Lux: A lightweight, statically typed XML update language

James Cheney University of Edinburgh jcheney@inf.ed.ac.uk

ABSTRACT

Several proposals for updating XML have been introduced. Many of them have a rather complicated semantics due to the interaction of side-effects and updates, and some proposals also complicate the semantics of XQuery because arbitrary side-effecting update statements are allowed inside queries. Moreover, static typechecking has not been studied for any proposed XML update language.

In this paper, we survey prior work on XML update languages and motivate an alternative approach to updating XML. We introduce an update language called Lux, which stands for Lightweight Updates for XML. Lux can perform relational database-style updates, has a simple, deterministic operational semantics, and has a sound static type system based on regular expression types and structural subtyping.

1. INTRODUCTION

Query and transformation languages for tree-structured data, such as XML, have been extensively investigated. XML, in particular, is now widely used for streaming, exchanging and storing data, and standards for validating (DTDs, XML Schema [15]), transforming (XSLT [12]), and querying (XQuery [4]) XML data have reached an advanced stage; formal semantics and type systems for such languages are well-understood, and mature, efficient implementations exist or are in development. However, dedicated update languages for XML have received much less attention.

Several languages for updating XML have been proposed over the last few years [22, 30, 5, 31, 29, 7, 19, 9, 20, 8], with a trend towards increasing expressiveness and semantic complexity. Some design choices and other features¹ of the various proposals can be characterized by the answers to a number of questions shown in Figure 1, many of which correspond directly to **MUST**- or **SHOULD**-requirements in the current draft of the W3C's XQuery Update Facility

Proceedings of the 5th ACM SIGPLAN Workshop on Programming Language Technologies for XML (PLAN-X 2007). January 20, 2007, Nice, France. Requirements [10]. Early XML update language proposals considered only sequences of atomic operations without iteration, but more recent proposals include sophisticated forms of iteration and control of side-effects so are much more expressive. Expressiveness is certainly a reasonable design goal, but it is important to weigh the benefits of added expressiveness against their costs in increased semantic and implementation complexity: sophisticated update operations are more difficult to optimize, typecheck, and analyze.

The situation in the relational/SQL world is very different. In SQL, side-effecting update statements cannot be embedded within arbitrary queries; moreover, the update language is considerably weaker than the query language. For example, SQL's INSERT-DELETE-UPDATE statements cannot simulate joins. This simplicity has a number of advantages. First, update optimization appears much easier than query optimization. Second, since updates cannot occur inside queries, query optimization is unburdened by the obligation to prove the absence (or irrelevance) of side-effects. A third benefit of SQL's simplistic approach to updates is that it is possible to statically typecheck them before execution. None of this is to say that SQL's update facility is perfect. But understanding what made it successful despite its limitations may be helpful when designing an update language for XML, especially given that updating relational data stored as XML is an important use case for XML update languages [10, 24].

In recent XML update proposals, however, updates introduce considerable complexity over XQuery. Some proposals permit updates as ordinary XQuery expressions, which means that many of XQuery's advantages as a purely functional language disappear due to the presence of side-effects. Other proposals separate queries and updates, but permit updating sets of nodes selected using arbitrary XPath or XQuery expressions. In either case, the results of iterative update operations may be nondeterministic. To avoid nondeterminism arising from updates and iteration, many proposals use a two-phase "snapshot" semantics. Even so, predicting the behavior of an update based on this semantics seems difficult. As a result, nontrivial static analyses appear necessary to determine whether updates can be reordered or applied eagerly [1, 2, 18].

Static typechecking also appears likely to pose significant challenges for most of the extant update languages. To the best of our knowledge, static typing has not been investigated for *any* such proposal. This may seem surprising because type systems for XML transformation and query

¹Note that many of these features are subject to change as the result of further study of the languages.

Atomic: What atomic updates are allowed?

I: Inserting data before or after a node

A: Appending or prepending data to a node's child list

D: Deleting a node (and its subtree)

 ${\bf R}:$ Replacing a node's value

 $\mathbf{N:}$ Renaming a node's label

Selection: How can parts of a document be selected for updating?

XPath: Using XPath expressions

XQuery: Using XQuery expressions

Ad hoc: Using ad hoc formalism

- **Style:** What style of programming is used for update statements, and how do such statements interact with XQuery expressions?
 - Ad hoc: Update statements use a separate, ad hoc syntax unrelated to XQuery (e.g. sequences of atomic updates)
 - **Separate:** Update statements use an XQuery-like syntax and contain XQuery expressions but are not allowed as XQuery expressions themselves.
 - **Combined:** Update statements are allowed as XQuery expressions of type ()
- **Nondeterminism:** How is nondeterminism arising from the interaction of side-effects and evaluation order handled?

Unspecified: issue not discussed; no formal semantics **Checked:** pending update lists with determinism en-

- forced by dynamic checks
- **Fixed order:** pending update lists with determinism enforced by fixing an evaluation order
- **User control:** nondeterminism controlled by user annotations

Mutability: Can an update change the value of a variable?

Transforms: Can values be constructed by "transformation" (applying an update to a copy of some data)?

Formal: Has a formal semantics been developed?

Types: Has a static type system been developed?

Figure 1: Some design choices and other characteristics of XML update languages.

languages have received considerable attention [14, 26, 17, 21]. However, XML values are immutable in all of these languages. In most XML update proposals, mutability is incorporated using standard imperative programming features; in particular, current proposals allow for updating XML document nodes using their identities as references.

While imperative programming techniques are well-understood and familiar, they also have many drawbacks from the point of view of typechecking. Updates and variables can "alias"; that is, an update may change the value of a variable in ways that may be hard to predict statically. Thus, when an update takes place, a type system would have to determine how to modify the types of variables to reflect the changes. This means that variables are essentially mutable, even though many proposals do not provide for direct assignment to variables. For these reasons, we will adopt an approach in which variables are always immutable; mutability is carefully controlled by the type system. A third difficulty encountered in many proposals is the issue of node identity. The XQuery and XPath data models expose node identifiers for all nodes in the document tree; as noted above, many update languages uses these identifiers as update targets. Many update languages therefore view the XML data as a "heap" mapping node identifiers to records containing other identifiers or base values (such as integers or strings). Their semantics accounts for "in-place" updates; thus, an update might change a child of a node while leaving the node intact. Also, two updates which differ with respect to the node identity semantics may be equivalent if node identities are ignored.

However, nothing in XQuery or XPath's existing semantics specifies how node identities should be transformed when a document is updated. Instead, these languages only specify their behavior within queries. Consequently, it seems reasonable to consider other policies for dealing with node identity during updates. One attractive policy is to ignore node identity altogether for the purpose of updating: that is, to define a value-based semantics for updates, and disallow the use of node identifiers when defining updates. This approach leads to a much simpler semantics which satisfies more equational laws that can be used to optimize updates.

The purpose of this paper is to investigate an XML update language which is less expressive than most proposals but has a simple, deterministic semantics and can be statically typechecked. Our language, LUX (Lightweight Updates for XML), has the following properties:

- All standard update operations (insert, append, delete, rename, replace) are supported.
- Nodes can be selected using only restricted XPath expressions; only horizontal iteration is allowed.
- Side-effecting update statements are not allowed as queries.
- Variables are immutable.
- The semantics is value-based, one-pass, and deterministic.
- Updates can be statically typechecked using regular expression types.

Our approach owes much to an update language called CPL+ for a generalized form of the relational data model, developed by Liefke and Davidson [23]. Lux can be viewed as an adaptation of their approach to an XML setting.

Lux is work in progress. Currently, we have developed a core language that consists of a small number of simple, orthogonal features, with a deterministic operational semantics and sound type system. We have also implemented a proof-of-concept interpreter and typechecker for the core language. In this paper, we will present a motivating high-level language; however, this high-level language is only meant to be illustrative and is subject to change.

The structure of this paper is as follows. Section 2 introduces the high-level version of Lux via a sequence of examples. Section 3 defines the operational semantics of the core Lux language, and gives translations of the high-level queries. Section 4 defines a type system for Lux updates, proves its soundness with respect to the operational semantics, and discusses the typechecking decision problem. Section 5 relates our approach to previous database and XML

	Atomic	Selection	Style	Nondeterminism	Transforms	Mutability	Formal	Types
XML:DB XUpdate [22]	IADRN	XPath	Ad hoc	Unspecified	No	No	No	No
Updating XML [30]	IDRN	XPath	Separate	Unspecified	No	Yes	No	No
Bruno et al. [5]	IADRN	XPath	Ad hoc	Unspecified	Yes	No	No	No
XML-RL Update [31]	IADR	Ad hoc	Ad hoc	Unspecified	No	Yes	Yes	No
UpdateX [29]	IADR	XQuery	Separate	Fixed order	No	Yes	Yes	No
Context logic [7]	IADR	XPath(child)	Ad hoc	Fixed order	No	Yes	Yes	No
XQuery [!] [19]	IADRN	XQuery	Combined	User control	No	Yes	Yes	No
XQuery Update [9]	IADRN	XQuery	Combined	Checked	Yes	Yes	No	No
LiXQuery ⁺ [20]	IADRN	XQuery	Combined	User control	Yes	Yes	Yes	No
XQueryP [8]	IADRN	XQuery	Combined	User control	Yes	Yes	No	No
Lux	IADRN	XPath(child)	Separate	N/A	No	No	Yes	Yes

Table 1: Summary of XML update language characteristics. See Figure 1 for explanations of the column entries and Section 5 for further discussion.

update languages and work on static type systems for XML; Section 6 discusses extensions and future directions and Section 7 concludes.

2. OVERVIEW AND EXAMPLES

2.1 Full language syntax

As with CPL+, XQuery and many other languages, we introduce a high-level, readable syntax with complex, overlapping operations which can be translated to a much simpler core language. Update statements allow us to create and drop (delete) top-level documents, insert XQuery values into trees based on paths (at the beginning or end of a child list or before or after a node), update subtrees at a given path, delete from a path all subtrees satisfying a condition, rename subtrees at a given path, or replace subtrees. Statements can also be sequenced and restricted using conditional expressions, and let-expressions can be used to bind variables to the values of expressions.

Path expressions in updates are used in two distinct senses, "tree-oriented" and "sequence-oriented". In tree-oriented updates, the update expects a singleton tree and is executed once for each subtree selected by the path expression; in sequence-oriented updates, the update operates on an arbitrary sequence and is executed on the child-list of each node selected by the path expression.

Tree-oriented insertions (INSERT BEFORE/AFTER) insert a value before or after each node selected by the path expression. Similarly, tree-oriented deletes (DELETE) assume that we have selected a single node and want to delete it. Plain replacement REPLACE Path WITH is tree-oriented by default.

Sequence-oriented insertions (INSERT INTO) insert a value into the child-list of each selected node, at the beginning (AS FIRST) or end (AS LAST). Sequence-oriented deletes (DELETE FROM) delete nodes from the child-list of each selected node. The REPLACE Path CONTENT WITH operator is sequence-oriented; it replaces the content of a path with new content.

The "dot" path expression . refers to the currently selected part of the tree; its value depends on context. Path expressions always begin with ./, but in the examples we usually omit this prefix. In what follows, we assume familiarity with XQuery and XPath syntax, and with XDuce-style regular expression types for XML data.

2.2 Creating a database and loading data

Suppose we start an XML database with no pre-loaded data; its type is the empty sequence (). We want to create

UpdStmt	::=	CREATE DOCUMENT Name
		DROP DOCUMENT Name
		INSERT (AS (FIRST LAST))? INTO Path
		VALUE Expr
		INSERT (BEFORE AFTER) Path VALUE Expr
		UPDATE (IN?) Path BY UpdStmt
		LET \$var := Expr IN UpdStmt
		DELETE (FROM)? Path
		RENAME Path TO Name
		REPLACE Path (CONTENT?) WITH Expr
		UpdStmt WHERE Cond
		UpdStmt ; UpdStmt
		{ UpdStmt }
Path	::=	. Path/Name Path/*
Expr	::=	XQuery expressions
Cond	::=	XQuery/XPath conditional expressions

Figure 2: Concrete syntax of Lux updates.

a small database listing books and authors. The following LUX updates accomplish this:

U1 : CREATE DOCUMENT books[]; CREATE DOCUMENT authors[]

After this update, the database has type

```
books[],authors[]
```

Now we want to load some XML data into the database. Since XML text is included in XQuery's expression language, we can just do the following:

This results in a database with type

2.3 Updating data

The data we initially inserted was incomplete and incorrect. Now we want to fill in some of the missing data and correct an error.

This update leaves the structure of the database unchanged.

2.4 Restructuring the database

SQL updates include ALTER TABLE statements that make it possible to add or remove columns. Similarly, we can add an element to each book in books as follows:

```
U5 : INSERT AS LAST INTO books/book
VALUE publisher["Grinch"]
```

After U5, the books database has type

Now perhaps we want to add a co-author; for example, perhaps Lewis Carroll collaborated on "Through the Looking-Glass" with Charles Dickens. This is not as easy as adding the publisher field to the end because we need to select a particular node to insert before or after. In this case we happen to know that there is only one author, so we can insert after that; however, this would be incorrect if there were multiple authors, and we would have to do something else (such as inserting before the title).

```
U6 : UPDATE books/book BY
INSERT AFTER author
VALUE <author>Charles Dickens</author>
WHERE name = "Through the Looking-Glass"
```

Now the **books** part of the database has the type:

Now that some books have multiple authors, we might want to change the flat author lists to nested lists:

```
U7 : REPLACE books/book WITH
        <book><authors>{author}</authors>
        {title}{year}{publisher}</book>
```

This visits each book and changes its structure so that the authors are grouped into an authors element. The resulting database has type:

2.5 Deleting data

Suppose we later decide that the publisher field is unnecessary after all. We can get rid of it using the following update:

U8 : DELETE books/book/publisher

This results in schema

Now suppose Lewis Carroll no longer interests us and we wish to remove all of his books from the database.

```
U9 : DELETE FROM books
    WHERE book/authors/author = "Lewis Carroll"
```

This update does not modify the type of the database. Finally, we can delete a top-level document as follows:

U10 : DROP DOCUMENT authors

2.6 Execution model

In general, an update operation with a path expression in it is evaluated as follows: The path expression is evaluated, yielding a sequence of *context nodes*. Then the corresponding basic update operation (insert, delete, etc.) is applied to each node in this set in turn, using it as the context node.

If arbitrary XPath or XQuery expressions could be used to select context nodes, it would be easy to construct examples for which the result of an update depends on traversal order. For example, suppose the document is of the form <a>b/>. If the following update were allowed:

```
UPDATE // BY { DELETE a/b; RENAME * TO c }
```

then the result depends on the order in which the updates are applied. The two possible results are $\langle a \rangle$ and $\langle a \rangle \langle c \rangle \langle a \rangle$. Not only is this behavior nondeterministic, it is difficult to typecheck (and likely also difficult to parallelize). For this reason, we place severe restrictions on the path expressions which may be used define a context node selection.

We identify two key properties which help to ensure that updates are deterministic and can be typechecked. First, an update can only modify data at or beneath its current context node. We call this the side-effect isolation property. For example, navigating to the context node's parent and then modifying another sibling is not allowed. In addition, whenever an iterative update occurs (that is, one involving a loop traversing a number of nodes), we require that the iteration is deterministic; that is, the result of an iterative update is independent of the order in which the nodes are updated. We call this the traversal-order independence property.

To ensure isolation of side effects and traversal-order independence, it is sufficient (but maybe not necessary) to restrict the XPath expressions that can be used to select a context node. Specifically, only the child axis² and path composition are allowed. This ensures that only descendants of a given context node can be selected as the new context node and that a selection of new context nodes contains no pairs of nodes in an ancestor-descendant relationship. Consequently, the side effects of an update are confined to the subtree of its context node, and the result of an iteration is

 $^{^2{\}rm The}$ attribute axis can also be handled, but the descendant, parent, and sibling axes seem nontrivial to handle.

independent of the traversal order. This keeps the semantics deterministic and helps make typechecking feasible.

In Section 3, we will define a core language that satisfies the side-effect isolation and traversal-order independence properties. Our implementation currently compiles a high-level language similar to the one outlined in this section to the core language. Some of the compilation rules used by the implementation are shown in Section 6.5; however, treatment of the full language and compilation and optimization issues is left for future work.

2.7 Non-design goals

Although we believe LUX works well for database-style updates (especially for SQL-like updates over XML data), there are several things that other proposals do that we make no attempt to do. Whether this is an advantage or disadvantage depends on the application.

Node identity: The XQuery data model provides identifiers for all nodes. Many XML update proposals take node identities into account and can use them as to update parts of the tree "by reference". In contrast, LUX's semantics is purely value-based: node identities are inaccessible to updates. It does not appear difficult to add support for node identity within XQuery expressions in LUX; however, this seems to make static typechecking extremely difficult, because of the possibility of aliasing and nondeterminism.

Pattern matching: Many transformation/query languages (e.g. [21, 12]) and some update languages (e.g. [23, 31]) allow defining transformations by *pattern matching*, that is, matching tree patterns against the data. Pattern matching is crucial for transforming XML, but we believe it is not as important for updates. We have not considered general pattern matching in Lux, in order to keep the type system and operational semantics as simple as possible.

Side-effects in queries: Several motivating examples for XQuery! [19] depend on the ability to perform sideeffects within queries. Examples include logging accesses to particular data or profiling or debugging an XQuery program. Lux cannot be used for these applications because side-effects are not allowed within queries.

Deep updates/recursion: LUX is only capable of horizontal iteration: iteration over all of the children of a node. LUX does not include recursively defined updates or recursive types, and cannot meaningfully update data arbitrarily deeply in the tree. Thus, a single LUX update cannot, for example, update all *person* elements everywhere in the tree in-place—even relative to a non-recursive DTD. LUX is not suitable for performing XML transformations such as rendering XML data to XHTML, or updating genuinely recursive data. While it does not seem problematic to support the descendant axis in query expressions within LUX (using the same typing rules as Collazzo et al. [13]), it appears difficult to extend our approach to typecheck updates that make use of recursion or the descendant, sibling or parent axes.

Joins: Lux updates cannot, by themselves, iterate superlinearly over the data to perform a join. However, updates can use embedded XQuery expressions to restructure the database, for example:

INSERT INTO tmp VALUE
FOR \$x in books,\$y in authors
WHERE \$x/author = \$y/name
RETURN <pair>{\$x}{\$y}</pair>

3. FORMALIZATION

3.1 Core query language

Because Lux uses queries to construct values, we need to introduce a query language and define its semantics before doing the same for Lux. We will use a small fragment of XQuery called μ XQ, introduced by Colazzo, Ghelli, Manghi and Sartiani [13].

Following [13], we distinguish between tree values $t \in Tree$, which include strings $w \in \Sigma^*$ (for some alphabet Σ), boolean values, and singleton trees n[v]; and (forest) values $v \in Val = Tree^*$, which are sequences of tree values:

```
(Forest) values v ::= () | t, v
Tree values t ::= n[v] | w | true | false
```

We overload the set membership symbol \in for trees and forests: that is, $t \in v$ means that t is a member of v considered as a list. Two forest values can be concatenated by concatenating them as lists; abusing notation, we identify trees t with singleton forests t, () and write v, v' for forest concatenation. We define a comprehension operation on forest values as follows:

$$[f(x) \mid x \in ()] = ()$$

$$[f(x) \mid x \in t, v] = f(t), [f(x) \mid x \in v]$$

This operation takes a forest (t_1, \ldots, t_n) and a function f(x) from trees to forests and applies f to each tree t_i , concatenating the resulting forests in order. Comprehensions satisfy basic Vonda laws as well as some additional equations (see [16]). We write \approx for equality of tree or forest values; this is just structural equality.

Again following [13], we distinguish between tree variables $\bar{x} \in TVar$, introduced by for, and forest variables, $x \in Var$, introduced by let. We write $X \in Var \cup TVar$ for an arbitrary tree or forest variable. The other syntactic classes of our variant of μXQ include labels $l, m, n \in Lab$ and expressions $e \in Exp$; the abstract syntax of expressions is defined by the following BNF grammar:

The distinguished variables \bar{x} in for $\bar{x} \in e$ return e'(x) and x in let x = e in e'(x) are bound in e'(x). Here and elsewhere, we employ common conventions such as considering expressions containing bound variables equivalent up to α -renaming and employing a richer concrete syntax including, for example, parentheses.

To simplify the presentation, we split μXQ 's projection operation \bar{x} child :: l into two expressions: child projection (\bar{x}/child) which returns the children of \bar{x} , and node name filtering (e :: n) which evaluates e to an arbitrary sequence and selects the nodes labeled n. We also define the children of a value other than n[v] to be the empty sequence rather than leaving it undefined. Thus, the ordinary child axis expression \bar{x} child :: n is syntactic sugar for (\bar{x}/child) :: n and the "wildcard" child axis is definable as \bar{x} child :: $* = \bar{x}/\text{child}$. We also consider only one built-in operation, equality.

An environment is a pair of functions $\sigma : (Var \to Val) \times (TVar \to Tree)$. Abusing notation, we write $\sigma(x)$ for $\pi_1(\sigma)(x)$ and $\sigma(\bar{x})$ for $\pi_2(\sigma)(\bar{x})$; similarly, $\sigma[x := v]$ and $\sigma[\bar{x} := t]$

$$\begin{array}{rcl} children(n[f]) &=& f\\ children(v) &=& () & (v \not\approx n[v']) \end{array}$$

$$\begin{bmatrix} \texttt{true} \end{bmatrix} \sigma &=& \texttt{true} \\ \llbracket() \rrbracket \sigma &=& () \\ \llbracketn[e] \rrbracket \sigma &=& n[\llbracket e \rrbracket \sigma] \\ \llbracketx \rrbracket \sigma &=& \sigma(x) \end{array} \qquad \begin{bmatrix} \texttt{false} \rrbracket \sigma &=& \texttt{false} \\ \llbrackete, e' \rrbracket \sigma &=& \llbrackete \rrbracket \sigma, \llbrackete' \rrbracket \sigma \\ \llbracketw \rrbracket \sigma &=& w \\ \llbracketx \rrbracket \sigma &=& \sigma(x) \end{array} \qquad \begin{bmatrix} \texttt{false} \rrbracket \sigma &=& \texttt{false} \\ \llbrackete, e' \rrbracket \sigma &=& \texttt{false} \\ \llbracketw \rrbracket \sigma &=& w \\ \llbracketv \rrbracket \sigma &=& w \\ \llbracketv \rrbracket \sigma &=& mtion \\ \llbrackete = e^{t} \rrbracket \sigma &=& \begin{cases} \llbrackete_{1} \rrbracket \sigma & \llbracketc \rrbracket \sigma \approx \texttt{true} \\ \llbrackete_{2} \rrbracket \sigma & \llbracketc \rrbracket \sigma \approx \texttt{false} \\ \llbrackete = e^{t} \rrbracket \sigma &=& \begin{cases} \verb"true} & \llbrackete \rrbracket \sigma \approx \texttt{false} \\ \llbrackete = e^{t} \rrbracket \sigma &=& \begin{cases} \verb"true} & \llbrackete \rrbracket \sigma \approx \texttt{false} \\ \verb"true} & \llbrackete \rrbracket \sigma \approx \texttt{false} \\ \llbrackete \rrbracket \sigma \not\approx \texttt{false} \end{bmatrix} \sigma \\ \llbrackete :: n \rrbracket \sigma &=& [n[v] \mid n[v] \in \llbrackete \rrbracket \sigma] \\ \llbracket\bar{x}/\texttt{child} \rrbracket \sigma &=& children(\sigma(\bar{x})) \\ \llbracketfor \ \bar{x} \in e_1 \ \texttt{return} e_2 \rrbracket \sigma &=& \llbrackete_{2} \rrbracket \sigma \ \bar{x} :=t \rrbracket \mid t \in \llbrackete_{1} \rrbracket \sigma \end{bmatrix}$$

Figure 3: Semantics of query expressions.

denote the corresponding environment updating operations. The semantics of expressions is defined as shown in Figure 3. The function *children* : $Val \rightarrow Val$ maps a singleton tree n[f] to its child list and maps other values to the empty sequence. We write $\sigma \vdash e \Rightarrow v$ when $v = [\![e]\!]\sigma$.

3.2 Core update language

We now introduce the core Lux update language, which includes statements $s \in Stmt$, tests $\phi \in Test$, and directions $d \in Dir$:

- d ::= left | right | children | iter

Updates include standard constructs such as the no-op skip, sequential composition, conditionals, and let-binding. The basic update operations include insertion insert e, which inserts a value provided the context is empty; deletion delete, which deletes the selected value; and the "snapshot" operation snapshot x in s, which binds x to the current value of the context node and then applies an update s, which may refer to x. Also, snapshot is not equivalent to XQuery!'s snap operator; snapshot binds x to an immutable value which can be used in s, whereas snap delimits a "snapshot semantics" execution of a side-effecting query. There is no way to refer to the context node of an update within a μ XQ query without using snapshot.

Updates also include *tests* ϕ ?s which allow us to examine the local structure of a tree value and perform an update if the structure matches. The node label test n?s checks whether the tree is of the form n[v], and if so executes s, otherwise is a no-op; the wildcard test only checks that the value is a singleton tree. Similarly, **bool**?s and **string**?s test whether a tree value is a boolean or string value.

Finally, updates include *navigation* operators that change the selected part of the tree, and perform an update on the sub-selection. The **left** and **right** operators perform an update on the empty sequence located to the left or right

$\sigma; v \vdash s \Rightarrow^{\tt U} v'$	
---	--

$$\begin{array}{c} \overline{\sigma; v \vdash \operatorname{skip} \Rightarrow^{\operatorname{U}} v} & \overline{\sigma; v \vdash s \Rightarrow^{\operatorname{U}} v_1 \quad \sigma; v_1 \vdash s' \Rightarrow^{\operatorname{U}} v_2} \\ \overline{\sigma; v \vdash \operatorname{skip} \Rightarrow^{\operatorname{U}} v} & \overline{\sigma; v \vdash s; s' \Rightarrow^{\operatorname{U}} v_2} \\ \hline \overline{\sigma; v \vdash \operatorname{let} x = e \text{ in } s \Rightarrow^{\operatorname{U}} v_2} \\ \hline \overline{\sigma; v \vdash \operatorname{let} x = e \text{ in } s \Rightarrow^{\operatorname{U}} v_2} \\ \hline \overline{\sigma; v \vdash \operatorname{if} e \text{ then } s_1 \text{ else } s_2 \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{if} e \text{ then } s_1 \text{ else } s_2 \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{if} e \text{ then } s_1 \text{ else } s_2 \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{if} e \text{ then } s_1 \text{ else } s_2 \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{if} e \text{ then } s_1 \text{ else } s_2 \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{if} e \text{ then } s_1 \text{ else } s_2 \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{insert} e \Rightarrow^{\operatorname{U}} v} \quad \overline{\sigma; v \vdash \operatorname{delete} \Rightarrow^{\operatorname{U}} ()} \\ \hline \overline{\sigma; v \vdash \operatorname{snapshot} x \text{ in } s \Rightarrow^{\operatorname{U}} v'} \quad \overline{\sigma; v \vdash \operatorname{delete} \Rightarrow^{\operatorname{U}} v} \\ \hline \overline{\sigma; v \vdash \operatorname{snapshot} x \text{ in } s \Rightarrow^{\operatorname{U}} v'} \quad \overline{\sigma; v \vdash s \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{shapshot} x \text{ in } s \Rightarrow^{\operatorname{U}} v'} \quad \overline{\sigma; v \vdash s \Rightarrow^{\operatorname{U}} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{left}[s] \Rightarrow^{\operatorname{U}} v', v} \quad \overline{\sigma; v \vdash \operatorname{shap} v'} \\ \hline \overline{\sigma; v \vdash \operatorname{left}[s] \Rightarrow^{\operatorname{U}} v', v} \quad \overline{\sigma; v \vdash \operatorname{right}[s] \Rightarrow^{\operatorname{U}} v, v'} \\ \hline \overline{\sigma; v \vdash \operatorname{left}[s] \Rightarrow^{\operatorname{U}} v', v'} \quad \overline{\sigma; v \vdash \operatorname{right}[s] \Rightarrow^{\operatorname{U}} v, v'} \\ \hline \overline{\sigma; v \vdash \operatorname{left}[s] \Rightarrow^{\operatorname{U}} v'_1, v'_2} \quad \overline{\sigma; () \vdash \operatorname{iter}[s] \Rightarrow^{\operatorname{U}} ()} \\ \hline \end{array}$$

Figure 4: Operational semantics of updates.

of a value. The **children** operator applies an update to the child list of a tree value. Conversely, the **iter** operator applies an update to each tree value in a forest.

We distinguish between singular (unary) updates which apply only when the context is a tree value and plural (multiary) updates which apply to a sequence. Tests ϕ ?s are always singular. The **children** operator applies a plural update to all of the children of a single node; the **iter** operator applies a singular update to all of the elements of a sequence. Other updates can be either singular or plural in different situations. Our type system (presented in Section 4) tracks arity as well as input and output types in order to ensure that updates are well-behaved.

Figure 4 shows the operational semantics of Core Lux. We write $\sigma; v \vdash s \Rightarrow^{U} v'$ to indicate that given environment σ and context value v, statement s evaluates to value v'. The rules for tests are defined in terms of the following semantic interpretation of tests:

$$\begin{split} \llbracket \texttt{bool} \rrbracket &= \ \{\texttt{true},\texttt{false} \} \\ \llbracket \texttt{string} \rrbracket &= \ \Sigma^* \\ \llbracket n \rrbracket &= \ \{n[v] \mid v \in Val \} \\ \llbracket * \rrbracket &= \ \{n[v] \mid n \in Lab, v \in Val \} \end{split}$$

Note that while we described the evaluation of LUX updates informally in terms of nodes, we can define the semantics entirely in terms of forest and tree values, without needing to define an explicit store. This would not be the case if we considered full XQuery, which includes node identity comparison operations.

It is straightforward to show that an update run against a given environment and input value has most one output value:

THEOREM 1. Let σ, v, s, v_1, v_2 be given such that $\sigma; v \vdash$

$\texttt{repl} \ e$	=	$\texttt{delete}; \texttt{insert} \ e$
${\tt delete} \ n$	=	iter[n?delete]
$\texttt{if} \ c \ \texttt{then} \ s$	=	if c then s else skip
update n by $\left[s ight]$	=	$\mathtt{iter}[n?s]$
update in n by $\left[s ight]$	=	update n by $[\texttt{children}[s]]$
update x in n by $[s]$	=	update n by [snapshot x in children $[s]$]
$\begin{array}{l} \texttt{delete} \ x \ \texttt{from} \ n \\ \texttt{where} \ c \end{array}$	=	update x in n by [if c then delete]
$\texttt{insert}^{\leftarrow} e$	=	left[insert e]
$\texttt{insert}^{\rightarrow} \ e$	=	$\texttt{right}[\texttt{insert} \ e]$

Figure 5: Abbreviations for derived update expressions.

 $s \Rightarrow^{\mathsf{U}} v_1 \text{ and } \sigma; v \vdash s \Rightarrow^{\mathsf{U}} v_2.$ Then $v_1 = v_2.$

PROOF. Straightforward by induction on the structures of the two derivations. The interesting cases are those for conditionals and iteration, since they are the only statements that have more than one applicable rule. However, in each case, only matching pairs of rules are applicable. \Box

Figure 6 shows the example updates from Section 2 translated to core Lux. To simplify the presentation, we use full XQuery syntax for subqueries, rather than more verbose core μ XQ queries, and use several abbreviations introduced in Figure 5.

4. TYPE SYSTEM

As noted earlier, some updates expect that the input value is a singleton (for example, **children**, n?s, etc.) while others work for an arbitrary sequence of trees. Singular (treeoriented) updates may fail if applied to a sequence. This makes programming and debugging updates tricky, so our type system should prevent such run-time errors. Moreover, as with all XML transformation languages, we often would like to ensure that when given an input tree of some type τ , an update is guaranteed to produce an output tree of some other type τ' . For example, updates made by non-privileged users are usually required to preserve the database schema.

4.1 Types, subtyping and typing rules

We consider a regular expression type system with structural subtyping, similar to those considered in several transformation and query languages for XML [21, 13, 16].

$$\begin{aligned} \tau & ::= & \alpha \mid () \mid \tau \mid \tau' \mid \tau, \tau' \mid \tau^* \\ \alpha & ::= & \texttt{bool} \mid \texttt{string} \mid n[\tau] \end{aligned}$$

We call types of the form α singular types (or tree types), and types of all other forms plural (or forest types). It should be obvious that a value of singular type must always be a sequence of length one (that is, a tree); plural types may have values of any length. There exist plural types with only values of length one, but which are not syntactically singular (for example int|bool). As usual, the + and ? quantifiers can be defined as follows: $\tau^+ = \tau, \tau^*$ and $\tau^? = \tau |()$. Type variables and recursive type definitions are not included.

$$\begin{array}{rcl} U_1 &:& \mathrm{insert}\;(books[], authors[])\\ U_2 &:& \left\{\begin{array}{l} \mathrm{update}\;\mathrm{in}\;books\;\mathrm{by\;insert}^{\rightarrow}\;(\ldots);\\ \mathrm{update}\;\mathrm{in}\;authors\;\mathrm{by\;insert}^{\rightarrow}\;(\ldots);\\ \mathrm{update}\;\mathrm{in}\;authors\;\mathrm{by\;insert}^{\rightarrow}\;(\ldots);\\ \end{array}\right.\\ U_3 &:& \left\{\begin{array}{l} \mathrm{update}\;\mathrm{in}\;books\;\mathrm{by\;[}\\ \mathrm{update}\;\mathrm{in}\;book\;\mathrm{by\;[}\\ \mathrm{update}\;\mathrm{in}\;author\;\mathrm{by\;[repl~"CD"]};\\ \mathrm{update\;in\;}update\;\mathrm{in\;}update\;\mathrm{books\;by\;[}\\ \mathrm{update\;in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}author\;\mathrm{by\;[}\\ \mathrm{update\;}author\;\mathrm{by\;[}\\ \mathrm{update\;}author\;\mathrm{by\;[}\\ \mathrm{update\;}author\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x\;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x;\mathrm{in\;}book\;\mathrm{by\;[}\\ \mathrm{update\;}x;\mathrm{in\;}book\;\mathrm{by\;}x;\mathrm{unthor\;}x;\mathrm{unthor\;}x;\mathrm{untho$$

 U_{10} : delete *authors*

Figure 6: Translations of example queries into core Lux.

As in XDuce, a type denotes a set of matching values:

$$\begin{split} \llbracket \text{string} &= \Sigma^* \qquad \llbracket \text{bool} \rrbracket = \{ \text{true, false} \} \\ \llbracket () \rrbracket &= \{ () \} \qquad \llbracket n[\tau] \rrbracket = \{ n[v] \mid v \in \llbracket \tau \rrbracket \} \\ \llbracket \tau, \tau' \rrbracket &= \{ v, v' \mid v \in \llbracket \tau \rrbracket, v' \in \llbracket \tau' \rrbracket \} \\ \llbracket \tau \mid \tau' \rrbracket &= \llbracket \tau \rrbracket \cup \llbracket \tau' \rrbracket \\ \llbracket \tau^* \rrbracket &= \{ () \} \cup \{ v_1, \dots, v_n \mid v_1 \in \llbracket \tau \rrbracket, \dots, v_n \in \llbracket \tau \rrbracket \} \end{split}$$

In addition, we define a binary subtyping relation on types. A type τ_1 is a subtype of τ_2 ($\tau_1 <: \tau_2$), by definition, if $[\![\tau_1]\!] \subseteq [\![\tau_2]\!]$. Since our system simply re-uses the subtyping judgment from XDuce's type system, we know that subtyping is decidable; although XDuce subtyping is EXPTIME-complete in general, the algorithm of Hosoya, Vouillon and Pierce is well-behaved in practice [21]. Therefore, we shall not give explicit inference rules for checking or deciding subtyping, but treat it as a "black box". We also consider subtyping relations between tree types and tests: we say that $\alpha <: \phi$ if $[\![\alpha]\!] \subseteq [\![\phi]\!]$. This is decidable using the following rules:

 $\overline{\texttt{bool} <:\texttt{bool}} \quad \overline{\texttt{string} <:\texttt{string}} \quad \overline{n[\tau] <:n} \quad \overline{n[\tau] <:*}$

We consider contexts Γ of the form

$$\Gamma ::= \cdot \mid x : \tau \mid \bar{x} : \alpha$$

 $\Gamma \vdash e : \tau$

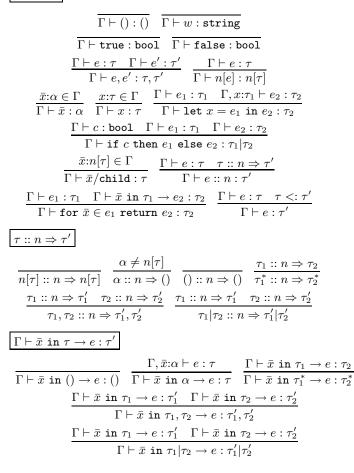


Figure 7: Query well-formedness.

that is, tree variables may only be bound to tree types. As usual, we assume that variables in contexts are distinct; this convention implicitly constrains the inference rules. We write $\llbracket \Gamma \rrbracket$ for the set of all environments σ such that $\sigma(X) \in \Gamma(X)$ for all $X \in dom(\Gamma)$.

The typing judgment for queries is $\Gamma \vdash e : \tau$ (that is, in context Γ , expression e has type τ); following [13], there are two auxiliary judgments, $\Gamma \vdash \bar{x}$ in $\tau \to s : \tau'$, used for typechecking for-expressions, and $\tau :: n \Rightarrow \tau'$, used for typechecking label matching expressions e :: n. The rules for these judgment are shown in Figure 7. There are two typing judgments for updates: singular well-formedness $\Gamma \vdash^{-1} \{\alpha\} \ s \ \{\tau'\}$ (that is, in context Γ , update s maps tree type α to type τ'), and plural well-formedness $\Gamma \vdash^{*} \ \{\tau\} \ s \ \{\tau'\}$ (that is, in context Γ , update s maps type τ to type τ'). Several of the rules are polymorphic with respect to the arity $a \in \{1, *\}$. In addition, there is an auxiliary judgment $\Gamma \vdash_{\text{iter}} \ \{\tau\} \ s \ \{\tau'\}$ for typechecking iterations. The rules for update well-formedness are shown in Figure 8.

4.2 Discussion

Before discussing the metatheoretic properties of the type system, we discuss its behavior at a high level.

The query typechecking judgment $\Gamma \vdash e : \tau$ has its usual meaning: assuming all of the free variables x in e have val-

 $\Gamma \vdash^a \{\tau\} \ s \ \{\tau'\}$

$$\frac{\Gamma \vdash^{a} \{\tau\} \operatorname{skip} \{\tau\}}{\Gamma \vdash^{a} \{\tau\} \operatorname{skip} \{\tau\}} \frac{\Gamma \vdash^{a} \{\tau\} \operatorname{s} \{\tau'\} \Gamma \vdash^{a} \{\tau\} \operatorname{s}' \{\tau''\}}{\Gamma \vdash^{a} \{\tau\} \operatorname{skip} \{\tau\}} \frac{\Gamma \vdash^{a} \{\tau\} \operatorname{skip} \{\tau\} \Gamma \vdash^{a} \{\tau\} \operatorname{s} \{\tau\} \Gamma \vdash^{a} \{\tau\} \operatorname{s}' \{\tau\}}{\Gamma \vdash^{a} \{\tau\} \operatorname{if} e \operatorname{then} s \operatorname{else} s' \{\tau_{1} \mid \tau_{2}\}} \frac{\Gamma \vdash e: \tau}{\Gamma \vdash^{a} \{\tau_{1}\} \operatorname{let} x = e \operatorname{in} s \{\tau_{2}\}} \frac{\Gamma \vdash e: \tau}{\Gamma \vdash^{a} \{\tau_{1}\} \operatorname{let} x = e \operatorname{in} s \{\tau_{2}\}} \frac{\Gamma \vdash e: \tau}{\Gamma \vdash^{a} \{\tau_{1}\} \operatorname{let} x = e \operatorname{in} s \{\tau_{2}\}} \frac{\Gamma \vdash e: \tau}{\Gamma \vdash^{a} \{\tau_{1}\} \operatorname{let} x = e \operatorname{in} s \{\tau_{2}\}} \frac{\Gamma \vdash^{a} \{\tau\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{a} \{\tau_{1}\} \operatorname{s} \{\tau\} \operatorname{s} \{\tau'\}} \frac{\alpha < : \phi \Gamma \vdash^{1} \{\alpha\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{a} \{\tau\} \operatorname{snapshot} x \operatorname{in} s \{\tau'\}} \frac{\alpha < : \phi \Gamma \vdash^{1} \{\alpha\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{1} \{\alpha\} \phi? \operatorname{s} \{\tau\}} \frac{\alpha < : \phi \Gamma \vdash^{1} \{\alpha\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{1} \{\alpha\} \phi? \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau'\}}{\Gamma \vdash^{1} \{\alpha\} \phi? \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{1} \{\alpha\} \phi? \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau'\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau'\}} \frac{\tau_{1} < \tau_{1}' \Gamma \vdash^{*} \{\tau_{1}'\} \operatorname{s} \{\tau_{2}'\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau'\}} \frac{\tau_{1} < \tau_{1}' \Gamma \vdash^{*} \{\tau_{1}\} \operatorname{s} \{\tau_{2}'\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau'\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}}{\tau} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\}}{\tau} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}}{\Gamma \vdash^{*} \{\tau\} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\}}{\tau} \operatorname{s} \{\tau\}} \operatorname{s} \{\tau\}} \frac{\Gamma \vdash^{*} \{\tau\}}{\tau} \operatorname{s} \tau} \frac{\Gamma \vdash^{*$$

Figure 8: Update well-formedness.

ues of type $\Gamma(x)$, the result of evaluating e will have type τ . In many functional languages, and several XML update proposals, side-effects have been typechecked by treating side-effecting update operations as expressions of some "unit" type, such as (). This works fine as long as the types of values reachable from the free variables in Γ can never change; however, if the side-effects do change the types of such values, then Γ needs to be updated to take these changes into account.

One possibility is to typecheck updates using a judgment $\Gamma \vdash s$: () | Γ' ; here, Γ' is the updated context reflecting the types of the variables after update s. This approach quickly becomes unmanageable, especially if it is possible for different variables to "alias", or refer to overlapping parts of the data.

This is *not* the approach to updates that is taken in LUX. Instead, a LUX update typechecking judgment such as $\Gamma \vdash^a \{\tau\}$ s $\{\tau'\}$ assigns an update much richer type information that describes the type of the context value before and after running s. The values of variables bound in Γ cannot ever change, so their types do not need to be updated.

The most unusual rules are those involving the iter, test, and children, left/right, and insert/delete operators. The following example should help show how the rules work for these constructs. Consider the update:

```
iter {a?children {iter {b? right insert c[]}}}
```

Intuitively, this update inserts a c after every b under a top-level a. Now consider the input type $a[b[]^*, c[], b[]^*], d[]$. Clearly, the output type *should* be $a[(b[], c[])^*, c[], (b[], c[])^*], d[]$. To see why this is the case, first note that the following can be derived for any τ, τ', s :

$$\frac{\vdash^1 \{a[\tau]\} \ s \ \{a[\tau']\}}{\vdash^* \{a[\tau], d\} \ \text{iter} \ \{a?s\} \ \{a[\tau'], d\}}$$

Using the rule for children, we can see that it suffices to check that iter $\{b?$ left insert $c[]\}$ maps type $b[]^*, c[], b[]^*$ to $(b[], c[])^*, c[], (b[], c[])^*$. This is also an instance of a derivable rule

$$\frac{\vdash^1 \{b[]\} \ s \ \{\tau\}}{\vdash^* \{b[]^*, c[], b[]^*\} \ \texttt{iter} \ \{b?s\} \ \{\tau^*, c[], \tau^*\}}$$

Hence, we now need to show only that right insert c[] maps type b to b, c, which is immediate:

$$\frac{\vdash c[]:c[]}{\vdash^1 \{()\} \text{ insert } c \ \{b[],c[]\}}$$
$$\vdash^1 \{b[]\} \text{ right insert } c[] \ \{b[],c[]\}$$

4.3 Metatheory

We take for granted the following type soundness property for queries (a similar property is shown for μ XQ in Colazzo et al. [13]).

THEOREM 2 (QUERY SOUNDNESS). If $\Gamma \vdash e : \tau, \sigma \in [\![\Gamma]\!]$ and $\sigma \vdash e \Rightarrow v$ then $v \in [\![\tau]\!]$.

The corresponding result also holds for updates, by a straightforward structural induction argument:

THEOREM 3 (UPDATE SOUNDNESS).

- 1. If $\Gamma \vdash^{1} \{\alpha\} e \{\tau'\}, t \in \llbracket \alpha \rrbracket, \sigma \in \llbracket \Gamma \rrbracket, and \sigma; t \vdash e \Rightarrow^{U} v' then v' \in \llbracket \tau' \rrbracket.$
- 2. If $\Gamma \vdash^* \{\tau\} e \{\tau'\}, v \in \llbracket \tau \rrbracket, \sigma \in \llbracket \Gamma \rrbracket, and \sigma; v \vdash e \Rightarrow^{\tt U} v'$ then $v' \in \llbracket \tau' \rrbracket$.

We now consider the decision problem of checking types for queries and updates. The subsumption rule makes typechecking non-syntax-directed. However, given a context and query, we can calculate a unique syntax-directed result type, if it can be typed:

LEMMA 1. Given Γ and e, there exists at most one τ such that $\Gamma \vdash e : \tau$ has a subsumption-free derivation.

Similarly, given a context, update, and initial type, we can calculate a unique syntax-directed output type:

LEMMA 2. Given Γ , a, τ , and s, there exists at most one τ' such that $\Gamma \vdash^a {\tau} e {\tau'}$ has a subsumption-free derivation.

Since there are no functions or function types in the query and update languages, we believe it is possible to show that all uses of the subsumption rules can be permuted down in a derivation tree:

CONJECTURE 1 (SUBSUMPTION). If a query or update typing judgment is derivable, it is derivable using at most one instance of the subsumption rule, at the root of the tree. Conjecture 1 would imply the completeness of the following algorithm for typechecking a query or update against a fixed context and input/output types: first, calculate the result type using only the syntax-directed rules (without subsumption), then test whether this is a subtype of the desired output type. This algorithm is sound whether or not Conjecture 1 holds; however, if the conjecture is not true, then more work will be necessary to design a complete typechecking algorithm.

Type inference appears trickier, but is not crucial in a database update setting, because of the absence of function types and because we always know the schema of an XML database statically (assuming the database has a schema). Nevertheless, the ability to typecheck a query or update in isolation from a database schema is potentially useful, and type inference and principal types should be investigated.

5. RELATED WORK

5.1 Updates in other data models

SQL's update facility has been standard for a long time; its features are essentially the same as the update facilities in early databases such as INGRES [28]. Research on database updates has focused on efficiently maintaining or updating views of a database or checking that integrity constraints are preserved, not in increasing expressiveness. Objectoriented databases generally provide for updates via generalpurpose programming in an object-oriented language such as C++ or Java. The Lore semistructured database [25] provides a more sophisticated update language interface which permits bulk creation and deletion of edges based on subqueries. Liefke and Davidson's update language CPL+ [23] extends CPL [6], a typed language for querying data consisting of arbitrary combinations of record and monadic collection types, with path-based insert, update, and delete operations. High-level CPL+ updates were formalized by translation to a simpler core language with orthogonal operations for iteration, navigation, insertion/deletion, and replacement.

5.2 Static typing for XML processing

XML processing systems can be characterized as domainspecific transformation languages (e.g. XSLT [12], XDuce [21]), query languages (e.g. XQuery, XPath [14]), or generalpurpose programming languages with native XML support (e.g. Xtatic [17]). We will focus on only the most closely related work; Møller and Schwartzbach [26] provide a much more complete survey of type systems for XML transformation languages.

Hosoya, Vouillon, and Pierce [21] introduced XDuce, the first statically typed XML transformation language. They used a structural type system based on regular expressions in which types can be interpreted as regular tree languages and subtyping is language inclusion. Subtyping is EXPTIMEcomplete, but Hosoya et al. developed an algorithm that appears well-behaved in practice.

XQuery has a type system based on XML Schema, with regular expression types and nominal subtyping; its operational semantics and type system have been standardized in the W3C's XPath 2.0/XQuery 1.0 Formal Semantics [14]. Siméon and Wadler [27] investigated XML Schema's nominal subtyping and established some basic formal properties of validation. Their work helped clarify inconsistencies in the informal XQuery type system specification.

Colazzo, Ghelli, Manghi and Sartiani [13] considered the problem of typechecking an XQuery core language called μ XQ, from the point of view of both *result analysis* (the result of the query always matches a type) and *correctness analysis* (every subexpression which is statically guaranteed to evaluate to () is syntactically ()).

5.3 XML updates

We will briefly describe each of the XML update languages of which we are aware. Since many features are common, we will focus on the distinctive features only; Table 1 summarizes the current state of knowledge about the languages in terms of the design space in Figure 1.

Laux and Martin [22] developed an early working draft of an XML update language called XUpdate. Tatarinov, Ives, Halevy and Weld [30] considered XQuery language extensions for XML updates, focusing on the performance of updates to XML data stored in a relational database. Bruno, Le Maitre, and Murisasco [5] proposed an extension to XQuery to include *transformation operators*, which apply an update to a copy of a value. Wang, Liu and Lu [31] investigated adding XML update primitives to XML-RL, a rule-based query language with pattern-matching query operations.

UpdateX, due to Sur, Hammer and Siméon [29], is an XML update language based on XQuery. UpdateX permits arbitrary XQuery expressions as update content, but updates cannot occur within queries and the syntax of updates has some arbitrary-seeming restrictions. Benedikt, Bonifati, Flesca, and Vyas [1, 2] studied the semantics of UpdateX, and developed analyses that support various optimizations.

Calcagno, Gardner and Zarfaty [7] investigated XML updates in the setting of *context logic*, a logic of "trees with holes". They considered a Hoare logic for reasoning about low-level update operations comparable to those provided by the the DOM interface. This work currently applies only to unordered, deterministic trees, a simplification of full XML.

Ghelli, Re and Siméon [19] have developed XQuery[!], an XML query language with side-effects. In XQuery[!], update statements are allowed as expressions returning (). Side-effects in expressions are treated lazily; they are collected into a *pending update list* and by default only performed at the end of the query; the **snap** operator can be used to force pending updates to be performed earlier. Ghelli, Rose, and Siméon [18] have developed a commutativity analysis for XQuery[!].

The W3C XQuery Update Facility [9] has been under development for several years. The current version is similar to XQuery[!]; the two main differences are the absence of the **snap** operator and the presence of transformations.

LiXQuery⁺ is an update language introduced by Hidders, Paredaens, and Vercammen [20] in order to study the expressive power of XQuery-based update languages.

XQueryP, due to Chamberlin, Carey, Florescu, Kossman, and Robie [8], is an extension to XQuery that includes both updates and traditional imperative programming features, including block structure, assignment, an **updating** keyword for declaring impure functions, and **while**-loops. XQueryP also includes an **atomic** keyword which groups operations into atomic blocks (transactions).

The design goals and constraints of many of these proposals differ from those that motivated this work. LUX is not meant to be a full-fledged programming language for mutable XML data. Instead, it is meant to play a role for XML and XQuery similar to that SQL's update facilities relative to relational databases and SQL. It seems expressive enough for common cases (including many of the relational use cases in the XQuery Update Facility Requirements [10]) while remaining simple and statically typecheckable.

6. EXTENSIONS AND FUTURE WORK

6.1 **Recursive types and nominal subtyping**

Our type system is weaker than other XML type systems in several respects. It does not include type variables or recursive type definitions, so although it permits *horizontal* iteration using Kleene star, *vertical* recursion is not supported. The type system also employs structural subtyping on unannotated XML data, in contrast to the XML Schema/XQuery type system which employs nominal subtyping on validated XML documents in which each element is tagged with a type name according to the XML Schema validation algorithm.

Although these are significant restrictions relative to full DTDs, XML Schemas, or XDuce regular expression types, we believe they are a reasonable given that LUX is meant to be used for database-style updates, not queries, stylesheet transformations, or general-purpose programming with XML. Many DTDs and XML Schemas encountered in database applications of XML are nonrecursive and "shallow" [11, 3].

On the other hand, there are use cases for XML updates involving recursion and XML Schemas [24] which we would like to handle. Extending the system to include recursive queries and recursive regular expression types appears to be straightforward; adding recursive *updates* may not be. Adapting our approach to updates to be compatible with XML Schema-validated documents and the XQuery nominal subtyping system seems difficult.

6.2 Transformations

Some of the existing update proposals include a facility for running an update operation within an XQuery expression in a side-effect-free way. Such a facility can easily be added using LUX updates:

 $e ::= \cdots \mid \texttt{transform} \ e \ \texttt{by} \ s$

with static and dynamic semantics given by

$$\frac{\sigma \vdash e \Rightarrow v \quad \sigma; v \vdash s \Rightarrow^{\mathrm{U}} v'}{\sigma \vdash \operatorname{transform} e \text{ by } s \Rightarrow v'} \quad \frac{\Gamma \vdash e: \tau_1 \quad \Gamma \vdash^* \{\tau_1\} \ s \ \{\tau_2\}}{\Gamma \vdash \operatorname{transform} e \text{ by } s: \tau_2}$$

6.3 Correctness analysis

Colazzo et al. [13] studied both "result analysis" (determining the result type of a query given types for its input variables) and "correctness analysis" (ensuring that no subexpression of the query is statically empty, except for ()). They introduced a type system that statically performs both result and correctness analysis; moreover, the correctness analysis is *complete*: it precisely characterizes queries as either correct or incorrect. This problem is decidable for μXQ since it lacks recursion.

Correctness analysis is potentially very useful in debugging queries, since often a bug manifests itself as an empty subquery; such bugs are generally not caught by result analysis since the empty sequence is a subtype of every type.

$$\begin{array}{rcl} \phi?u_1;\phi?u_2 &\equiv& \phi?(u_1;u_2)\\ d[u_1];d[u_2] &\equiv& d[u_1;u_2]\\ \\ \text{snapshot } x \text{ in snapshot } y \text{ in } e &\equiv& \text{snapshot } x \text{ in } e[x/y]\\ \\ &\text{snapshot } x \text{ in } e &\equiv& e \quad (x \notin FV(e))\\ \\ &s \text{delete} &\equiv& \text{delete}\\ d[\text{if } c \text{ then } s_1 \text{ else } s_2] &\equiv& \text{if } c \text{ then } d[s_1] \text{ else } d[s_2] \end{array}$$

Figure 9: Some equational laws for updates.

Colazzo et al.'s type system for query correctness analysis can clearly also be used in LUX updates to help avoid bugs in subqueries. Moreover, it seems reasonable to extend their concept of correctness to updates by judging an update to be correct only if every sub-statement that is statically guaranteed to have no effect on the database is literally a no-op (skip). Updates can easily be incorrect: for example, if a test update bool?s is only run against data of type string, then s will never be executed; also s will never be typechecked, so might contain nonsense. Such "dead" code is likely a bug in the update, and should be flagged as a warning to the programmer. It would be very interesting to develop a type system for update correctness, especially a complete one.

6.4 Update optimization

Liefke and Davidson [23] investigated a number of optimization techniques for CPL+, including rewriting updates using equational laws that provably improve performance (relative to an appropriate cost model), and transforming query-based updates to more efficient in-place updates ("deltafication"). These techniques should also be applicable to LUX. For example, the equational laws shown in Figure 9 are clearly valid for LUX.

6.5 Designing a high-level update language

In the first half of this paper we motivated and informally described a high-level update language that can be systematically translated to the core Lux language, just as full XQuery is defined by translation to a much smaller core language. In addition we have shown how a number of examples can be so translated. At present our proof-of-concept implementation provides high-level updates similar to those in Section 2; they are expanded to core Lux updates during parsing. Figure 10 illustrates a few "compilation" rules which are used in the current system (extending the abbreviations of Figure 5). However, this design is subject to change and is provided as an example only.

7. CONCLUSIONS

In this paper, we have introduced a simple approach to database-style updates for XML, called LUX (Lightweight Updates for XML). Our approach is inspired in large part by the update language CPL+ for nested relational data introduced by Liefke and Davidson. Our approach, like CPL+, factors complicated operations similar to SQL's INSERT, DELETE, and UPDATE statements into a small, orthogonal set of operations which perform basic updates, navigate, examine the local database structure, or combine updates. LUX is carefully designed so that each update's side-effects are confined

update x in n/p by u	\longrightarrow	update in n by
		update x in p by u
update n/p by u	\longrightarrow	update in n by
		update p by u
insert before n value e	\longrightarrow	update n by insert $\stackrel{\leftarrow}{-} e$
insert after n value e	\longrightarrow	update n by insert $^{ ightarrow} e$
insert as first into n	\longrightarrow	update in n by insert ^{\leftarrow} e
value e		
insert as last into n	\longrightarrow	update in n by insert $^{\rightarrow} e$
value e		
insert $\dots n/p$	\longrightarrow	update in n by
value e		insert p
		value e
delete n/p	\longrightarrow	update in n by
		delete p
delete $x \text{ from } n/p$	\longrightarrow	update x in n by
where e		if c then delete p

Figure 10: Some "compilation" rules for high-level updates

to a particular subtree, so that the result of an iterative update is independent of the order in which the updates are performed. Variables are immutable; updates can only mutate the database state. Side-effecting updates are syntactically distinguished from purely functional queries. Node identifiers are not available for use in defining updates.

These design decisions make it possible to give updates a simple, deterministic semantics. In contrast, the semantics of more sophisticated update languages involve two passes, one to collect atomic updates evaluated against a snapshot of the database, and another to perform the updates. In addition, it is possible to provide a type system for the core update language which can be used to check that an update has desired input-output behavior. We believe that this is the first static type system to be developed for an XML update language.

Besides proving type soundness, we have described a simple, sound typechecking algorithm that computes an upper bound on the type of a query or update, then checks whether the bound is included in the desired result type. We believe this algorithm is complete, but have not fully verified this conjecture.

Lux is work in progress, and in the near future we plan to prove Conjecture 1, improve the design of the high-level interface, refine the high-level language design, and experiment with large-scale examples (including some of the W3C's use cases [24]) inside an XML database management system. Additional possible areas for future work include static typechecking relative to XQuery's type system; considering recursive types and updates; and developing a type system capturing correctness.

Acknowledgments. Thanks to Peter Buneman, Irini Fundulaki, Sam Lindley, and Phil Wadler for comments and discussions, and to the anonymous reviewers for suggesting several clarifications.

8. **REFERENCES**

[1] Michael Benedikt, Angela Bonifati, Sergio Flesca, and Avinash Vyas. Adding updates to XQuery: Semantics, optimization, and static analysis. In Daniela Florescu and Hamid Pirahesh, editors, XIME-P, 2005.

- [2] Michael Benedikt, Angela Bonifati, Sergio Flesca, and Avinash Vyas. Verification of tree updates for optimization. In Kousha Etessami and Sriram K. Rajamani, editors, CAV, volume 3576 of Lecture Notes in Computer Science, pages 379–393. Springer, 2005.
- [3] Geert Jan Bex, Frank Neven, and Jan Van den Bussche. DTDs versus XML schema: a practical study. In WebDB '04: Proceedings of the 7th International Workshop on the Web and Databases, pages 79–84, New York, NY, USA, 2004. ACM Press.
- [4] Scott Boag, Don Chamberlin, Mary F. Fernndez, Daniela Florescu, Jonathan Robie, and Jérôme Siméon. XQuery 1.0: An XML query language. W3C Candidate Recommendation, June 2006. http://www.w3.org/TR/xquery.
- [5] Emmanuel Bruno, Jacques Le Maitre, and Elisabeth Murisasco. Extending XQuery with transformation operators. In *DocEng '03: Proceedings of the 2003 ACM Symposium on Document Engineering*, pages 1–8, New York, NY, USA, 2003. ACM Press.
- [6] Peter Buneman, Shamim A. Naqvi, Val Tannen, and Limsoon Wong. Principles of programming with complex objects and collection types. *Theor. Comp. Sci.*, 149(1):3–48, 1995.
- [7] Cristiano Calcagno, Philippa Gardner, and Uri Zarfaty. Context logic and tree update. SIGPLAN Not., 40(1):271–282, 2005.
- [8] Don Chamberlin, Mike Carey, Daniela Florescu, Donald Kossmann, and Jonathan Robie. XQueryP: Programming with XQuery. In *XIME-P*, 2006.
- [9] Don Chamberlin, Daniela Florescu, and Jonathan Robie. XQuery update facility. W3C Working Draft, July 2006. http://www.w3c.org/TR/xqupdate/.
- [10] Don Chamberlin and Jonathan Robie. XQuery update facility requirements. W3C Working Draft, June 2005. http://www.w3.org/TR/xquery-update-requirements/.
- [11] Byron Choi. What are real DTDs like? In WebDB, pages 43–48, 2002.
- [12] J. Clark. XSL transformations (XSLT). W3C Recommendation, November 1999. http://www.w3.org/TR/xslt.
- [13] Dario Colazzo, Giorgio Ghelli, Paolo Manghi, and Carlo Sartiani. Types for path correctness of XML queries. In *ICFP '04: Proceedings of the ninth ACM SIGPLAN International Conference on Functional Programming*, pages 126–137, New York, NY, USA, 2004. ACM Press.
- [14] D. Chamberlin et al. XQuery 1.0 and XPath 2.0 formal semantics. W3C Candidate Recommendation, June 2006.

http://www.w3.org/TR/xquery-semantics/.

- [15] David C. Fallside (Ed). XML Schema Part 0: Primer. W3C Recommendation, October 2004. http://www.w3.org/TR/xmlschema-0.
- [16] Mary F. Fernandez, Jérôme Siméon, and Philip Wadler. A semi-monad for semi-structured data. In *ICDT '01: Proceedings of the 8th International Conference on Database Theory*, pages 263–300, London, UK, 2001. Springer-Verlag.

- [17] Vladimir Gapeyev, François Garillot, and Benjamin C. Pierce. Statically typed document transformation: An Xtatic experience. In Giuseppe Castagna and Mukund Raghavachari, editors, *PLAN-X*, pages 2–13. BRICS, Department of Computer Science, University of Aarhus, 2006.
- [18] G. Ghelli, K. Rose, and J. Siméon. Commutativity analysis in XML update languages. In *ICDT*, 2007. To appear.
- [19] Giorgio Ghelli, Christopher Re, and Jérôme Siméon. XQuery!: An XML query language with side effects. In Torsten Grust, Hagen Höpfner, Arantza Illarramendi, Stefan Jablonski, Marco Mesiti, Sascha Müller, Paula-Lavinia Patranjan, Kai-Uwe Sattler, Myra Spiliopoulou, and Jef Wijsen, editors, *EDBT* Workshops, volume 4254 of Lecture Notes in Computer Science, pages 178–191. Springer, 2006.
- [20] Jan Hidders, Jan Paredaens, and Roel Vercammen. On the expressive power of XQuery-based update languages. In Sihem Amer-Yahia, Zohra Bellahsene, Ela Hunt, Rainer Unland, and Jeffrey Xu Yu, editors, XSym, volume 4156 of Lecture Notes in Computer Science, pages 92–106. Springer, 2006.
- [21] Haruo Hosoya, Jérôme Vouillon, and Benjamin C. Pierce. Regular expression types for XML. ACM Trans. Program. Lang. Syst., 27(1):46–90, 2005.
- [22] Andreas Laux and Lars Martin. XUpdate XML update language. http://xmldb-org.sourceforge.net/xupdate/xupdate-wd.html, September 2000. Work in progress.
- [23] Hartmut Liefke and Susan B. Davidson. Specifying updates in biomedical databases. In SSDBM, pages 44–53, 1999.
- [24] Ioana Manolescu and Jonathan Robie. XQuery update facility use cases. W3C Working Draft, May 2006. http://www.w3.org/TR/xqupdateusecases.
- [25] Jason McHugh, Serge Abiteboul, Roy Goldman, Dallan Quass, and Jennifer Widom. Lore: A database management system for semistructured data. *SIGMOD Record*, 26(3):54–66, 1997.
- [26] Anders Møller and Michael I. Schwartzbach. The design space of type checkers for XML transformation languages. In Thomas Eiter and Leonid Libkin, editors, *ICDT*, volume 3363 of *Lecture Notes in Computer Science*, pages 17–36. Springer, 2005.
- [27] Jérôme Siméon and Philip Wadler. The essence of XML. SIGPLAN Not., 38(1):1–13, 2003.
- [28] Michael Stonebraker, Gerald Held, Eugene Wong, and Peter Kreps. The design and implementation of INGRES. ACM Trans. Database Syst., 1(3):189–222, 1976.
- [29] Gargi Sur, Joachim Hammer, and Jérôme Siméon. UpdateX - an XQuery-based language for processing updates in XML. In *PLAN-X*, 2004.
- [30] Igor Tatarinov, Zachary G. Ives, Alon Y. Halevy, and Daniel S. Weld. Updating XML. In SIGMOD Conference, pages 413–424, 2001.
- [31] Guoren Wang, Mengchi Liu, and Li Lu. Extending XML-RL with update. In *IDEAS*, pages 66–75. IEEE Computer Society, 2003.