_2 \rightarrow (\text{if } q \text{ then } q_1), q_2$$

and then translating the query to a first-order formula as shown in Figure 14.

Path abstraction. Figure 15 defines a function $\text{Pt}(q)$ that constructs a set of XPath paths that covers all of the nodes that might be returned by q .

PROPOSITION 2. *For any selection query q and store σ , we have $\text{Sel}[\![q]\!] \sigma \subseteq \bigcup_{p \in \text{Pt}(q)} \text{Sel}[\![p]\!] \sigma$. Hence, if every pair of paths $p_1 \in \text{Pt}(q_1)$ and $p_2 \in \text{Pt}(q_2)$ are disjoint then so are q_1 and q_2 .*

C.4 Heuristics for path disjointness

The SMT-based approach can be used with path abstraction, since the paths translate into simple positive queries as described in Section 4. In addition we employ several heuristics for either simplifying or solving certain classes of path disjointness problems.

1. *Suffix incompatibility.* If two paths end with suffixes of child steps that are obviously incompatible (e.g. $/*/a/b$ and $/*/c/d/e$) then the paths cannot overlap.
2. *Displacement tests.* Define the *minimum displacement* of a path $\delta_{\min}(p)$ to be the smallest difference in height between the root of the tree and the node selected by the path, if this exists; otherwise $\delta_{\min}(p) = -\infty$. Similarly, define the *maximum displacement* $\delta_{\max}(p)$ to be the maximum difference in height between the root of the tree and the node selected by the path, if this exists; otherwise $\delta_{\max}(p) = +\infty$. For example, the minimum and maximum displacement of $/child :: */desc :: a$ are 2 and ∞ , respectively. Define the *displacement interval* $\delta(p)$ of p to be the set $\{n \in \mathbb{Z} \mid \delta_{\min}(p) \leq n \leq \delta_{\max}(p)\}$. Given paths p_1, p_2 , if $\delta_{\max}(p_i) < 0$ for either p_1 or p_2 then p_i is unsatisfiable (starting from the root of a tree) and so p_1 and p_2 are trivially disjoint. Otherwise, if $\delta(p_1) \cap \delta(p_2) = \emptyset$ then p_1 and p_2 are disjoint.
3. *Prefix-incompatibility for downward paths.* For two downward-only paths, if the prefixes of child steps are incompatible then the paths are disjoint. For example, $/a/b/c/*//d$ and $/a/b/e/*$ are prefix-incompatible.

Although this is a special case that is also solvable by Hammerschmidt et al.’s algorithm, this check is faster than their automaton construction and detects many common cases where the paths are disjoint.

If the paths are downward-only and their disjointness is not resolved by the heuristics, we employ the automaton construction in Hammerschmidt et al. [20]. Otherwise, we attempt the approaches via SMT-solving described in the body of the paper.

D. ADDITIONAL RELATED AND FUTURE WORK

Although view maintenance for XML queries has been explored in much prior work (see for example [7, 12]) the focus there has been the use of data structures that aid efficient recomputation when the view must be refreshed. Fast static techniques can effectively complement incremental techniques, especially in distributed settings, as argued here and in [27, 2].

Our query-rewriting algorithms for computing destabilizers are reminiscent of algebraic approaches to view-maintenance, e.g. those based on finite differencing [25, 19]. However, there is no formal connection between these problems.

Disjointness testing for XPath/XQuery is not well understood either in theory or practice. Hidders [21] established that XPath

satisfiability is in PTIME, and intersection testing is NP-hard, but the complexity of binary intersection of simple paths containing arbitrary axes remains open. Hammerschmidt et. al. [20] showed that it is in PTIME for downward-only paths, while Benedikt et al. [4] show that even satisfiability (and hence binary intersection) of downward paths is NP-hard in the presence of a schema. Genevès and Layaïda [16] present a translation from XPath containment problems to *MSO(Tree)*. Our translation to *mona* is similar. Subsequently, Genevès, Layaïda and Schmitt [17] have developed a system for analyzing XPath expressions with respect to schemas. Our approach to intersection testing using *yices* appears more efficient than their prototype but only solves a special case of this problem.

One of our main conceptual contributions is the idea of finitely representing collections of updates that can change a query result. This is motivated by prior work on XML projection [24, 6], which concentrates on sets of nodes as a representation system. We explore a more precise representation system here via the use of op-based destabilizers, but we intend to explore much richer systems in the future. For example, one can gain precision by utilizing descriptions that alternate conjunction and disjunction, stating a sequence to *simultaneously* overlap with several sets of atomic updates. In Example 2, a finer representation system would allow us to specify that for an update sequence to change Q_2 , it must either replace node 1, or it must both modify nodes 2 or 3 and also modify one of nodes 4 and 5. While we use these representation systems in conjunction with static analysis, they could also be used at runtime – either in the context of incremental view maintenance, or simply as a way of understanding how an anomalous query result could be changed.

E. BENCHMARK QUERIES AND UPDATES

The XMark queries are written in plain XQuery using complex XPath expressions and other high-level constructs, and as such need to be translated to the core language before static analysis. We translated them by hand following the rules in the XQuery standard. We also inlined let-binding expressions (this is not always safe in full XQuery but appears safe for the XMark queries). Finally, XMark queries also use several features that we do not directly handle, such as built-in functions (`count` or `equality`). We translated these queries to queries that do not use these features, but have similar dependencies. Thus, for a query such as Q10:

```
for $b in $auction//site/regions
return count($b//item)
```

we used the simplified query:

```
for $b in $auction//site/regions
return $b//item
```

The XPathMark queries were used as the basis for updates derived in a systematic way. For example, XPathMark query A_2 is:

```
$auction//closed_auction//keyword
```

and from this we generated an insert:

```
insert nodes <foo/>
into $auction//closed_auction//keyword
```