

# Propagating XML View Updates to a Relational Database \*

Vanessa P. Braganholo<sup>1</sup>  
vanessa@inf.ufrgs.br

Susan B. Davidson<sup>2</sup>  
susan@cis.upenn.edu

Carlos A. Heuser<sup>1</sup>  
heuser@inf.ufrgs.br

<sup>1</sup>Instituto de Informática  
Universidade Federal do Rio Grande do Sul - UFRGS  
Brazil

<sup>2</sup>Department of Computer and Information Science  
University of Pennsylvania  
USA

Technical Report Number RP-341  
Universidade Federal do Rio Grande do Sul - UFRGS  
Instituto de Informática  
Porto Alegre - RS - Brazil

February 2004

## Abstract

This paper addresses the question of updating relational databases through XML views. Using a notion of *query trees* to capture the notions of selection, projection, nesting, grouping, and heterogeneous sets found throughout most XML query languages, we show how XML views expressed using query trees can be mapped to a set of corresponding relational views. We then show how updates on the XML view are mapped to updates on the corresponding relational views. Existing work on updating relational views can then be leveraged to determine whether or not the relational views are updatable with respect to the relational updates, and if so, to translate the updates to the underlying relational database.

## 1 Introduction

XML is frequently used as an interface to relational databases. In this scenario, XML documents (or views) are exported from relational databases and published, exchanged, or used as the internal representation in user applications. This fact has stimulated much research in exporting and querying relational data as XML views [23, 33, 32, 16]. However, the problem of updating a relational database through an XML view has not received as much attention: Given an update on an XML view of a relational database, how should it be translated to updates on the relational database? Since the problem of updates through relational views has been studied for more than 20 years by the database community, it would be good to use all that work to solve the new arising problem of updates through XML views. Specifically, is there a way to leverage existing work on updating through relational views to map view updates to the underlying relational database?

In the relational case, attention has focused on updates through select-project-join views since they represent a common form of view that can be easily reasoned about using key and foreign key information. Similarly, we focus on a common form of XML views that allows nesting, composed attributes, heterogeneous sets and repeated elements. An example of such a view is shown in figure 2, which was defined over the database of figure 1. In this XML view, *books* are nested under the *products* node, and the *address* node composes attributes in a nested record format. The *products* node is composed of tuples of two different types, *book* and *dvd*.

We represent XML view expressions as *query trees*. Query trees can be thought of as the intermediate representation of a query expressed by some high-level XML query language, and provide a language independent framework in which to study how to map updates to an underlying relational database. They are expressive enough to capture the XML views that we have encountered in practice, yet are simple to understand and manipulate. Their expressive power is equivalent to that of DB2 DAD files [17]. Throughout the paper, we will use the term “XML view” to mean those produced by query trees.

---

\*Research partially supported by CNPq as well as NSF DBI-9975206.

Vendor(vendorId, vendorName, url, state, country)  
 Book(isbn, title, publisher, year)  
 DVD(asin, title, genre, nrDisks)  
 Sell-Book(vendorId, isbn, price)  
 – foreign key(vendorId) references Vendor  
 – foreign key(isbn) references Book  
 Sell-DVD(vendorId, asin, price)  
 – foreign key(vendorId) references Vendor  
 – foreign key(asin) references DVD

Figure 1: Sample database

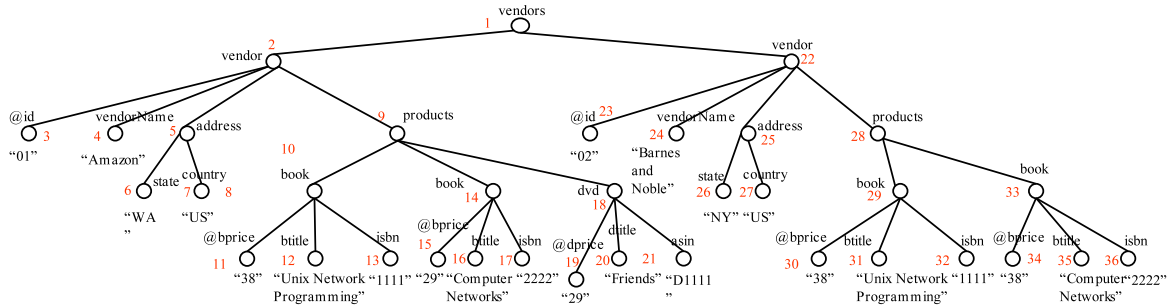


Figure 2: View 2: vendors, books and dvds

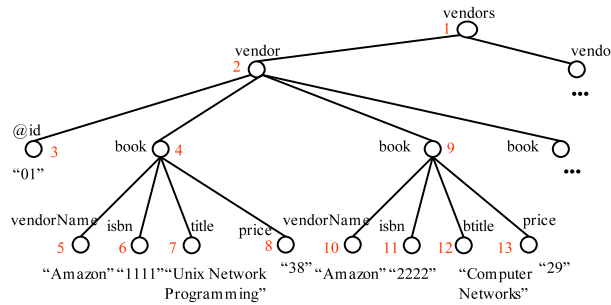


Figure 3: View 1: books and vendors

The strategy we adopt is to map an XML view to a set of underlying relational views. Similarly, we map an update against the XML view to a set of updates against the underlying relational views. It is then possible to use any existing technique on updates through relational views to both translate the updates to the underlying relational database and to answer the question of whether or not the XML view is updatable with respect to the update.

This strategy is similar to that adopted in [13] for XML views constructed using the nested relational algebra (NRA), however, our view and update language are far more general. In particular, nested relations cannot handle heterogeneity. Thus, the NRA is capable of representing the XML view of Figure 3 but not that of Figure 2, and maps an XML view to exactly one underlying relational view.

The outline and contributions of this paper are as follows:

- Section 2 defines query trees, their abstract types, and the resulting XML view DTD.
- Section 3 presents the algorithm for mapping an XML view to a set of underlying relational views, and proves its correctness.
- Section 4.1 defines a simple XML update language and algorithms to detect whether or not an update is correct with respect to the XML view DTD.
- Section 4.2 gives an algorithm for mapping insertions, modifications and deletions on XML views to updates on the underlying relational views, and Section 4.3 proves its correctness.

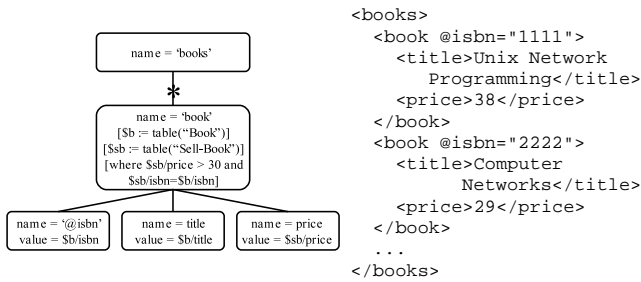


Figure 4: Example of query tree and its resulting XML view

- Section 4.4 illustrates our approach by showing how to use the techniques of [21] detect if an XML view is updatable with respect to a given update.
- Section 5 shows how query trees can be easily extended to add grouping capabilities.
- Section 6 discusses the expressive power of our language, and evaluates our technique with respect to existing proposals on extracting XML views of relational databases.

Related work can be found in section 7. We conclude in section 8 with a discussion of future work.

## 2 Query Trees

Query trees are used as a representation of the XML view extraction query. We use this abstract representation rather than an XML query language syntax for several reasons: First, reasoning about updates and the updatability of an XML view is performed at this level. Second, they are easy to understand yet expressive enough to capture several important aspects of XQuery such as nesting, composed attributes, and heterogeneous sets.<sup>1</sup> They can therefore be thought of as the intermediate processing form for a subset of many different XML query languages. For example, we have developed an implementation of our technique which uses a subset of XQuery as the top-level language [14].

After defining query trees, we introduce a notion which will be used to describe the mapping to relational queries, the abstract type of a query tree node. We use this notion of typing to define the semantics of query trees, and then present their result type DTD.

### 2.1 Query Trees Defined

An example of a query tree can be found in Figure 4, which retrieves books that are sold for prices greater than \$30. The query tree resembles the structure of the resulting XML view. The root of the tree corresponds to the root element of the result. Leaf nodes correspond to attributes of relational tables, and interior nodes whose incoming edges are starred capture repeating elements. The result of this query is also presented in figure 4.

Query trees are very similar to the *view forests* of [23] and *schema-tree queries* presented in [11]. The difference is that, instead of annotating all nodes with the relational queries that are used to build the content model of a given node, we annotate interior nodes in the tree using only the selection criteria (not the entire relational query). An annotation can be a *source* annotation or a *where* annotation. Source annotations bind variables to relational tables, and *where* annotations impose restrictions on the relational tables making use of the variables that were bound to the tables.

In the definitions that follow, we assume that  $\mathcal{D}$  is a relational database over which the XML view is being defined.  $\mathbb{T}$  is the set of table names of  $\mathcal{D}$ .  $\mathbb{A}_T$  is the set of attributes of a given table  $T \in \mathbb{T}$ .

**Definition 2.1** A query tree defined over a database  $\mathcal{D}$  is a tree with a set of nodes  $\mathbb{N}$  and a set of edges  $\mathbb{E}$  in which: **Edges** are simple or starred ("\*-edge"). An edge is simple if, in the corresponding XML instance, the child node appears exactly once in the context of the parent node, and starred otherwise. **Nodes** are as follows:

1. All nodes have a name that represents the tag name of the XML element associated with this node in the resulting XML view.
2. Leaf nodes have a value (to be defined). Names of leaf nodes that start with "@" are considered to be XML attributes.

<sup>1</sup>They can also capture grouping, and we present such extension in section 5.

3. Starred nodes (*nodes whose incoming edge is starred*) may have one or more source annotations and zero or more where annotations (to be defined).

Since we map XML views to relational views, nodes with the same name in the query tree may cause ambiguities in the mapping. This problem can easily be solved by associating with each node name a number corresponding to its position in the query tree, and using it internally in the mapping. For simplicity, in this paper we will ignore this problem and use unique names for nodes in the query trees.

Returning to the example in Figure 4, there is a \*-edge from the root (named *books*) to its child named *book*, indicating that in the corresponding XML instance there may be several *book* subelements of *books*. There is a simple edge from the node named *book* to the node named *title*, indicating that there is a single *title* subelement of *book*. The node named *@isbn* will be mapped to an XML attribute instead of an element.

Before giving an example of how values are associated with nodes, we define *source* and *where* annotations on nodes of a query tree.

**Definition 2.2** A source annotation  $s$  within a starred node  $n$  is of the form  $[\$x := \text{table}(T)]$ , where  $\$x$  denotes a variable and  $T \in \mathbb{T}$  is a relational table. We say that  $\$x$  is bound to  $T$  by  $s$ .

**Definition 2.3** A where annotation on a starred node  $n$  is of the form  $[\text{where } \$x_1/A_1 \text{ op } Z_1 \text{ AND } \dots \text{ AND } \$x_k/A_k \text{ op } Z_k]$ ,  $k \geq 1$ , where  $A_i \in \mathbb{A}_{T_i}$  and  $\$x_i$  is bound to  $T_i$  by a source annotation on  $n$  or some ancestor of  $n$ . The operator  $\text{op}$  is a comparison operator  $\{=, \neq, >, <, \leq, \geq\}$ .  $Z_i$  is either a literal (integer, string, etc.) or an expression of the form  $\$y/B$ , where  $B \in \mathbb{A}_T$  and  $\$y$  is bound to  $T$  by a source annotation on  $n$  or some ancestor of  $n$ .

**Definition 2.4** The value of a node  $n$  is of form  $\$x/A$ , where  $A \in \mathbb{A}_T$  and  $\$x$  is bound to table  $T$  by a source annotation on  $n$  or some ancestor of  $n$ .

In Figure 4, the node *book* has source annotations and where annotations. The source annotations bind variable  $\$b$  to the relational table *Book*, and variable  $\$sb$  to the relational table *Sell-Book*. The where annotations restrict the books that appear in the view to those with price greater than \$30, and specify the join condition of tables *Book* and *Sell-Book*. The value of the node *@isbn* is specified as  $\$b/\text{isbn}$ , indicating that the content of the XML view attribute *isbn* will be generated using attribute *isbn* of the table *Book*.

A more complex example of a query tree can be found in Figure 5 (ignore for now the types  $\tau$  associated with nodes). This query tree retrieves *vendors*, and for each *vendor*, its *@id*, *vendorName*, *address* and a set of *books* and *dvds* within *products*. The root *vendors* has a set of *vendor* child nodes (\*-edge). The *vendor* node is annotated with a binding for  $\$v$  (to table *Vendor*), and has several children at the end of simple edges (*@id*, *vendorName*, and *address*). The value of its *id* attribute is specified by the path  $\$v/\text{vendorId}$ , and that of *vendorName* is specified by the path  $\$v/\text{vendorName}$ . The node *address* is more complex, and is composed of *state* and *country* subelements.

The node *products* has two \*-edge children, *book* and *dvd*. Source annotations of the *book* node include bindings for  $\$b$  (*Book*) and  $\$sb$  (*Sell-Book*) and its where annotations connect tuples in *Sell-Book* to tuples in *Book*, and tuples in *Sell-Book* with tuples in *Vendor* (join conditions). Node *dvd* has source annotations for  $\$d$  (*DVD*) and  $\$sd$  (*Sell-DVD*). Its where annotation connects tuples in *Sell-DVD* to tuples in *DVD* and tuples in *Sell-DVD* with tuples in *Vendor*. The result of this query tree is View 2, shown in Figure 2.

From now on, we assume that a query tree is *non-empty*, i.e. that it consists of more than a root node.

## 2.2 Abstract Types

In our mapping strategy, it will be important to recognize nodes that play certain roles in a query tree. In particular, we identify five abstract types of nodes:  $\tau$ ,  $\tau_T$ ,  $\tau_N$ ,  $\tau_C$  and  $\tau_S$ . We call them *abstract types* to distinguish them from the type or DTD of the XML view elements.

Nodes in the query tree are assigned abstract types as follows:

1. The root has abstract type  $\tau$ .
2. Each leaf has abstract type  $\tau_S$  (Simple).
3. Each non-leaf node with an incoming simple edge has abstract type  $\tau_C$  (Complex).
4. Each starred node which is either a leaf node or whose subtree has only simple edges has an abstract type of  $\tau_N$  (Nested).

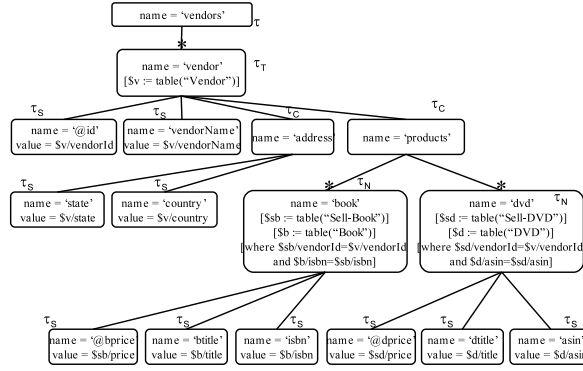


Figure 5: Query tree for View 2

5. All other starred nodes have abstract type  $\tau_T$  (Tree).

Note that each node has exactly one type unless it is a starred leaf node, in which case it has types  $\tau_S$  and  $\tau_N$ .

As an example of this abstract typing, consider the query tree in figure 5, which shows the type of each of its nodes. Since *book* and *dvd* are repeating nodes whose descendants are non-repeating nodes, their types are  $\tau_N$  rather than  $\tau_T$ .

The motivation behind abstract types is as follows. To map updates in the XML view to updates in the underlying relational database, we must be able to identify a mapping from the attribute of a tuple in the relational database to an element or attribute in the XML view. Ideally, this mapping is 1:1, i.e. each attribute of a tuple occurs at most once in the XML view and can therefore be updated without introducing side-effects into the view. In general, however, it may be a 1:n mapping. The class of views allowed by our query trees and its associated abstract type views captures this mapping intrinsically.

Specifically:

- $\tau_T/\tau_N$  identifies potential tuples in the underlying relational database. Nodes of type  $\tau_T/\tau_N$  are mapped to tuples, and the node itself serves as a tuple delimiter. A node of type  $\tau_T$  may have children of type  $\tau_T$ , i.e. nesting is allowed.
- $\tau_S$  identifies relational attributes (columns). A node of type  $\tau_S$  must have a node of type  $\tau_T$  or  $\tau_N$  as its ancestor. Starred leaf nodes are an exception to this rule: they need not to have such ancestor.
- $\tau_C$  identifies complex XML elements. Since they do not carry a value, this type of node is not mapped to anything in the relational model. Nodes of type  $\tau_C$  are present in our model to allow more flexible XML views, but are not important in the mapping process.

We call the XML views produced by query trees and their associated abstract types *well-behaved* because, as we will show in the next section, they can be easily mapped to a set of corresponding relational views. However, before turning to the mapping we prove two facts about query trees that will be used throughout the paper.

**Proposition 2.1** *There is at least one  $\tau_N$  node in the abstract type of a query tree  $qt$ .*

*Proof:* Since query trees are assumed to be non-empty,  $qt$  must have at least one leaf. This means that  $qt$  must have at least one starred node  $n$ , since the leaf node has a value which involves at least one variable which must be defined in some source annotation attached to a starred node. Since the tree is finite, at least one of these starred nodes is either a leaf node or has a subtree of simple edges, i.e. the starred node is a  $\tau_N$  node. ■

**Proposition 2.2** *There is at most one  $\tau_N$  node along any path from a leaf to the root in the abstract type of a query tree  $qt$ .*

*Proof:* Suppose there are two  $\tau_N$  nodes,  $n_1$  and  $n_2$ , along the path from some leaf to the root of  $qt$ . Without loss of generality, assume that  $n_1$  is the ancestor of  $n_2$ . By definition of  $\tau_N$ ,  $n_2$  must be a starred node. Therefore  $n_1$  has a \*-edge in its subtree, a contradiction. ■

We will refer to the abstract type of an element by the abstract type that was used to generate it followed by the element name. As an example, the abstract type of the element *dvd* in Figure 5 is referred to as  $\tau_N(dvd)$ , and its type (DTD) is  $\langle !ELEMENT\ dvd\ (dtitle,\ asin)\rangle$ .

```

eval(qt, d)
evaluate(root(qt, d))

evaluate(n, d)
{Assume a node type has functions abstract_type(n), name(n), value(n), children(n), sources(n), and where(n) (with the obvious meanings.)}
Let bindings{ } be a hash array of bindings of variable attributes to values, initially empty.
case abstract_type(n)
   $\tau_C$ : buildElement(n)
   $\tau_T$  |  $\tau_N$ : table(n)
   $\tau_S$ : print "<name(n)>value(n)</name(n)>"
end case

buildElement(n)
let tag = "name(n)"
for each attribute c in children(n) do
  add "name(c) = value(c)" to tag
end for
print "< tag >"
for each non-attribute c in children(n) do
  evaluate(c)
end for
print "</name(n)>"

table(n)
let w be a list of conditions in sources(n)
for each w[i] do
  if w[i] involves a variable v in bindings{ } then
    substitute the value binding{v} for v
  end if
end for
calculate the set B of all bindings for variables in sources(n) that makes the conjunction of the modified w[i]'s true, using d
for each b in B do
  add b to bindings{ }
  buildElement(n)
  remove b from bindings{ }
end for

```

**Algorithm 1:** Eval algorithm

## 2.3 Semantics of Query Trees

The semantics of a query tree follows the abstract type of its nodes, and can be found in algorithm 1. The algorithm constructs the XML view resulting from a query tree recursively, and starts with  $n$  being the root of the query tree. The basic idea is that the source and where annotations in each starred node  $n$  are evaluated, producing a set of tuples. The algorithm then iterates over these tuples, generating one element corresponding to  $n$  in the output for each of these tuples and evaluating the children of  $n$  once for each tuple.

The *bindings*{ } hash array keeps the values of variables, taken from the underlying relational database. We assume that values in *bindings*{ } are represented as  $\$x/A = 1$ ,  $\$x/B = 2$ , where  $\$x$  is a variable bound to a relational table  $T$ ,  $A$  and  $B$  are the attributes of  $T$  and 1 and 2 are the values of attributes  $A$  and  $B$  in the current tuple of  $T$ .

## 2.4 DTD of a Query Tree

Query tree views defined over a relational database have a well-defined schema (DTD) that is easily derived from the tree. Given a query tree, its DTD is generated as follows:

1. For each attribute leaf node named  $@A$  with parent named  $E$ , create an attribute declaration `<!ATTLIST  $E$   $@A$  CDATA #REQUIRED>`
2. For each non-attribute leaf node named  $E$ , create an element declaration `<!ELEMENT  $E$  (#PCDATA)>`
3. For each non-leaf node named  $E$ , create an element declaration `<!ELEMENT  $E$  ( $E_1, \dots, E_k, E_{k+1}^*, \dots, E_n^*)>$` , where  $E_1, \dots, E_k$  are non-attribute child nodes of  $E$  connected by a simple edge, and  $E_{k+1}^*, \dots, E_n^*$  are child nodes of  $E$  connected by a  $*$ -edge. In case  $n = 0$ , then create an element declaration `<!ELEMENT  $E$  EMPTY>`

As an example, the DTD of the view produced by the query tree shown in figure 5 is:

```

<!ELEMENT vendors (vendor*)>
<!ELEMENT vendor (vendorName, address, products)>
<!ATTLIST vendor id CDATA #REQUIRED>
<!ELEMENT vendorName (#PCDATA)>
<!ELEMENT address (state, country)>

```

```

<!ELEMENT state (#PCDATA)>
<!ELEMENT country (#PCDATA)>
<!ELEMENT products (book*,dvd*)>
<!ELEMENT book (btitle, isbn)>
<!ATTLIST book bprice CDATA #REQUIRED>
<!ELEMENT btitle (#PCDATA)>
<!ELEMENT isbn (#PCDATA)>
<!ELEMENT dvd (dtitle, asin)>
<!ATTLIST dvd dprice CDATA #REQUIRED>
<!ELEMENT asin (#PCDATA)>
<!ELEMENT dtitle (#PCDATA)>

```

Note that all (#PCDATA) elements are required. When the value of a relational attribute is null, we produce an element with a distinguished null value.

### 3 Mapping to Relational Views

In our approach, updates over an XML view are translated to SQL update statements on a set of corresponding relational view expressions. Existing techniques such as [21, 25, 29, 9, 38] can then be used to accept, reject or modify the proposed SQL updates. In this section, we discuss how an XML view constructed by a query tree is mapped to a set of corresponding relational view expressions.

**Map.** Given a query tree  $qt$  with only one  $\tau_N$  node, the corresponding SQL view statement is generated as follows. Join together all tables found in source annotations (called *source tables*) in a given node  $n$  in  $qt$ , using the where annotations that correspond to joins on source tables in  $n$  as inner join conditions. If no such join condition is found then use “true” (e.g. 1=1) as the join condition, resulting in a cartesian product. Call these expressions *source join expressions*. Use the hierarchy implied by the query tree to left outer join source join expressions in an ancestor-descendant direction, so that ancestors with no children still appear in the view. The conditions for the outer joins are captured as follows: If node  $a$  is an ancestor of  $n$  and a where annotation in  $n$  specifies a join condition on a table in  $n$  with a table in  $a$ , then use this annotation as the join condition for the outer join. Similar to inner joins, if no condition for the outer join is found, then use “true” as the join condition so that if the inner relation is empty, the tuples of the outer will still appear. Use the remaining where annotations (the ones that were not used as inner or outer join conditions) in an SQL where-clause and project the values of leaf nodes. The resulting SQL view statement represents an unnested version of the XML view.

Source join expressions are as follows:

```

<source table> AS <source variable> INNER JOIN
<source table> AS <source variable> INNER JOIN ...
ON <inner joincond>

```

The resulting SQL expression is:

```

SELECT <leaf value> AS <leaf name>, ...,
       <leaf value> AS <leaf name>
FROM (<source join expression> LEFT JOIN
      <source join expression> ON <outer joincond>) LEFT JOIN ...
WHERE <remaining "where" annotation> AND ...
      AND <remaining "where" annotation>

```

For example, the relational view corresponding to the query tree in Figure 4 is:

```

SELECT b.isbn AS isbn, b.title AS title, sb.price AS price
FROM (Book AS b INNER JOIN Sell-Book AS sb
      ON sb.isbn=b.isbn) WHERE sb.price > 30

```

The mapping algorithm is presented in details by algorithm 2. The auxiliary functions used in this algorithm have obvious meanings. The one that is not so obvious is function  $variable(n)$ . It returns the variable that was used in the value of a leaf node, without the \$. For example, if the value of node  $n$  is  $\$x/A$ , then  $variable(n)$  returns  $x$ . When the parameter is a source annotation  $s$ , then the function returns the variable referenced in this source annotation, without the \$ (e.g. with  $s = \$x$  in  $table("X")$ , function  $variable(s)$  returns  $x$ ). Function  $attribute(n)$  returns the relational attribute that was used to specify the value of a leaf node. Using the example of value of leaf node  $n$ ,  $attribute(n)$  returns  $A$ .

**Split.** For a query tree with more than one  $\tau_N$  node, this process is incorrect. As an example, consider the query tree of Figure 5 which has two  $\tau_N$  nodes (*book* and *dvd*). If we follow the mapping process described above, the tables DVD and Book will be joined, resulting in a cartesian product. In this expression, a book is repeated for

```

map(qt[])
Let sql[] be an array of strings, initially empty; Let numberqt be the number of split trees in qt[]
for k from 1 to numberqt do
  Let n be the node of type  $\tau_N$  in qt[k]
  sql[k] = "CREATE VIEW " + name(n) + " AS "
  sql[k] = sql[k] + "SELECT "
  Let N be the list of leaf nodes in qt[k]
  for i from 1 to size(N) do
    get next n in N
    if i > 1 then
      sql[k] = sql[k] + "," + variable(n) + "." + attribute(n) + " AS " + name(n)
    else
      sql[k] = sql[k] + variable(n) + "." + attribute(n) + " AS " + name(n)
    end if
    i = i + 1
  end for
  sql[k] = sql[k] + " FROM "; Let from = ""; Let N be the set of starred nodes in qt[k]
  for each n in N do
    Let join = ""; Let S be the list of source annotations in n; Let W be the list of where annotations in n
    for i = 1 to size(S) do
      get next s in S
      join = join + table(s) + " AS " + variable(s)
      if i < size(S) then
        join = join + " INNER JOIN "
      end if
      i = i + 1
    end for
    Let count = 0
    for i = 1 to size(W) do
      get next w in W
      if w is of the form  $\$x/A \text{ op } \$y/B$  AND  $\$x$  is bound to table X by a source annotation  $s \in S$  AND  $\$y$  is bound to table Y by a source annotation  $s' \in S$  then
        if count = 0 then
          join = join + " ON " + x.A op y.B
        else
          join = join + " AND " + x.A op y.B
        end if
        i = i + 1; count = count + 1
      end if
    end for
    if count = 0 then
      join = join + " ON (1=1) "
    end if
    if size(S) > 1 then
      join = "(" + join + ")"
    end if
    Let A be the set of starred ancestors of n; Let count = 0
    if n has a starred ancestor then
      join = " LEFT JOIN " + join
      for i = 1 to size(W) do
        get next w in W
        if w is of the form  $\$x/A \text{ op } \$y/B$  AND (( $\$x$  is bound to table X on node n AND  $\$y$  is bound to table Y on a node a in A) OR ( $\$x$  is bound to table X on a node a in A AND  $\$y$  is bound to table Y on node n)) then
          if count = 0 then
            join = join + " ON " + x.A op y.B
          else
            join = join + " AND " + x.A op y.B
          end if
          end if
          i = i + 1; count = count + 1
        end for
        if count = 0 then
          join = join + " ON (1=1) "
        end if
        from = "(" + from + join + ")"
      end if
    end for
    sql[k] = sql[k] + from; Let W' be the set of all where annotations on nodes of qt[k]. Let count = 0
    for each w' in W' do
      if w' is of the form  $\$x/A \text{ op } Z$  AND Z is an atomic value then
        if count = 0 then
          sql[k] = sql[k] + " WHERE " + x.A op Z
        else
          sql[k] = sql[k] + " AND " + x.A op Z
        end if
      end if
    end for
  end for
end for
return sql[]

```

**Algorithm 2:** The *map* algorithm

```

split( $qt$ )
Let  $qt[]$  be an array of query trees, initially empty
Let  $i = 0$ 
Let  $N$  be the set of nodes of type  $\tau_N$  in  $qt$ 
for each node  $n$  in  $N$  do
  inc  $i$ 
  {initialize  $t[i]$  with  $qt$ }
   $qt[i] = qt$ 
  repeat
    delete from  $qt[i]$  all subtrees rooted at a node  $z$  of type  $\tau_N$ , where  $z \neq n$ 
    retype the ancestors of the deleted nodes
  until  $n$  is the only node of type  $\tau_N$  in  $qt[i]$ 
end for
return  $qt[]$ 

```

**Algorithm 3:** The *split* algorithm

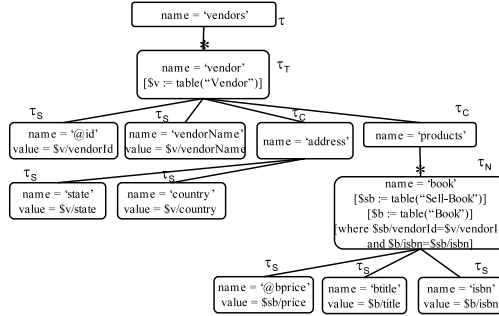


Figure 6: Partitioned query tree for  $\tau_N(\text{book})$

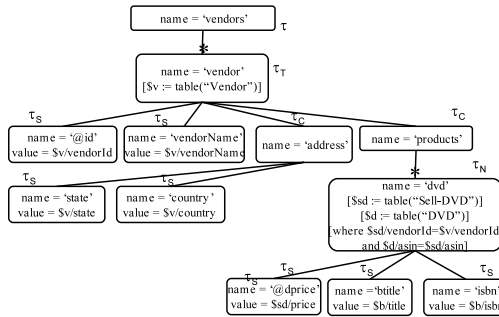


Figure 7: Partitioned query tree for  $\tau_N(\text{dvd})$

each DVD, violating the semantics of the query tree. We must therefore split a query tree into sub-query trees containing exactly one  $\tau_N$  node each before generating the corresponding relational views. After the splitting process, each sub-query tree produced is mapped to a relational view as explained above.

The splitting process consists in isolating a node  $n$  of type  $\tau_N$  in the query tree  $qt$ , and taking its subtree as well as its ancestors and their non-repeating descendants (types  $\tau_C$  and  $\tau_S$ ) to form a new tree  $qt_i$ . Recall that  $qt$  must have at least one  $\tau_N$  node by Proposition 2.1.

The first step to generate  $qt_i$  is to copy  $qt$  to  $qt_i$ . Then, delete from  $qt_i$  all subtrees rooted at nodes of type  $\tau_N$ , except for the subtree rooted at  $n$ . Observe that deleting a subtree  $r$  may change the abstract type of the ancestors of  $r$ . Specifically, if  $r$  has an ancestor  $a$  with type  $\tau_T$ , and  $r$  is  $a$ 's only starred descendant, the type of  $a$  becomes  $\tau_N$  after the deletion of  $r$ . Continue to delete subtrees rooted at nodes of type  $\tau_N$  in  $qt_i$  and retype ancestors until  $n$  is the only node of type  $\tau_N$  in  $qt_i$ . The process is repeated for every node of type  $\tau_N$  in  $qt$  and results in exactly one  $\tau_N$  node per split tree.

Formally, the *split* algorithm (algorithm 3) splits a query tree  $qt$ , producing one split tree  $qt_i$  for each node of type  $\tau_N$  in  $qt$ .

The result of this process for the query tree of Figure 5 is shown in Figures 6 and 7. Using these split trees, the corresponding relational views *ViewBook* and *ViewDVD* are (we name these views so we can refer to them in the examples of section 4):

	id	vendorName	state	country	bprice	btitle	isbn	dprice	dtitle	asin
$t_1$	1	Amazon	WA	US	38	Unix Network Programming	1111	NULL	NULL	NULL
$t_2$	1	Amazon	WA	US	29	Computer Networks	2222	NULL	NULL	NULL
$t_3$	1	Amazon	WA	US	NULL	NULL	NULL	29	Friends	D1111
$t_4$	2	Barnes and Noble	NY	US	38	Unix Network Programming	1111	NULL	NULL	NULL
$t_5$	2	Barnes and Noble	NY	US	38	Computer Networks	2222	NULL	NULL	NULL

Table 1: Tuples resulting from  $evalRel(eval(qt, d))$  for the query tree of Figure 5

	id	vendorName	state	country	bprice	btitle	isbn	dprice	dtitle	asin
$t_1$	1	Amazon	WA	US	38	Unix Network Programming	1111	NULL	NULL	NULL
$t_2$	1	Amazon	WA	US	29	Computer Networks	2222	NULL	NULL	NULL
$t_3$	2	Barnes and Noble	NY	US	38	Unix Network Programming	1111	NULL	NULL	NULL
$t_4$	2	Barnes and Noble	NY	US	38	Computer Networks	2222	NULL	NULL	NULL
$t_5$	1	Amazon	WA	US	NULL	NULL	NULL	29	Friends	D1111
$t_6$	2	Barnes and Noble	NY	US	NULL	NULL	NULL	NULL	NULL	NULL

Table 2: Tuples resulting from  $relEval(\{ViewBook, ViewDVD\}, d)$

```
CREATE VIEW VIEWBOOK AS
SELECT v.vendorId AS id, v.vendorName AS vendorName,
       v.state AS state, v.country AS country,
       sb.price AS bprice, b.isbn AS isbn, b.title AS btitle
FROM (Vendor AS v LEFT JOIN (Sell-Book AS sb INNER JOIN
Book AS B ON b.isbn=sb.isbn) ON v.vendorId=sb.vendorId);
```

```
CREATE VIEW VIEWDVD AS
SELECT v.vendorId AS id, v.vendorName AS vendorName,
       v.state AS state, v.country AS country,
       sd.price AS dprice, d.asin AS asin, d.title AS dtitle
FROM (Vendor AS v LEFT JOIN (Sell-DVD AS sd INNER JOIN
DVD AS d ON d.asin=sd.asin) ON v.vendorId=sd.vendorId)
```

As described above, *split* takes as input the original query tree  $qt$  and produces as output a set of query trees  $\{qt_1, \dots, qt_n\}$ , each of which has one  $\tau_N$  node; *map* takes  $\{qt_1, \dots, qt_n\}$  as input and produces a set of relational view expressions  $\{V_1, \dots, V_n\}$ , where each  $V_i$  is produced from  $qt_i$  as described above. It follows directly from these algorithms that:

**Proposition 3.1** *The number of relational view expressions in  $map(split(qt))$  is the number of  $\tau_N$  nodes in  $qt$ .*

The correctness of the set of relational view expressions resulting from *map* and *split* can be understood in the following sense: Each tuple in the bindings relations for the XML view is in one or more instances of the corresponding relational views. To be more precise, we define the following:

**Definition 3.1** *The evaluation schema  $S$  of a query tree  $qt$  is the set of all names of leaf nodes in  $qt$ .*

**Definition 3.2** *Let  $x$  be an XML instance of a query tree  $qt$  with evaluation schema  $S$ , in which the instance nodes are annotated by the query tree type from which they were generated. Let  $n$  be the deepest  $\tau_N$  or  $\tau_T$  instance nodes for some root to leaf path in  $x$ . Let  $p$  be the set of nodes in the path from  $n$  to the root of  $x$ . An evaluation tuple of  $x$  is created from  $n$  by associating the value of each leaf node  $l$  that is a descendant of  $n$  or of some node in  $p$  with the attribute in  $S$  corresponding to the name of  $l$ , and leaving the value of all other attributes in  $S$  null.*

*The multi-set of all evaluation tuples of  $x$  is called its evaluation relation and is denoted  $evalRel(x)$ .*

For example, table 1 shows the result of  $evalRel(x)$  for the query tree of Figure 5.

**Definition 3.3** *Let  $\{V_1, \dots, V_n\}$  be defined over a relational schema  $\mathcal{D}$ , and  $d$  be an instance of  $\mathcal{D}$ . Then  $relEval(\{V_1, \dots, V_n\}, d)$  denotes the set of relational instances that result from taking the outer union of the evaluation of each  $V_i$  over  $d$ :  $relEval(\{V_1, \dots, V_n\}, d) = evalV(V_1, d) \cup \dots \cup evalV(V_n, d)$ , where  $\cup$  denotes outer union, and  $evalV(V, d)$  instantiates  $V$  over  $d$ .*

For example,  $relEval(\{ViewBook, ViewDVD\}, d)$  is the outer union of  $evalV(ViewBook, d)$  and  $evalV(ViewDVD, d)$ , whose result is shown on table 2.

The correctness of the set of relational views resulting from *map* and *split* can now be understood in the following sense:

**Theorem 3.1** *Given a query tree  $qt$  defined over a database  $\mathcal{D}$  and an instance  $d$  of  $\mathcal{D}$ , then  $evalRel(eval(qt, d)) \subseteq relEval(map(split(qt)), d)$ .*

*Proof:* The  $\subseteq$  operation needs the two multi-sets being compared to be union compatible. By definition, the schema of  $evalRel$  is the *evaluation schema*  $S$ , which is composed of all leaf node names in  $qt$ . The execution of  $map(split(qt), d)$  results in a set of relational views  $\{V_1, \dots, V_n\}$ . Each view  $V_i$  is a schema composed of names of leaf nodes in  $qt_i$  (which is produced by  $split(qt)$ ). By definition of  $split$ , each split tree  $qt_i$  contains a single  $\tau_N$  node  $n_i$ : the subtrees rooted at  $\tau_N$  nodes different from  $n_i$  are deleted from  $qt_i$ . However, nodes deleted in  $qt_i$  are preserved in  $qt_j$ , so that each node  $n$  in  $qt$  is in at least one of the  $qt_1, \dots, qt_n$ . Consequently, the schema of  $V_1 \cup \dots \cup V_n$  equals  $S$ .

Assume  $t$  is in  $evalRel(eval(qt, d))$ , but not in  $relEval(map(split(qt), d))$ . Let  $x$  be the XML view resulting from  $eval(qt, d)$ . Since  $t$  is in  $evalRel(eval(qt, d))$ , it was constructed by taking values from the leaf nodes in a given path  $p$ . The path  $p$  starts in a node  $n$  which is the deepest node of type  $\tau_N$  or  $\tau_T$  in a given subtree and goes up to the root of  $x$ . If  $n$  is of type  $\tau_N$ , and  $V_i$  is the view corresponding to  $n$ , then  $t$  is in  $evalV(V_i, d)$ , and consequently,  $t$  is in  $relEval(map(split(qt), d))$ , a contradiction. If  $n$  is of type  $\tau_T$ , then the node that originated  $n$  in the query tree has at least one node of type  $\tau_N$  in its subtree. Assume  $V_j, \dots, V_k$  are the relational views corresponding to those  $\tau_N$  nodes. Consequently,  $t$  is in  $V_j \cup \dots \cup V_k$ , and thus in  $relEval(map(split(qt), d))$ , a contradiction. ■

Furthermore, the tuples in  $relEval(map(split(qt), d)) - evalRel(eval(qt, d))$  represent starred nodes with an empty evaluation (which we call “stubbied” nodes). More precisely:

**Definition 3.4** Let  $x$  be an XML instance of a query tree  $qt$  with evaluation schema  $S$ , and  $n$  be a  $\tau_N$  or  $\tau_T$  instance node in  $x$ . A stubbed tuple of  $x$  is created from  $n$  by associating the value of each leaf node  $l$  that is an ancestor of  $n$  with the attribute in  $S$  corresponding to the name of  $l$ , and leaving the value of all other attributes in  $S$  null.

The set of all stubbed tuples of  $x$  is denoted  $stubs(x)$ .

As an illustration of a stubbed tuple, consider tuple  $t_6$  in table 2. Since the XML instance of Figure 2 does not have any book sold by vendor *Barnes and Noble*, there is a tuple  $[2, \text{Barnes and Noble}, \text{NY}, \text{US}, \text{null}, \text{null}]$  in *ViewBook* which was added by the LEFT join. This is correct, since *vendor* is in a common part of the view, so its information appears both in *ViewBook* and *ViewDVD*. However,  $t_6$  is not in table 1, since when the entire view is evaluated, this vendor joins with a DVD.

**Theorem 3.2** Given a query tree  $qt$  defined over a database  $\mathcal{D}$  and an instance  $d$  of  $\mathcal{D}$ , then every tuple  $t$  in  $relEval(map(split(qt), d)) - evalRel(eval(qt, d)) \subseteq stubs(x)$ .

*Proof:* Tuples in  $relEval(map(split(qt), d))$  that are not in  $evalRel(eval(qt, d))$  are those resulting from left outer joins with no match in a given relational view  $V_i \in map(split(qt), d)$ . Since  $stubs(x)$  contains tuples that has nulls in attributes related to descendant nodes, and a LEFT JOIN always keeps information the ancestor, then  $relEval(map(split(qt), d)) - evalRel(eval(qt, d)) \subseteq stubs(x)$ . ■

Note that the statement of correctness is *not* that the XML view can be constructed from instances of the underlying relational views. The reason is that we do not know whether or not keys of relations along the path from  $\tau_N$  nodes to the root are preserved, and therefore do not have enough information to group tuples from different relational view instances together to reconstruct the XML view. When keys at all levels *are* preserved, then the query tree can be modified to a form in which the variables iterate over the underlying relational views instead of base tables, and used to reconstruct the XML view. We call this algorithm *replace*.

**Replace.** For query trees that preserve the keys of each source table in the resulting XML view, we can use the corresponding relational views to reconstruct the view. Assuming that  $map(qt) = \{V_1, \dots, V_n\}$ , the algorithm *replace* replaces references to relational tables in the source and where annotations of  $qt$  by references to the set of relational views  $V_1, \dots, V_n$ . This must be done in a way such that for any instance  $d$  of the underlying relational database  $D$ , evaluating  $qt$  over  $d$  is the same as applying  $eval$  over  $replace(qt)$  using the evaluation of the relational views  $V_1, \dots, V_n$  produced by  $map(split(qt))$ .

Before presenting the algorithm, we present some definitions that will be used within this section. We say  $view(n)$  is the relational view corresponding to node  $n$ , if  $n$  is of type  $\tau_N$ .  $Atts(n)$  is the set of leaf node names whose values are specified using the variables declared on node  $n$ . For example, in Figure 5,  $Atts(\text{book}) = \{bprice, btitle, isbn\}$ . Notice that we exclude the “@” from attribute names. In the same way, the function  $AttsAncestrals(n)$  returns a set of leaf nodes whose values are specified using variables declared in the ancestors of node  $n$ . As an example,  $AttsAncestrals(\text{book}) = \{id, vendorName, state, country, url\}$ . The function  $var(n)$

```

replace(qt)
Let  $qt_r = qt$ 
for each node  $n$  of type  $\tau_N$  or  $\tau_T$  in  $qt_r$  do
  remove all source and where annotations from  $n$ 
  if abstract_type( $n$ ) =  $\tau_N$  then
     $def = createSourceDef(n, view(n))$ 
  else
    Let  $n_1, \dots, n_m$  be the set of nodes with abstract type  $\tau_N$  in the subtree rooted at  $n$ 
     $def = createSourceDef(n, (view(n_1) UNION \dots UNION view(n_m)))$ 
  end if
  Let  $s$  be a source annotation
   $s = [variable(n) := Table(X)]$ , where  $X$  is defined as  $def$ 
  Annotate  $n$  with  $s$ 
  if  $n$  has a starred ancestor  $a$  then
    Let  $w$  be a where annotation
    Let  $W$  be a string
    for each  $el$  in  $AttsAncestors(n)$  do
       $W = W + var(a)/el = var(n)/el$ 
      if  $el$  is not the last element in  $AttsAncestors(n)$  then
         $W = W + " AND "$ 
      end if
    end for
    Annotate  $n$  with  $w = [where W]$ 
  end if
end for
for each node  $n$  of type  $\tau_S$  in  $qt_r$  do
  Let the value of  $n$  be of the form  $\$x/A$ 
  Let  $a$  be the starred ancestor of  $n$  (if  $n$  is starred, then  $a = n$ )
  Replace the value of  $n$  by  $var(a)/A$ 
end for
return  $qt_r$ 

createSourceDef( $n, viewName$ )
 $def = "SELECT DISTINCT" + AttsAncestors(n) + ", " + Atts(n) + "FROM " + viewName$ 
return  $def$ 

```

**Algorithm 4:** The *replace* algorithm

returns a unique variable name to be used in node  $n$ . Note that every call of  $var(n)$  for the same node  $n$  returns the same variable name. The function  $var(\$x)$  finds the node  $n$  in which variable  $\$x$  was defined, and return the result of  $var(n)$ .

The *replace*( $qt$ ) algorithm (algorithm 4) takes each starred node  $n$  in  $qt$  and analyze it. It first removes all source and where annotations from  $n$ . Then, it creates a single source annotation that bounds a new unused variable  $\$x$  to  $X$ , where  $X$  is defined over the relational views that carries values for node  $n$ . The complete algorithm is shown on algorithm 4. It returns a modified query tree  $qt_r$ , which references  $V_1, \dots, V_n$  instead of base tables.

## 4 Updates

Given an update against a well-behaved view, we translate it to a set of SQL update statements against the corresponding relational view expressions, so existing work on updates through relational views can be used to translate the updates to the underlying relational database. In this section, we start by defining XML updates and then describe the translation. We also summarize how to determine whether or not an update is side-effect free.

Although no standard has been established for an XML update language, several proposals have appeared [7, 36, 12, 28]. The language described below is much simpler than any of these proposals, and in some sense can be thought of as an internal form for one of these richer languages (assuming a static translation of updates [12]). The simplicity of the language allows us to focus on the key problem we are addressing.

### 4.1 Update language

Updates are specified using path expressions to point to a set of target nodes in the XML tree at which the update is to be performed. For insertions and modifications, the update must also specify a  $\Delta$  containing the new values.

**Definition 4.1** An update operation  $u$  is a triple  $\langle t, \Delta, \text{ref} \rangle$ , where  $t$  is the type of operation (insert, delete, modify);  $\Delta$  is the XML tree to be inserted, or (in case of a modification) an atomic value; and  $\text{ref}$  is a simple

path expression in XPath [18] which indicates where the update is to occur.

The path expression  $ref$  is evaluated from the root of the tree and may yield a set of nodes which we call *update points*. In the case of *modify*, it must evaluate to a set of leaf nodes. We restrict the filters used in  $ref$  to conjunctions of comparisons of attributes or child elements with atomic values, and call the expression resulting from removing filters in  $ref$  the *unqualified portion* of  $ref$ . For example, the unqualified portion of  $/vendors/vendor[@id="01"]$  is  $/vendors/vendor$ .

**Definition 4.2** An update path  $ref$  is valid with respect to a query tree  $qt$  iff the unqualified portion of  $ref$  is non-empty when evaluated on  $qt$ .

For example,  $/vendors/vendor[@id="01"]/vendorName$  is a valid path expression with respect to the query tree of Figure 5, since the path  $/vendors/vendor/vendorName$  is non-empty when evaluated on that query tree.

The semantics of *insert* is that  $\Delta$  is inserted as a child of the nodes indicated by  $ref$ ; the semantics of *modify* is that the atomic value  $\Delta$  overwrites the values of the leaf nodes indicated by  $ref$ ; and the semantics of a *delete* is that the subtrees rooted at nodes indicated by  $ref$  are deleted.

The following examples refer to figure 2:

**Example 4.1** To insert a new book selling for \$38 under the vendor with  $id="01"$  we specify:  $t = \text{insert}$ ,  $ref = /vendors/vendor[@id="01"]/products$ ,

```
 $\Delta = \{ \langle \text{book } bprice = "38" \rangle$ 
   $\langle \text{btitle} \rangle \text{New Book} \langle \text{btitle} \rangle \langle \text{isbn} \rangle 9999 \langle \text{isbn} \rangle$ 
   $\langle \text{book} \rangle \}$ .
```

**Example 4.2** To change the  $vendorName$  of the vendor with  $id = "01"$  to  $Amazon.com$  we specify:  $t = \text{modify}$ ,  $ref = /vendors/vendor[@id="01"]/vendorName$ ,  $\Delta = \{Amazon.com\}$ .

**Example 4.3** To delete all books with title "Computer Networks" we specify:  $t = \text{delete}$ ,  $ref = /vendors/vendor/products/book[btitle="Computer Networks"]$ .

Note that not all insertions and deletions make sense since the resulting XML view may not conform to the DTD of the query tree (for details, see section 2.4). For example, the deletion specified by the path  $/vendors/vendor/vendorName$  would not conform to the DTD of figure 5 since  $vendorName$  is a required subelement of  $vendor$ . We must also check that  $\Delta$ 's inserted and subtrees deleted are correct.

**Definition 4.3** An update  $\langle t, \Delta, ref \rangle$  against an XML view specified by a query tree  $qt$  is correct iff

- $ref$  is valid with respect to  $qt$ ;
- if  $t$  is a modification, then the unqualified portion of  $ref$  evaluated on  $qt$  arrives at a node whose abstract type is  $\tau_S$ ;
- if  $t$  is an insertion, then the unqualified portion of  $ref +$  the root of  $\Delta$  evaluated on  $qt$  arrives at a node whose incoming edge is starred (equivalently, its abstract type is  $\tau_T$  or  $\tau_N$ );
- if  $t$  is a deletion, then the unqualified portion of  $ref$  evaluated on  $qt$  arrives at a node whose incoming edge is starred;
- if nonempty, then  $\Delta$  conforms to the DTD of the element arrived at by  $ref$ .

For example, the deletion of example 4.3 is correct since  $book$  is a starred subelement of  $products$ . However, the deletion specified by the update path  $/vendors/vendor/vendorName$  is not correct since  $vendorName$  is of abstract type  $\tau_S$ , as would the deletion specified by the invalid update path  $/vendors/vendor/dvd$ .

As another example, the insertion of example 4.1 is correct since  $book$  (arrived at by  $/vendors/vendor/products$ ) is a starred subelement of  $products$ , the DTD for  $book$  is  $\langle !ELEMENT book (btitle, isbn) \rangle$ , and  $\Delta$  conforms to this DTD. However, the following insertion would not be correct for the update path  $/vendors/vendor[@id="01"]/products$  and

```
 $\Delta = \{ \langle \text{book } bprice="38" \rangle \langle \text{rating} \rangle \text{Children} \langle \text{rating} \rangle \langle \text{book} \rangle \}$ 
```

since the  $isbn$  and  $btitle$  subelements are missing, and  $book$  does not have a  $rating$  subelement.

```

translateUpdate( $x, qt, u$ )
case  $u.t$ 
  insert: translateInsert( $x, qt, u.ref, u.\Delta$ )
  delete: translateDelete( $x, qt, u.ref$ )
  modify: translateModify( $x, qt, u.ref, u.\Delta$ )
end case

```

**Algorithm 5:** The *translateUpdate* algorithm

```

translateInsert( $V, qt, ref, \Delta$ )
{Insert  $\Delta$  in the XML view  $V$  using  $ref$  as insertion point.  $\Delta$  must be inserted under every node resulting from the evaluation of  $ref$  in  $V$ .  $qt$  is the query tree.} {Assume that  $view(n)$  returns the name of the rel. view associated with node  $n$ }
Let  $p$  be the unqualified portion of  $ref$  concatenated with the root of  $\Delta$ 
Let  $m$  be the node resulting from the evaluation of  $p$  against  $qt$ 
Let  $N$  be the set of nodes resulting from the evaluation of  $ref$  in  $V$ 
for each  $n$  in  $N$  do
  if  $abstract\_type(m) = \tau_N$  then
    generateInsertSQL( $view(m), root(\Delta), n, V$ )
  else
    Let  $X$  be the set of nodes of abstract type  $\tau_N$  in  $\Delta$ 
    for each  $x$  in  $X$  do
      generateInsertSQL( $view(x), x, n, V$ )
    end for
  end if
end for

generateInsertSQL( $RelView, r, InsertionPoint, V$ )
{Inserts the subtree rooted at  $r$  into  $RelView$ }
 $sql = \text{"INSERT INTO"} + RelView + getAttributes(RelView)$ 
 $sql = sql + \text{" VALUES ("}$ 
for  $i = 0$  to  $getTotalNumberAttributes(RelView) - 1$  do
   $att = getAttribute(RelView, i)$ 
  if  $att$  is a child  $n$  of  $r$  then
     $sql = sql + getValue(n)$ 
  else
    Find  $att$  in  $V$ , starting from  $InsertionPoint$  examining the leaf nodes until  $V$ 's root is found
    Let the node found be  $m$ 
     $sql = sql + getValue(m)$ 
  end if
if  $i < getTotalNumberAttributes(RelView) - 1$  then
   $sql = sql + ", "$ 
else
   $sql = sql + ")"$ 
end if
end for

```

**Algorithm 6:** The *translateInsert* algorithm

## 4.2 Mapping XML updates to relational views

We now discuss how correct updates to an XML view are translated to SQL updates on the corresponding relational views produced in the previous section.

Throughout this section, we will use the XML view 2 of figure 5 as an example. The relational views *ViewBook* and *ViewDVD* corresponding to this XML view were presented in section 3.

The translation algorithm for insertions, deletions and modifications, *translateUpdate*, is given in algorithm 5. What it does is to check the type of update operation and call the corresponding algorithm to translate the update. All the three algorithms (*translateInsert*, *translateDelete* and *translateModify*) assume that the update specification  $u$  was already checked for schema conformance.

### 4.2.1 Insertions

To translate an insert operation on the XML view to the underlying relational views we do the following: First, the unqualified portion of the update path  $ref$  is used to locate the node in the query tree under which the insertion is to take place. Together with  $\Delta$ , this will be used to determine which underlying relational views are affected. Second,  $ref$  is used to query the XML instance and identify the update points. Third, SQL insert statements are generated for each underlying relational view affected using information in  $\Delta$  as well as information about the labels and values in subtrees rooted along the path from each update point to the root of the XML instance.

Observe that by proposition 2.2, there is at most one node of type  $\tau_N$  along the path from any node to the root of the query tree and that insertions can never occur below a  $\tau_N$  node, since all nodes below a  $\tau_N$  node are of type  $\tau_S$  or  $\tau_C$  by definition.

The algorithm *translateInsert* is presented in algorithm 6.

For example, to translate the insertion of example 4.1, we use the unqualified update path  $/vendors/vendor/products$  on the query tree of figure 5, and find that the type of the update point is  $\tau_C(products)$ . Continuing from  $\tau_C(products)$  using the structure of  $\Delta$ , we discover that the only  $\tau_N$  node in  $\Delta$  is its root, which is of type  $\tau_N(book)$ . The underlying view affected will therefore be *ViewBook*. We then use the update path  $ref=/vendors/vendor[id="01"]/products$  to identify update points in the XML document. In this case, there is one node (8). Therefore, a single SQL insert statement against view *ViewBook* will be generated.

To generate the SQL insert statement, we must find values for all attributes in the view. Some of these attribute-value pairs are found in  $\Delta$ , and others must be taken from the XML instance by traversing the path from each update point to the root and collecting attribute-value pairs from the leaves of trees rooted along this path. In example 4.1,  $\Delta$  specifies  $bprice="38"$ ,  $btitle="New Book"$  and  $isbn="9999"$ . Along the path from the node 8 to the root in the XML instance of figure 2, we find  $id="01"$ ,  $vendorName="Amazon"$ ,  $state="WA"$  and  $country="US"$ . Combining this information, we generate the following SQL insert statement:

```
INSERT INTO VIEWBOOK (id, vendorName, state, country,
  bprice, isbn, btitle)
VALUES ("01", "Amazon", "WA", "US", 38, "9999", "New Book")
```

As another example, consider the following insertion against the view 2:  $t = insert, ref = /vendors,$

```
 $\Delta = \{ \langle vendor \ id="03" \rangle$ 
   $\langle vendorName \rangle$ New Vendor $\langle /vendorName \rangle$ 
   $\langle address \rangle$ 
     $\langle state \rangle$ PA $\langle /state \rangle$ 
     $\langle country \rangle$ US $\langle /country \rangle$ 
   $\langle /address \rangle$ 
   $\langle products \rangle$ 
     $\langle book \ bprice="30" \rangle$ 
       $\langle btitle \rangle$ Book 1 $\langle /btitle \rangle \langle isbn \rangle$ 9111 $\langle /isbn \rangle \langle /book \rangle$ 
     $\langle book \ bprice="30" \rangle$ 
       $\langle btitle \rangle$ Book 2 $\langle /btitle \rangle \langle isbn \rangle$ 9222 $\langle /isbn \rangle \langle /book \rangle$ 
     $\langle dvd \ dprice="30" \rangle$ 
       $\langle dtitle \rangle$ DVD 1 $\langle /dtitle \rangle \langle asin \rangle$ D9333 $\langle /asin \rangle \langle /dvd \rangle$ 
   $\langle /products \rangle$ 
 $\langle /vendor \rangle \}$ .
```

The unqualified update path  $ref$  evaluated against the query tree of figure 5 yields a node  $\tau(vendors)$ , which is the root. Continuing from here using labels in  $\Delta$ , we discover two nodes of type  $\tau_N$ :  $\tau_N(book)$  and  $\tau_N(dvd)$ . We will therefore generate SQL insert statements to *ViewBook* and as well as *ViewDVD*.

Evaluating  $ref$  against the XML instance of figure 2 yields one update point, node 1. Traversing the path from this update point to the root yields no label-value pairs (since the update point is the root itself). We then identify each node of type  $\tau_N$  in  $\Delta$ , and generate one insertion for each of them. As an example, traversing the path from the first  $\tau_N(book)$  node in  $\Delta$  yields label-value pairs  $bprice = "30"$ ,  $btitle = "Book 1"$ , and  $isbn = "9111"$ . Going up to the root of, we have  $id = "03"$ ,  $vendorName = "New Vendor"$ ,  $state = "PA"$  and  $country = "US"$ . This information is therefore combined to generate the following SQL insert statement:

```
INSERT INTO VIEWBOOK (id, vendorName, state, country,
  bprice, isbn, btitle)
VALUES ("03", "New Vendor", "PA", "US", 30, "9111", "Book 1");
```

In a similar way, information is collected from the remaining two  $\tau_N$  nodes in  $\Delta$  to generate:

```
INSERT INTO VIEWBOOK (id, vendorName, state, country,
  bprice, isbn, btitle)
VALUES ("03", "New Vendor", "PA", "US", 30, "9222", "Book 2");
INSERT INTO VIEWDVD (id, vendorName, state, country,
  dprice, asin, dtitle)
VALUES ("03", "New Vendor", "PA", "US", 30, "D9333", "DVD 1");
```

## 4.2.2 Modifications

By definition, modifications can only occur at leaf nodes. To process a modification, we do the following: First, we use the unqualified  $ref$  against the query tree to determine which relational views are to be updated. This is done by looking at the first ancestor of the node specified by  $ref$  which has type  $\tau_T$  or  $\tau_N$ , and finding all nodes of type  $\tau_N$  in its subtree. (At least one  $\tau_N$  node must exist, by definition.) If the leaf node that is being modified is of type  $\tau_N$  itself, then it is guaranteed that the update will be mapped only to the relational view corresponding to this node.

Second, we generate the SQL modify statements. The qualifications in  $ref$  are combined with the terminal label of  $ref$  and value specified by  $\Delta$  to generate an SQL update statement against the view. The corresponding algorithm is presented in algorithm 7.

```

translateModify( $V, qt, ref, \Delta$ )
Let  $p$  be the unqualified portion of  $ref$ 
Let  $m$  be the node resulting from the evaluation of  $p$  against  $qt$ 
if  $abstract\_type(m) = \tau_N$  then
     $r = m$ 
else
    Let  $r$  be the ancestor of  $m$  whose abstract type is  $\tau_T, \tau_G$  or  $\tau_N$ 
end if
if  $abstract\_type(r) = \tau_N$  then
    generateModifySQL(view( $r$ ),  $\Delta, ref$ )
else
    Let  $X$  be the set of nodes with abstract type  $\tau_N$  under  $r$ 
    for each  $x$  in  $X$  do
        generateModifySQL(view( $x$ ),  $\Delta, ref$ )
    end for
end if

generateModifySQL(RelView,  $\Delta, ref$ )
 $sql = "UPDATE " + RelView + " SET "$ 
Let  $t$  be the terminal node in  $ref$ 
 $sql = sql + t + "=" + \Delta$ 
for each filter  $f$  in  $ref$  do
    if  $f$  is the first filter in  $ref$  then
         $sql = sql + " WHERE " + f$ 
    else
         $sql = sql + " AND " + f$ 
    end if
end for

```

**Algorithm 7:** The *translateModify* algorithm

```

translateDelete( $V, qt, ref$ )
{Deletes the subtree rooted at  $ref$  from  $V$ }
Let  $p$  be the unqualified portion of  $ref$  concatenated with the root of  $\Delta$ 
Let  $m$  be the node resulting from the evaluation of  $p$  against  $qt$ 
if  $abstract\_type(m) = \tau_N$  then
    generateDeleteSQL(view( $m$ ),  $ref$ )
else
    Let  $X$  be the set of nodes of abstract type  $\tau_N$  under  $m$ 
    for each  $x$  in  $X$  do
        generateDeleteSQL(view( $x$ ),  $ref$ )
    end for
end if

generateDeleteSQL(RelView,  $ref$ )
 $sql = "DELETE FROM " + RelView$ 
for each filter  $f$  in  $ref$  do
    if  $f$  is the first filter in  $ref$  then
         $sql = sql + " WHERE " + f$ 
    else
         $sql = sql + " AND " + f$ 
    end if
end for

```

**Algorithm 8:** The *translateDelete* algorithm

For example, consider the update in example 4.2. The unqualified *ref* is */vendors/vendor/vendorName*. The  $\tau_N$  nodes in the subtree rooted at *vendor* (the first  $\tau_T$  or  $\tau_N$  ancestor of *vendorName*) are  $\tau_N(book)$  and  $\tau_N(dvd)$ , and we will therefore generate SQL update statements for both *ViewBook* and *ViewDVD*. We then use the qualification *id = "01"* from *ref = /vendors/vendor[@id = "01"]/ vendorName* together with the new value in  $\Delta$ , to yield the following SQL modify statements:

```

UPDATE VIEWBOOK SET vendorName="Amazon.com" WHERE id="01";
UPDATE VIEWDVD SET vendorName="Amazon.com" WHERE id="01"

```

### 4.2.3 Deletions

Deletions are very simple to process. First, the unqualified portion of the update path *ref* is used to locate the node in the query tree at which the deletion is to be performed. This is then used to determine which underlying relational views are affected by finding all  $\tau_N$  nodes in its subtree. Second, SQL delete statements are generated for each underlying relational view affected using the qualifications in *ref*. The corresponding algorithm is presented in algorithm 8.

As an example, consider the deletion in example 4.3. The unqualified update path expression is */vendors/vendor/products/book*. The only  $\tau_N$  node in the subtree indicated by this path in the query tree is  $\tau_N(book)$ .

This means that the deletion will be performed in *ViewBook*. Examining the update path `/vendors/vendor/products/book[btitle="Computer Networks"]` yields the label-value pair `btitle="Computer Networks"`. Thus the deletion on the XML view is translated to an SQL delete statement as:

```
DELETE FROM VIEWBOOK WHERE btitle="Computer Networks"
```

It is important to notice that if a tuple  $t$  in one relation “owns” a set of tuples in another relation via a foreign key constraint (e.g. a vendor “owns” a set of books), then deletions must cascade in the underlying relational schema in order for the deletion of  $t$  specified through the XML view to be allowed by the underlying relational system.

### 4.3 Correctness

Since we are not focusing on how updates over relational views are mapped to the underlying relational database, our notion of correctness of the update mappings is their effect on each relational view *treated as a base table*.

Let  $x = eval(qt, d)$  be the initial XML instance,  $u$  be the update as specified in Definition 4.1, and  $apply(x, u)$  be the updated XML instance resulting from applying  $u$  to  $x$ . The function  $translateUpdate(x, qt, u)$  (shown in section 4.2) translates  $u$  to a set of SQL update statements  $\{U_{11}, \dots, U_{1m_1}, \dots, U_{n1}, \dots, U_{nm_n}\}$ , where each  $U_{ij}$  is an update on the underlying view instance  $v_i = evalV(V_i, d)$  generated by  $map(split(qt))$ .

We use the notation  $v'_i = applyR(v_i, \{U_{i1}, \dots, U_{im_i}\})$  to denote the application of  $\{U_{i1}, \dots, U_{im_i}\}$  to  $v_i$ , resulting in the updated view  $v'_i$ . If the set of updates for a given  $v_i$  is empty, then  $v'_i = v_i$ .

**Theorem 4.1** *Given a query tree  $qt$  defined over database  $\mathcal{D}$ , then for any instance  $d$  of  $\mathcal{D}$  and correct update  $u$  over  $qt$ ,  $evalRel(apply(x, u)) \subseteq v'_1 \cup \dots \cup v'_n$ , where  $\cup$  denotes outer union.*

*Proof:* Since the update  $u$  does not change the view schema, and the application of an update  $U_{ij}$  over view  $v_i$  also does not change  $v_i$ 's schema, by theorem 3.1 we have that  $evalRel(apply(x, u))$  and  $v'_1 \cup \dots \cup v'_n$  have the same schema (are union compatible).

**Insertions** Suppose  $t$  is a tuple in  $evalRel(apply(x, u))$ , resulting from a insertion of a subtree in  $x$ . Assume  $t$  is not in  $v'_1 \cup \dots \cup v'_n$ . Assume update  $U_{ij}$  is the translation of  $u$ .

Consider a tuple  $t'$  which was inserted by update  $U_{ij}$  in  $v_i$ . Since  $U_{ij}$  is the translation of  $u$ ,  $t'$  has the values of one of the subtrees that were inserted in  $x$  by  $u$ , and also the values of  $x$  that were above the update point  $ref$  of  $u$ . As a consequence,  $t = t'$ , and  $t$  is in  $v'_1 \cup \dots \cup v'_n$ , a contradiction.

The same applies for the insertion of a more complex subtree. It will generate several tuples  $t_1, \dots, t_n$  to appear in  $evalRel(apply(x, u))$ . Each of these tuples will be inserted in the relational views by a set of updates  $U_{ij}, \dots, U_{kl}$ . So we have that  $evalRel(apply(x, u)) \subseteq v'_1 \cup \dots \cup v'_n$  holds for insertions.

**Modifications** Suppose  $t$  is a tuple in  $evalRel(apply(x, u))$ , resulting from a modification of a leaf value in  $x$ . Assume  $t$  is not in  $v'_1 \cup \dots \cup v'_n$ . Assume update  $U_{ij}$  is the translation of  $u$ .

Consider a tuple  $t'$  which was modified by update  $U_{ij}$  in  $v_i$ . Since  $U_{ij}$  is the translation of  $u$ ,  $t'$  had a single attribute modified - the one that was updated in  $x$ . As a consequence,  $t = t'$ , and  $t$  is in  $v'_1 \cup \dots \cup v'_n$ , a contradiction.

The same applies for modifications that affect more than one leaf in  $x$ , that is, when  $ref$  in  $u$  evaluates to more than one update point. For every node affected by the modification, will be generated one modification  $U_{ij}$ . Since by theorem 3.1 all tuples in  $evalRel(x)$  are in  $v_1 \cup \dots \cup v_n$ , then  $evalRel(apply(x, u)) \subseteq v'_1 \cup \dots \cup v'_n$ .

**Deletions** Following the inverse reasoning of insertions, every subtree deleted from  $x$  makes a tuple to disappear from  $evalRel(apply(x, u))$ s. Analogously, the translation  $U_{ij}$  of  $u$  will make that tuple to disappear from  $v'_1 \cup \dots \cup v'_n$ , so  $evalRel(apply(x, u)) \subseteq v'_1 \cup \dots \cup v'_n$  holds. ■

**Theorem 4.2** *Given a query tree  $qt$  defined over a database  $\mathcal{D}$  and an instance  $d$  of  $\mathcal{D}$ , then  $v'_1 \cup \dots \cup v'_n - evalRel(apply(x, u)) \subseteq stubs(apply(x, u))$ .*

*Proof:* Insertions of uncomplete subtrees or deletions of uncomplete subtrees may cause tuples to be filled in with nulls because of the LEFT JOINS in some  $v'_i$ . These tuples, however, will be in  $stubs$ . The reasoning is the same as in proof of theorem 3.2. ■

Note that a correctness definition like  $apply(eval(qt, d), u) \equiv eval(qt, d')$ , where  $d'$  is the updated relational database state resulting from the application of the translated view updates  $\{U_{11}, \dots, U_{1m_1}, \dots, U_{n1}, \dots, U_{nm_n}\}$

to updates on  $d$ , does not make sense due to the fact that we do not control the translation of view updates. Therefore we cannot claim that they are side-effect free.

In the next subsection, we discuss a scenario in which this claim can be made.

#### 4.4 Updatability

There are several choices of techniques that could be used to translate from updates on relational views to updates on the underlying relational database. Some consider a translation to be correct if it does not affect any part of the database that is outside the view [9, 29]. Others consider a translation to be correct as long as it corresponds exactly to the specified update, and does not affect anything else in the view [21]. Still others use additional information to build specific translators for each view [26, 31, 38]. Here, we choose [21] to illustrate how reasoning about *side-effect free* relational view updates can be extended to XML views.

In [13], we define conditions under which XML views constructed by “nest-last” nested relational algebra (NRA) expressions are updatable. Since nest-last NRA expressions perform nests over a relational algebra expression, our results are based on the ability to unnest the NRA expression to obtain a (single) corresponding relational view, and then build on the results of [21] to detect updatability. Since query trees also express nesting and are mapped to a *set* of corresponding relational views, we can use these results to reason about the updatability of XML views constructed by query trees. We assume the underlying relational database is in BCNF (as required by [21]), and impose three restrictions on query trees: (1) each table must be bound to at most one variable; (2) each value in a leaf node must be unique, that is, if the value of  $n$  is specified as  $\$x/A$ , then this value specification does not appear on any other node in the query tree; (3) comparisons on  $u.ref$  must be equalities. These restrictions are imposed so that the resulting relational views do not include joins of the same tables, and projections of the same attribute (as required by [21]). The restrictions on equalities are also required by [21].

**Theorem 4.3** *A correct update  $u$  to an XML view defined by a query tree  $qt$  is side-effect free if for all  $(U_i, V_i)$ , where  $V_i$  is the corresponding relational view of  $qt_i$  and  $U_i$  is the translation of  $u$  over  $V_i$ ,  $U_i$  is side-effect free in  $V_i$ .*

*Proof:* In our approach, a given XML update  $u$  can be mapped to a *set* of updates over the corresponding relational views. Formally,  $u = \bigcup_{i=n}^{i=1} (U_i, V_i)$ ,  $n \geq 1$ . The update  $u$  is an atomic operation, that is, it has to be executed completely, or aborted.

Suppose one of the updates  $(U_k, V_k)$ ,  $1 \leq k \leq n$  is not side-effect free. Since  $u$  is an atomic operation, it needs all of its  $n$  components to work correctly to be considered successful. Consequently, if update  $(U_k, V_k)$  fails,  $u$  also fails. Additionally, if update  $(U_k, V_k)$  is not side-effect free, then  $u$  is also not side-effect free. ■

Based on theorem 4.3, we can now answer a more general question: Is there a class of query tree views for which all possible updates are side-effect free? To answer this question, we summarize the results of [13] and [21] for conditions under which NRA views are updatable, and generalize them for XML views constructed by query trees.

**Insertions.** An insertion over an NRA view is side-effect free when the corresponding relational view  $V$  is a select-project-join view, the primary and foreign keys of the source relations of  $V$  are in the view and joins are made only through foreign keys. In terms of query trees, this means that the primary keys of the source relations of  $qt_i$  must appear as values in leaf nodes of  $qt_i$  and the *where* annotations in  $qt_i$  specifies joins using foreign keys, for all split trees  $qt_i$  corresponding to a query tree  $qt$ .

**Deletions and modifications.** Deletions and modifications over an NRA view  $V$  are side-effect free when the above conditions for insertions are met and  $V$  is *well-nested* [13]. By *well-nested*, we mean that the source relations in  $V$  must be nested according to key-foreign key constraints in the underlying relations. We rephrase this condition in terms of query trees as follows:

**Definition 4.4** *A query tree  $qt$  is well-nested if for any two source relations  $R$  and  $S$  in  $qt$ , if  $S$  is related to  $R$  by a foreign key constraint then the source annotation for  $R$  occurs in an ancestor of the node  $s$  containing the source annotation for  $S$ . Additionally, attributes of  $R$  must not appear as values in the descendants of  $s$ .*

The results above identify three classes of updatable XML views: one that is updatable for all possible insertions; one that is updatable for all possible insertions, deletions and modifications; and a general one whose updatability with respect to a given update can be reasoned about using theorem 4.3. Furthermore, we can now prove the following:

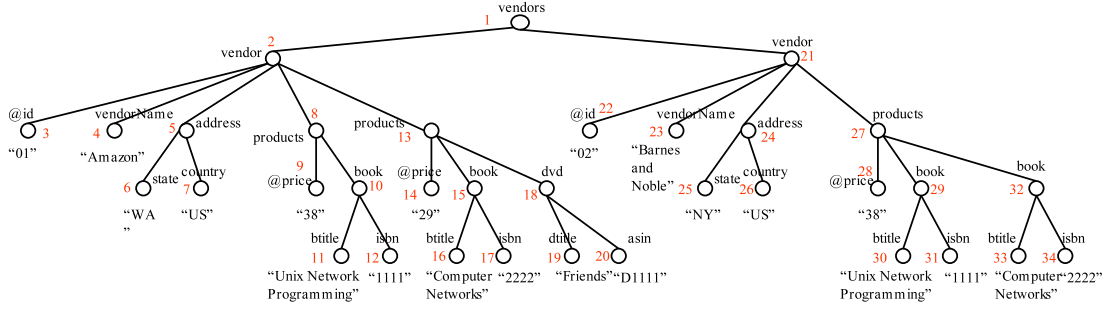


Figure 8: View 3: books and dvds clustered by price

**Theorem 4.4** Given a query tree  $qt$  with the restrictions mentioned above and defined over a BCNF database  $\mathcal{D}$ , then for any instance  $d$  of  $\mathcal{D}$  and correct update  $u$  over  $qt$ :  $\text{apply}(\text{eval}(qt,d), u) \equiv \text{eval}(qt, d')$ , where  $d'$  is the updated relational database state resulting from the application of the translated view updates  $\{U_{11}, \dots, U_{1m_1}, \dots, U_{n1}, \dots, U_{nm_n}\}$  using the techniques of [21].

We leave the study of updatability using other existing relational techniques for future work.

## 5 Extending Query Trees to Support Grouping

Figure 8 shown an example of XML view with group nodes. It is analogous of that of Figure 2, with the exception that now *books* and *dvds* are clustered by *price* under *product*.

To support grouping, we make the following changes to the definition of query trees:

**Definition 5.1** Nodes of query trees can be of three types:

**Leaf nodes** have a value (to be defined), which is either projected or grouped. Names of leaf nodes that start with “@” are considered to be XML attributes.

**Starred nodes** (nodes whose incoming edge is starred) may have one or more source annotations and zero or more where annotations (to be defined). An exception is made for starred nodes with group children, which must have no source annotation.

A **Group node** (one that has a grouped value) must have siblings that are starred nodes or group nodes of a restricted form (see definition 5.2).

**Definition 5.2** The value of a node  $n$  can be projected or grouped.

A **projected value** is of form  $\$x/A$ , where  $A \in \mathbb{A}_T$  and  $\$x$  is bound to table  $T$  by a source annotation on  $n$  or some ancestor of  $n$ . Projected values must be unique, e.g.  $\$x/A$  cannot occur twice within the tree.

A **grouped value** is of form  $\text{GROUP}(\$x_1/A_1 \mid \dots \mid \$x_m/A_m)$ , where  $m \geq 1$  and  $A_i \in \mathbb{A}_{T_i}$  and  $\$x_i$  is bound to  $T_i$  by a source annotation on a sibling node of  $n$ . The domains of  $A_1, \dots, A_m$  in  $\mathcal{D}$  must be the same. Group nodes with the same parent must be defined over the same set of variables  $x_1, \dots, x_m$ , and must have  $m$  siblings  $b_1, \dots, b_m$  whose incoming edges are starred<sup>2</sup>. Furthermore, the parent of node  $n$  must be starred, and it must have no source annotations.

The intuition behind multiple group nodes with the same parent is to allow tuples to be clustered based on a set of attributes rather than a single attribute.

Additionally, it is necessary to add another starred abstract type to our set of abstract types, so we can distinguish grouped nodes. We call this type  $\tau_G$ . Nodes of type  $\tau_G$  are identified as follows: each starred node which has one or more GROUP children has abstract type  $\tau_G$ .

As an example of query tree which uses group nodes, consider the query tree shown in Figure 9. It is the query tree that generates the XML view of Figure 8. Notice that there is a node *price* whose value is grouped:  $\text{GROUP}(\$sb/price \mid \$sd/price)$ .

<sup>2</sup>Notice that we do not require that  $(\$x_1/A_1 \mid \dots \mid \$x_m/A_m)$  in the group operation be in the same order as  $b_1, \dots, b_m$ .

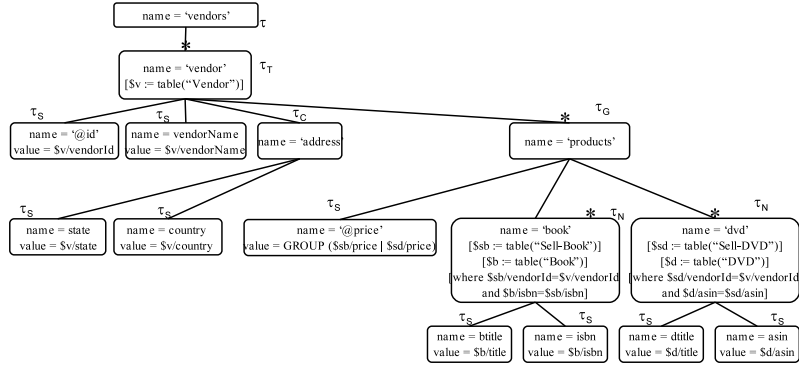


Figure 9: Query tree with grouped values

Given these definitions, we now show how the *eval* and *split* algorithms are modified to support grouping. The eval algorithm has to have an additional function called *group*, which deals with the generation of grouped nodes. Notice that this is the only difference between the algorithm below and algorithm 1.

```
eval(qt, d)
evaluate(root(qt, d))
```

```
evaluate(n, d)
Let bindings{} be a hash array of bindings of variable attributes to values, initially empty.
case abstract_type(n)
  τ|τ_C: buildElement(n)
  τ_T|τ_N: table(n)
  τ_G: group(n)
  τ_S: print "<name(n)>value(n)</name(n)>"
end case
```

```
buildElement(n)
let tag = "name(n)"
for each attribute c in children(n) do
  add "name(c) = value(c)" to tag
end for
print "<tag>"
for each non-attribute c in children(n) do
  eval(c)
end for
print "</name(n)>"
```

```
table(n)
let w be a list of conditions in sources(n)
for each w[i] do
  if w[i] involves a variable v in bindings{} then
    substitute the value binding{v} for v
  end if
end for
calculate the set B of all bindings for variables in sources(n) that makes the conjunction of the modified w[i]'s true
for each b in B do
  add b to bindings{}
  buildElement(n)
  remove b from bindings{}
end for
```

```
group(n)
let g1, ..., gs be the GROUP children of n
let w be a list of conditions in sources(m), for all starred nodes m that are children of n
for each w[i] do
  if w[i] involves a variable v in bindings{} then
    substitute the value binding{v} for v
  end if
end for
calculate the set B of all bindings for variables in sources(m) (for all starred nodes m that are children of n) that makes the conjunction of the modified w[i]'s true
let V1 = ∪_i values of i'th group term in g1, taken from B
let ...
let Vs = ∪_i values of i'th group term in gs, taken from B
for each v1 in V1 do
  add variable bindings xi/p= value(g1) for each group variable xi to bindings{}
  for each vs in Vs do
    add variable bindings xi/p= value(gs) for each group variable xi to bindings{}
  buildElement(n)
```

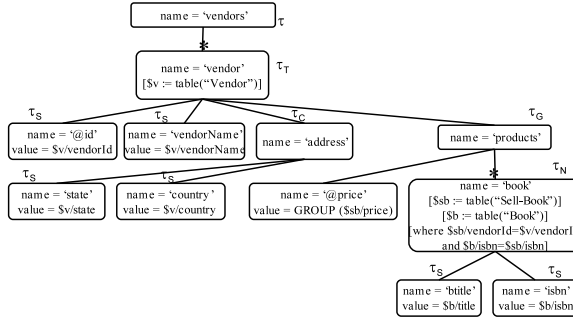


Figure 10: Partitioned query tree for  $\tau_N(\text{book})$

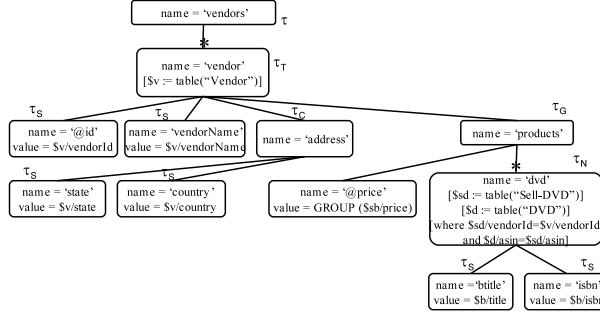


Figure 11: Partitioned query tree for  $\tau_N(\text{dvd})$

```

remove variable bindings  $x_i/p = \text{value}(g_s)$  for each group variable  $x_i$  in bindings{ }
end for
remove variable bindings  $x_i/p = \text{value}(g_1)$  for each group variable  $x_i$  in bindings{ }
end for

```

The *split* algorithm needs to take care of group values. It needs to remove parts of the value of group nodes, so that variable references are correct in each split tree. As an example, the query tree of Figure 8 has a group node *price* whose value is  $GROUP(\$sb/price / \$sd/price)$ . The two split trees generated by the split algorithm will have a node *price* referencing just one of the variables each ( $\$sb$  or  $\$sd$ ), as shown in Figures 10 and 11. The modified *split* algorithm is shown below:

```

split(qt)
Let  $t[]$  be an array of query trees, initially empty
Let  $i = 0$ 
Let  $N$  be the set of nodes of type  $\tau_N$  in  $qt$ 
for each node  $n$  in  $N$  do
  inc  $i$ 
  {initialize  $t[i]$  with  $qt$ }
   $t[i] = qt$ 
  repeat
    delete from  $t[i]$  all subtrees rooted at a node  $z$  of type  $\tau_N$ , where  $z \neq n$ 
    retype the ancestors of the deleted nodes
  until  $n$  is the only node of type  $\tau_N$  in  $t[i]$ 
  for each group node  $g$  in  $t[i]$  do
    delete from  $g$  all the variable references not declared as source annotations in its starred sibling
  end for
end for
end for

```

As for the *map* algorithm, the only change that needs to be made is to add lines 64, 65 and 66, since nodes of type  $\tau_G$  are starred nodes, but they do not carry any source or where annotation.

1: map( $qt[]$ )

```

2: Let  $sql[]$  be an array of strings, initially empty; Let  $numberqt$  be the number of split trees in  $qt[]$ 
3: for  $k$  from 1 to  $numberqt$  do
4:   Let  $n$  be the node of type  $\tau_N$  in  $qt[k]$ 
5:    $sql[k] = \text{"CREATE VIEW " + name}(n) + \text{" AS "}$ 
6:    $sql[k] = sql[k] + \text{"SELECT "}$ 
7:   Let  $N$  be the list of leaf nodes in  $qt[k]$ 
8:   for  $i$  from 1 to  $size(N)$  do
9:     get next  $n$  in  $N$ 
10:    if  $i > 1$  then
11:       $sql[k] = sql[k] + \text{" , " + variable}(n) + \text{" . " + attribute}(n) + \text{" AS " + name}(n)$ 
12:    else
13:       $sql[k] = sql[k] + variable}(n) + \text{" . " + attribute}(n) + \text{" AS " + name}(n)$ 
14:    end if
15:     $i = i + 1$ 
16:  end for
17:   $sql[k] = sql[k] + \text{" FROM "}$ ; Let  $from = \text{" "}$ ; Let  $N$  be the set of starred nodes in  $qt[k]$ 
18:  for each  $n$  in  $N$  do
19:    Let  $join = \text{" "}$ ; Let  $S$  be the list of source annotations in  $n$ ; Let  $W$  be the list of where annotations in  $n$ 
20:    for  $i = 1$  to  $size(S)$  do
21:      get next  $s$  in  $S$ 
22:       $join = join + table}(s) + \text{" AS " + variable}(s)$ 
23:      if  $i < size(S)$  then
24:         $join = join + \text{" INNER JOIN "}$ 
25:      end if
26:       $i = i + 1$ 
27:    end for
28:    Let  $count = 0$ 
29:    for  $i = 1$  to  $size(W)$  do
30:      get next  $w$  in  $W$ 
31:      if  $w$  is of the form  $\$x/A op \$y/B$  AND  $\$x$  is bound to table  $X$  by a source annotation  $s \in S$  AND  $\$y$  is bound to table  $Y$  by a source annotation  $s' \in S$  then
32:        if  $count = 0$  then
33:           $join = join + \text{" ON " + x.A op y.B}$ 
34:        else
35:           $join = join + \text{" AND " + x.A op y.B}$ 
36:        end if
37:         $i = i + 1$ ;  $count = count + 1$ 
38:      end if
39:    end for
40:    if  $count = 0$  then
41:       $join = join + \text{" ON (I=1) "}$ 
42:    end if
43:    if  $size(S) > 1$  then
44:       $join = "(" + join + ")"$ 
45:    end if
46:    Let  $A$  be the set of starred ancestors of  $n$ ; Let  $count = 0$ 
47:    if  $n$  has a starred ancestor then
48:       $join = \text{" LEFT JOIN " + join}$ 
49:      for  $i = 1$  to  $size(W)$  do
50:        get next  $w$  in  $W$ 
51:        if  $w$  is of the form  $\$x/A op \$y/B$  AND ( $\$x$  is bound to table  $X$  on node  $n$  AND  $\$y$  is bound to table  $Y$  on a node  $a$  in  $A$ ) OR ( $\$x$  is bound to table  $X$  on a node  $a$  in  $A$  AND  $\$y$  is bound to table  $Y$  on node  $n$ ) then
52:          if  $count = 0$  then
53:             $join = join + \text{" ON " + x.A op y.B}$ 
54:          else
55:             $join = join + \text{" AND " + x.A op y.B}$ 
56:          end if
57:          if
58:             $i = i + 1$ ;  $count = count + 1$ 
59:          end for
60:          if  $count = 0$  then
61:             $join = join + \text{" ON (I=1) "}$ 
62:          end if
63:           $from = "(" + from + join + ")"$ 
64:        else
65:          if  $abstract\_type}(n) \neq \tau_G$  then
66:             $from = from + join$ 
67:          end if
68:        end if
69:      end for
70:       $sql[k] = sql[k] + from$ ; Let  $W'$  be the set of all where annotations on nodes of  $qt[k]$ . Let  $count = 0$ 
71:      for each  $w'$  in  $W'$  do
72:        if  $w'$  is of the form  $\$x/A op Z$  AND  $Z$  is an atomic value then
73:          if  $count = 0$  then
74:             $sql[k] = sql[k] + \text{" WHERE " + x.A op Z}$ 
75:          else
76:             $sql[k] = sql[k] + \text{" AND " + x.A op Z}$ 
77:          end if
78:        end if
79:      end for
80:    end for
81:  return  $sql[]$ 

```

## 5.1 Updatability

Extending query trees with group values reflects in our updatability study. Specifically, we present a rewritten version of theorem 4.3 when group nodes are considered.

**Theorem 5.1** *A correct update  $u$  to an XML view defined by a query tree  $qt$  is side-effect free if:*

1.  $u$  does not modify the leaf child of a  $\tau_G$  node, or in other words,  $ref$  does not point to a group node;
2.  $u$  does not delete a starred child of a  $\tau_G$  node; and
3. For all  $(U_i, V_i)$ , where  $V_i$  is the corresponding relational view of  $qt_i$  and  $U_i$  is the translation of  $u$  over  $V_i$ ,  $U_i$  is side-effect free in  $V_i$ .

*Proof:* We divide the proof in three steps, one for each condition in the theorem.

**Condition 1:** To prove this condition, all we need to do is to show an example of a modification of a group node that causes side-effects. Suppose we specify a modification over the view in Figure 8 by  $ref = /vendor/vendor[id="1"]/products[@price="38"]/@price$  and  $\Delta = \{29\}$ . The evaluation of  $ref$  yields node 8. Although it seems fine to modify the value of this node, the reconstructed XML view would collapse the subtree rooted at node 13 with the subtree rooted at node 8. This happens because we are changing the value of node 8 to a value that was already in the view, and the semantics of GROUP requires that nodes that agree in the value of  $price$  should be collected together. As a consequence, the XML view modified by the user will be different from the reconstructed view – a side-effect.

**Condition 2:** This condition can also be proved by a contra example. Consider a deletion over the view in Figure 8 with update path  $ref = /vendor/vendor[id="1"]/products[@price="38"]/book$ , which evaluates to node 10. The deletion of this book will also make the subtree rooted at node 8 ( $products$ ) disappear. This is because node 10 was the only book being sold by this price under this vendor. Consequently, the update is not side-effect free.

**Condition 3:** This condition is the statement of theorem 4.3 itself. Please refer to that theorem for proof. ■

Condition 3 states that if an update passes the “grouping conditions” (conditions 1 and 2), then the update is side-effect free if all updates in its translation onto the underlying relational views are side-effect free. Hence, any side-effect free relational view update technique could be used.<sup>3</sup>

Recall that reasoning about whether or not an update is side-effect free involves the query tree rather than the resulting XML instance. There may therefore be instances for which updates to grouped nodes do *not* cause other updates to be introduced in the recomputed view. Examples include changing the value of node 9 to \$20 in figure 8 and deleting the subtree rooted at node 18. These (desirable) updates are outlawed in our approach.

There are two ways in which we could allow the desired updates on group nodes above: (1) perform instance analysis to catch exactly the cases that produce side-effects; or (2) allow side-effects in these special cases, or re-define side-effects to exclude empty groups or groups which collapse. We leave this for future work.

## 6 Evaluation

For purposes of presentation, the query tree language presented in this paper was kept simple to highlight how the mapping of the query tree and updates are performed.

Query trees can be extended in a number of ways, for example to deal with grouping, aggregates, function applications and so on. An example of such extension can be found in section 5, where we allow *grouped values* which allow tuples that agree on a given value to be clustered together, as well as leaf nodes with attributes.

However, another consideration that must be kept in mind when extending the language is whether or not the relational views resulting from the XML view are updatable. The language presented in this paper, with suitable restrictions on the way in which joins and nesting are performed with respect to keys and foreign keys in the underlying relational database, presents a subset of XQuery in which *side-effect free* updates can be defined as discussed in the previous section. While grouped values and leaf nodes with attributes do not affect these results, the addition of functions and aggregates would. Analogous to work on updating views in relational databases which restricts views to select-project-join queries, we have therefore initially decided against considering a richer language (although we plan to do so in future work).

<sup>3</sup>The term *exact translation* used in [25] is equivalent to our notion of side-effect free.

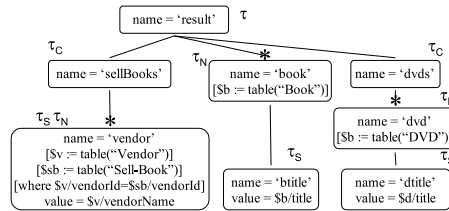


Figure 12: Example of query tree

The EBNF for the subset of XQuery corresponding to our language (with grouped values) can be found in the section 6.2.

To evaluate our language, we first discuss the restrictions in our form of queries, and what query trees can or cannot express. Second, we examine the power of expression of query trees, and compare it with existing proposals in literature. We have also analyzed the “practicality” of XML views constructed by query trees by collecting examples of real XML views extracted from relational databases and evaluating whether or not query trees can capture them. For these real XML views, query trees were sufficiently expressive.

## 6.1 Limitations of Query Trees

Although query trees are quite expressive, there are some restrictions.

**Values must come from the relational database.** We do not allow constants to be introduced as values in leaves, nor do we allow functions to calculate new values from values in the database. Allowing constant values in leaves is potentially useful (for example, to add a version number to the view), but they are not interesting from the perspective of updates to the relational database nor can they themselves be updated since they are not part of the database schema. Calculating a value from a set of values (e.g. taking the average of a relational column) creates a one to many mapping which cannot be updated; research on relational views also disallows this case. However, calculating a new value from a single value in the database (e.g. translating length in centimeters to length in inches) could be allowed as long the reverse function was also specified.

**Queries are trees rather than graphs.** This restriction disallows recursive queries, which are also disallowed in SilkRoute [23]. For example, suppose the relational database contained a relation Patriarchs(PName, CName) with instance {(John, Marc), (John, Chris), (Justin, John)}. An XML view of this that one might wish to construct would be:

```
<Patriarch>
  <Name>Justin</Name>
  <Children>
    <Name>John</Name>
    <Children> <Name>Marc</Name>
               <Name>Chris</Name>
    </Children>
  </Children>
</Patriarch>
```

Since recursive queries cannot be mapped to select-project-join queries, our technique would have to be extended significantly to reason about them.

On the other hand, query trees are flexible enough to represent heterogeneous structures (e.g. the view in Figure 5). It can also represent query trees with a repeating leaf node, as shown in Figure 12 (note that *vendor* is labeled with  $\tau_N$  and  $\tau_S$ ). The XML view resulting from this query tree is as follows:

```
<result>
  <sellBooks>
    <vendor>Amazon</vendor>
    <vendor>Barnes and Nobel</vendor>
  </sellBooks>
  <book><btitle>Unix Network Programming</btitle></book>
  <book><btitle>Computer Networks</btitle></book>
  ...
  <dvds>
    <dvd><dtitle>Friends</dtitle></dvd>
    ...
  </dvds>
</result>
```

It turns out that XML views with heterogeneous content and repeating leaves arise frequently in practice, but that recursive views are not common. We therefore believe that the above restrictions do not limit the usefulness of our approach.

## 6.2 Power of Expression

We have shown how to transform the XML view update problem into the well studied relational view update problem. However, since our proposal is based on *query trees*, the next question is: Are query trees expressive enough to be used in practice? To answer this question, we compare the power of expression of query trees with SilkRoutes' *view forests* [23]; XPERANTO [32]; and DB2 DAD files [17].

Since most of these proposals are based on the XQuery [10] query language, we first show how a query tree can be translated to an XQuery query. Based on the structuring rules of query trees and on these translation rules, we then present a BNF for the subset of XQuery that query trees are capable of expressing.

Given a query tree  $qt$ , an XQuery  $xq$  is generated by the *generateXQuery* algorithm (algorithm 9). The algorithm is recursive, and it starts with  $n$  being the root of  $qt$ . Function  $value(n)$  returns the value associated with a leaf node. For example, suppose node  $n$  has value  $\$x/A$ , then  $value(n)$  returns the expression  $\$x/A$ . Function  $name(n)$  returns the node name in  $qt$ . When  $n$  is an attribute, the function returns the name without "@". As an example, if  $n$  has name  $@id$ , then  $name(n)$  returns  $id$ .

As an example, the XQuery query corresponding to the query tree of Figure 5 is as follows:

```
<vendors>
  {for $v in document("Vendor.xml")//row
   return
   <vendor id='{ $v/vendorId/text() }'>
     <vendorName>{ $v/vendorName/text() }</vendorName>
     <address>
       <state>{ $v/state/text() }</state>
       <country>{ $v/country/text() }</country>
       <web>
         <url>{ $v/url/text() }</url>
       </web>
     </address>
     {let $sb' := document("Sell-Book.xml")//row
      let $b' := document("Book.xml")//row
      let $sd' := document("Sell-DVD.xml")//row
      let $d' := document("DVD.xml")//row
      for $price in distinct-values($sb'/price | $sd'/price)
      return
      <products price='{ $price/text() }'>
        {for $sb in document("Sell-Book.xml")//row
         for $b in document("Book.xml")//row
         where $sb/vendorId = $v/vendorId and
              $b/isbn = $sb/isbn and
              $price = $sb/price
          return
          <book>
            <bttitle>{ $b/title }</bttitle>
            <isbn>{ $b/isbn }</isbn>
          </book>
        }
        {for $sd in document("Sell-DVD.xml")//row
         for $d in document("DVD.xml")//row
         where $sd/vendorId = $v/vendorId and
              $d/asin = $sd/asin and
              $price = $sd/price
          return
          <dvd>
            <dttitle>{ $d/title }</dttitle>
            <asin>{ $d/asin }</asin>
          </dvd>
        }
      </products>
     }
   </vendor>
 }
</vendors>
```

This translation algorithm assumes that each relational table  $X$  with attributes  $A, B, C, \dots$  is exported to XML as follows:

```
<X>
  <row>
    <A> ... </A>
    <B> ... </B>
    <C> ... </C>
    ...
  </row>
  ...
</X>
```

```

generateXQuery(n)
case abstract_type(n)
   $\tau | \tau_G$ : buildElement(n)
   $\tau_T | \tau_N$ : table(n)
   $\tau_G$ : group(n)
   $\tau_S$ : leaf(n)
end case

buildElement(n)
let tag = "name(n)"
for each attribute c in children(n) whose value is grouped do
  add "name(c) = '$ name(c)/text()'" to tag
end for
for each other remaining attribute c in children(n) do
  add "name(c) = 'value(c)/text()'" to tag
end for
print "<tag >"
for each non-attribute c in children(n) do
  generateXQuery(c)
end for
print "</name(n)>"

leaf(n)
if n has a grouped value then
  print "<name(n)>{ $name(n)/text()}</name(n)>"
else
  print "<name(n)>{ value(n)/text()}</name(n)>"
end if

table(n)
print "{"
for each source annotation binding a table X to a variable $x in n do
  print "for $x in document('X.xml')/row"
end for
Let W be the set of where annotations in n
Let count = 1
for each w in W do
  if count > 1 then
    print "AND w"
  else
    print "WHERE w"
  end if
  count = count + 1
end for
if n is a child of a node a, and abstract_type(a) =  $\tau_G$  then
  let  $G = \{g_1, \dots, g_s\}$  be the GROUP children of n
  for each  $g_i$  in G do
    Let value( $g_i$ ) be of the form GROUP( $\$x_1/A_1 | \dots | \$x_k/A_k$ )
    Find each of the  $x_i$  in the value of  $g_i$  that is declared in a source annotation in n
    if count > 1 then
      print "AND $ name( $g_i$ ) =  $\$x_i/A_i$ "
    else
      print "WHERE $ name( $g_i$ ) =  $\$x_i/A_i$ "
    end if
    count = count + 1
  end for
end if
print "return"
buildElement(n)
print "}"
group(n)
let  $G = \{g_1, \dots, g_s\}$  be the GROUP children of n
let S be the set of source annotations in m, for all starred nodes m that are children of n
print "{"
for each s in S binding a variable $x to a table X do
  print "let  $\$x'$  := document('X.xml')/row" {Notice that variable $x is primed in the generated let expression}
end for
for each  $g_i$  in G do
  Let value( $g_i$ ) be of the form GROUP( $\$x_1/A_1 | \dots | \$x_k/A_k$ )
  print "for $ name( $g_i$ ) in distinct values ( $\$x'_1/A_1 | \dots | \$x'_k/A_k$ )" {Notice again the use of primed variables. They correspond to the variables bound by the let expression on line (67)}
end for
buildElement(n)
print "}"

```

**Algorithm 9:** The *generateXQuery* algorithm

According to the translation algorithm and to the structuring rules of query trees, the XQuery queries corresponding to query trees follows the following EBNF:

- [1] XQuery ::= QueryBody  
[2] QueryBody ::= ElmtConstructor

```

[3] ElmtConstructor ::= "<" QName AttList "/>" | "<" QName AttList? ">" ElmtContent+ "</" QName ">"
[4] ElmtContent   ::= ElmtConstructor | EnclosedExpr+
[5] AttList       ::= ((QName "=" AttValue)?) +
[6] AttValue      ::= (' ' AttValueContent ' ') | (' ' AttValueContent ' ')
[7] AttValueContent ::= "{ " PathExprAtt "}"
[8] PathExprAtt  ::= "$ VarName "/" QName "/" NodeTest
[9] VarName       ::= QName
[10] EnclosedExpr ::= "{ " (FWRExpr | LFWRExpr | PathExpr) "}"
[11] Expr         ::= OrExpr
[12] OrExpr       ::= AndExpr ( "or" AndExpr ) *
[13] AndExpr      ::= ComparisonExpr ( "and" ComparisonExpr ) *
[14] FWRExpr      ::= ((ForClause)+ WhereClause? OrderByClause? "return") * ElmtConstructor
[15] LFWRExpr     ::= ((LetClause)+ (ForClauseDistinct)+ WhereClause? OrderByClause? "return") * ElmtConstructor
[16] ComparisonExpr ::= ValueExpr (GeneralComp ValueExpr ) ?
[17] ValueExpr    ::= PathExpr | PrimaryExpr
[18] PathExpr     ::= "$ VarName "/" QName ("/" NodeTest)?
[19] NodeTest     ::= TextTest
[20] TextTest     ::= "text" "(" " )"
[21] ForClause    ::= "for" "$ VarName "in" DocExpr
[22] ForClauseDistinct ::= "for" "$ VarName "in distinct-values" UnionExpr
[23] DocExpr      ::= "table (" ' ' QName ' ' " )//row" | "table (" ' ' QName ' ' " )//row"
[24] WhereClause  ::= "where" Expr
[25] GeneralComp  ::= "=" | "!=" | "<" | "<=" | ">" | ">="
[26] OrderByClause ::= "order" "by" OrderSpecList
[27] OrderSpecList ::= OrderSpec ( " , " OrderSpec ) *
[28] OrderSpec    ::= PathExpr
[29] PrimaryExpr  ::= Literal | ParenthesizedExpr
[30] Literal       ::= NumericLiteral | StringLiteral
[31] NumericLiteral ::= IntegerLiteral | DecimalLiteral | DoubleLiteral
[32] ParenthesizedExpr ::= "(" Expr? ")"
[33] UnionExpr    ::= "(" "$ VarName "/" QName ( "union" | "|" ) "$ VarName "/" QName ) * ")"

```

XPERANTO [32] can express all queries in XQuery. View forests [23] are capable of expressing any query in the XQueryCore that does not refer to element order, use recursive functions or use *is/is not* operators. Query trees present the same limitations as [23], and are also not capable of expressing *if/then/else* expressions; sequence of expressions (since we require that the result of the query always be an XML document); function applications; arithmetic and set operations. Input functions are also a limitation of query trees. It is not possible to bind results of expressions to variables. Variables can only be bound to relational tables, while in SilkRoute, they can be bound to arbitrary expressions.

DB2 XML Extender provides mappings from relations to XML through DAD files. Mappings can be done in two ways: using a single SQL statement (by using the `SQL_stmt` element in the DAD file), or using the `RBD_node` mapping. The `SQL_stmt` method allows only a single SQL statement, so XML views with heterogeneous structures (like the one in Figure 5) can not be constructed. The `RBD_node` method allows heterogeneous structures, since instead of specifying a single SQL statement for the XML extraction, the user specifies, for each XML element or attribute in the XML view, the table and attribute name from which the data must be retrieved. It is also possible to specify conditions for each XML node in the DAD file (join conditions and selection conditions). DB2 DAD files with `RBD_node` method are equivalent to query trees in expressive power, since all the data come directly from the relational database and functions cannot be applied over the retrieved data. This is meaningful, since DB2 DAD files represent features that are useful in practice, and because this subset can easily be mapped to relational views.

### 6.3 Real applications

We were able to obtain three real world applications: the XBrain project [5], a Tobacco company [27] and the Mondial Database [30]. The views produced by these applications are available at [1].

The XBrain Project is an application of SilkRoute [23]. The application queries a brain mapping database, over which an XML *public view* is defined.

The Tobacco company example is the most interesting. The company needs to send monthly reports to a Tobacco Producer's Association<sup>4</sup>. The report shows data about the producers from whom the company bought tobacco, as well as prices, quantities, etc. The company stores all the transactions in a relational database, and at the end of the month it generates an XML report containing all the information solicited by the Producer's Association.

The Mondial database is a case study for information extraction and integration. Facts about global geography are extracted from the Web and integrated in a very large database. An XML view over this database is then provided.

These applications highlight several characteristics of real XML views: They are large and complex, typically involving several relational tables. They are also well structured, which is no surprise since the source data

<sup>4</sup>We omit the company and the association name for copyright reasons.

is relational and therefore structured. Additionally, joins are made through keys and foreign keys - a desired property of updatable views. All these views can be expressed using query trees except for the portion of the Mondial view shown below:

```
<religions percentage="70">Muslim</religions>
<religions percentage="10">Roman Catholic</religions>
<religions percentage="20">Albanian Orthodox</religions>
```

As presented so far, query trees cannot represent this XML view since *religions* has a child *percentage* but must be a leaf to have a value, a contradiction. The same problem occurs in the Mondial view for an element called *ethnicgroups*. Note that allowing attributes within atomic elements is similar to allowing mixed content elements.

It turns out that query trees can easily be extended to deal with this case. We can introduce a special attribute called *textContent* to represent the text content of the element *religions*. The tree will be processed with this additional attribute, and transformed back to its original structure before being presented to the user.

There are also three problematic attributes in the Mondial view: attribute *membership* of element *country*, and the attributes *id* and *is\_country\_capital*. These attributes require computation to construct their value, and are therefore not addressed by our work as discussed earlier in this section.

## 6.4 Normalized XML documents

A proposal to extract a nested normalized XML document [8] from a relational database is presented in [40]. The proposal explores keys and foreign key constraints to build a graph of dependencies between tables. By traversing this graph, it is possible to decide which table(s) will be a top-level element and how the remaining tables will be nested under this table. When nesting a given table leads to redundancy, it is placed directly below the root, and relationships are expressed using IDs and IDREFs. Such views can be easily expressed using query trees, and are updatable for all correct insertions, deletions and modifications.

## 6.5 XQuery use cases

For completeness, we also analyzed the relational use cases of XQuery [15]. Of the eighteen queries presented in [15], our query trees are capable of expressing only two (Q3 and Q4). This is mainly caused by the use of aggregate operations. We believe that these use cases highlight the difference between queries and view definitions, rather than demonstrating shortcomings of query trees. Aggregate operators and specialized functions are typically not considered in work on updating views.

## 6.6 XML documents stored in relations

XML documents are frequently mapped to relational databases for efficient storage. We analyzed the most popular approach, hybrid inlining [34], to check if the XML view definitions resulting from this mapping could be expressed by query trees.

We analyzed six different XML documents. Three of them represent information about courses of different universities [2]. We also analyzed the SIGMOD Record [4], the DBLP in XML [3] and the action XML file of XMark [6].

All of the three course documents are fully compatible with query trees. DBLP and XMark are also compatible except for one element that has mixed content (*title* in DBLP and *text* in XMark). Although hybrid inlining does not support mixed content, we mapped it by assuming an upper bound  $n$  on the number of fragments of text content for an element, and used special attributes called  $textContent_1, \dots, textContent_n$  to capture the fragmented text values. The resulting mapping could be expressed using query trees.

As for updates, since hybrid inlining introduces artificial primary keys, these keys must be projected in the view in order to update the underlying relational database through the reconstructed document. This can be handled by introducing an *id* attribute that holds the primary keys in elements that represent database tuples.

In summary, query trees are able to capture many of the "real" XML views of relational databases that we were able to find. The evaluations were also incredibly valuable to identify extensions that should be made to our definitions, such as those to handle atomic elements with attributes and mixed content elements.

## 7 Related Work

There are several proposals for exporting and querying XML views of relational databases [16, 23, 32, 33]. For updates, [39] presents a round trip case study, where XML documents are stored in relational databases, reconstructed and then updated. In this case, it is always possible to translate the updates back to the underlying relational database. Our approach differs since we address update *legacy* databases through XML views.

Commercial relational databases offer support for extracting XML data from relations as well as restricted types of updates. In SQL Server [19], an XML view generated by an annotated XML Schema can be modified using *updategrams*. To update, the user provides a before and after image of the XML view [20]. The system computes the difference between the images and generates SQL update statements. The views supported by this approach are very restricted: joins are through keys and foreign keys, and nesting is controlled to avoid redundancy. This corresponds to our well-nested query trees, which are therefore provably updatable with respect to all insertions, deletions and modifications. Oracle [22] offers the specification of an annotated XML Schema, but the only possible update is to insert an XML document that agrees with the schema. IBM DB2 XML Extender [17] requires that updates be issued directly in the relational tables.

Native XML databases also support updates [37, 24, 35]. The goal of all these systems differs from ours since they do not update through views.

## 8 Conclusions

In this paper, we present a technique for updating relational databases through XML views. The views are constructed using query trees, which allow nesting as well as heterogeneous sets of tuples, and can be used to capture mixed content, grouping, as well as repeating text elements and text elements with attributes.

The main contributions of this paper are the mapping of the XML view to a set of underlying relational views, and the mapping of updates on an XML view instance to a set of updates on the underlying relational views. By providing these mappings, the XML update problem is reduced to the relational view update problem and existing technique on updates through views [21, 25, 9, 29] can be leveraged. As an example, we show how to use the approach of [21] to produce side-effect free updates on the underlying relational database.

Another benefit of our approach is that query trees are agnostic with respect to a query language. Query trees represent an intermediate query form, and any (subset of an) XML query language that can be mapped to this form could be used as the top level language. In particular, we have implemented our approach in a system called *Pataxó* that uses a subset of XQuery to build the XML views and translates XQuery expressions into query trees as an intermediate representation [14]. Similarly, our update language represents an intermediate form that could be mapped into from a number of high-level XML update languages (using a static evaluation of which updates are to be performed). In our implementation, we use a graphical user interface which allows users to click on the update point or (in the case of a set oriented update) specify the path in a separate window and see what portions of the tree are affected.

In future work, we plan to study the updatability of XML views using other proposals of updates through relational views in literature. We also plan to extend the language to include other features such as aggregates, and to extend the model to include order.

**Acknowledgments** We would like to thank those who helped us getting data for our experiments: Marcelo Arenas, Byron Choi, Carina Dorneles, Juliana Freire, Alon Halevy, Zack Ives, Eduardo Kroth, Cristiano Leivas, Leonardo Murta, Dan Suciu, and Yifeng Zheng.

## References

- [1] <http://www.inf.ufrgs.br/vanessa/xml/evaluation>.
- [2] <http://www.cs.washington.edu/research/xmldatasets/www/>.
- [3] DBLP in XML. <http://dblp.uni-trier.de/xml/>.
- [4] SIGMOD record. <http://www.acm.org/sigs/sigmod/record/xml/>.
- [5] XBrain project. <http://quad.biostr.washington.edu/stang/>.

- [6] XMark - an XML benchmark project. <http://monetdb.cwi.nl/xml/downloads.html>.
- [7] ABITEBOUL, S., QUASS, D., MCHUGH, J., WIDOM, J., AND WIENER, J. The Lorel Query Language for Semistructured Data. *International Journal on Digital Libraries* 1, 1 (1997), 68–88.
- [8] ARENAS, M., AND LIBKIN, L. A normal form for XML documents. In *Proceedings of PODS 2002* (Madison, Wisconsin, Jun 2002).
- [9] BANCILHON, F., AND SPYRATOS, N. Update semantics of relational views. *ACM Transactions on Database Systems* 6, 4 (Dec 1981).
- [10] BOAG, S., CHAMBERLIN, D., FERNANDEZ, M. F., FLORESCU, D., ROBIE, J., AND SIMÉON, J. XQuery 1.0: An XML Query Language. W3C Working Draft, May 2003.
- [11] BOHANNON, P., GANGULY, S., KORTH, H., NARAYAN, P., AND SHENOY, P. Optimizing view queries in ROLEX to support navigable result trees. In *Proceedings of VLDB 2002* (Hong Kong, China, Aug. 2002).
- [12] BONIFATI, A., BRAGA, D., CAMPI, A., AND CERI, S. Active xquery. In *ICDE* (San Jose, California, feb 2002).
- [13] BRAGANHOLO, V., DAVIDSON, S., AND HEUSER, C. On the updatability of XML views over relational databases. In *Proceedings of WEBDB 2003* (San Diego, CA, June 2003).
- [14] BRAGANHOLO, V., DAVIDSON, S., AND HEUSER, C. UXQuery: building updatable XML views over relational databases. In *Brazilian Symposium on Databases* (Manaus, AM, Brazil, 2003), pp. 26–40.
- [15] CHAMBERLIN, D., FANKHAUSER, P., FLORESCU, D., MARCHIORI, M., AND ROBIE, J. XML Query Use Cases. W3C Working Draft, May 2003.
- [16] CHAUDHURI, S., KAUSHIK, R., AND NAUGHTON, J. On relational support for XML publishing: Beyond sorting and tagging. In *Proceedings of SIGMOD 2003* (San Diego, CA, Jun 2003).
- [17] CHENG, J., AND XU, J. XML and DB2. In *Proceedings of ICDE'00* (San Diego, CA, 2000).
- [18] CLARK, J., AND DEROSE, S. XML Path Language (XPath) Version 1.0. W3C Recommendation, Nov 1999.
- [19] CONRAD, A. A Survey of Microsoft SQL Server 2000 XML Features. MSDN Library. <http://msdn.microsoft.com/library/en-us/dnxml/html/xml07162001.asp>. Jul 2001.
- [20] CONRAD, A. Interactive microsoft SQL Server & XML online tutorial. <http://www.topxml.com/tutorials/main.asp?id=sqlxml>.
- [21] DAYAL, U., AND BERNSTEIN, P. A. On the correct translation of update operations on relational views. *ACM Transactions on Database Systems* 8, 2 (Sep 1982), 381–416.
- [22] EISENBERG, A., AND MELTON, J. SQL/XML is making good progress. *SIGMOD RECORD* 31, 2 (2002).
- [23] FERNÁNDEZ, M., KADIYSKA, Y., SUCIU, D., MORISHIMA, A., AND TAN, W.-C. Silkroute: A framework for publishing relational data in XML. *ACM Transactions on Database Systems (TODS)* 27, 4 (Dec 2002), 438–493.
- [24] JAGADISH, H. V., AL-KHALIFA, S., CHAPMAN, A., LAKSHMANAN, L. V., NIERMAN, A., PAPARIZOS, S., PATEL, J. M., SRIVASTAVA, D., WIWATWATTANA, N., WU, Y., AND YU, C. TIMBER: A native XML database. *The VLDB Journal* 11, 4 (2002), 274–291.
- [25] KELLER, A. M. Algorithms for translating view updates to database updates for views involving selections, projections, and joins. In *Proceedings of SIGMOD* (Portland, Oregon, Mar. 1985), ACM, pp. 154–163.
- [26] KELLER, M. The role of semantics in translating view updates. *IEEE Computer* 19, 1 (1986), 63–73.

- [27] KROTH, E. Personal communication, Jun 2003.
- [28] LAUX, A., AND MARTIN, L. XUpdate WD. Working Draft. <http://www.xmldb.org/xupdate/xupdate-wd.html>, Sep 2000.
- [29] LECHTENBÖRGER, J. The impact of the constant complement approach towards view updating. In *Proceedings of PODS 2003* (San Diego, CA, Jun 2003), pp. 49–55.
- [30] MAY, W. Information extraction and integration with FLORID: The MONDIAL case study. Tech. Rep. 131, Universität Freiburg, Institut für Informatik, 1999. <http://www.informatik.uni-freiburg.de/may/Mondial/>.
- [31] ROWE, L. A., AND SHOENS, K. A. Data abstraction, views and updates in rigel. In *SIGMOD* (Boston, Massachusetts, 1979), pp. 71–81.
- [32] SHANMUGASUNDARAM, J., KIERNAN, J., SHEKITA, E., FAN, C., AND FUNDERBURK, J. Querying XML views of relational data. In *Proceedings of VLDB 2001* (Roma, Italy, Sept. 2001).
- [33] SHANMUGASUNDARAM, J., SHEKITA, E. J., BARR, R., CAREY, M. J., LINDSAY, B. G., PIRAHESH, H., AND REINWALD, B. Efficiently publishing relational data as XML documents. *The VLDB Journal* (2000), 65–76.
- [34] SHANMUGASUNDARAM, J., TUFTE, K., ZHANG, C., HE, G., DEWITT, D. J., AND NAUGHTON, J. F. Relational databases for querying XML documents: Limitations and opportunities. *The VLDB Journal* (1999), 302–314.
- [35] SOFTWARE AG. Tamino XML Server. <http://www.softwareag.com/tamino/details.htm>, 2002.
- [36] TATARINOV, I., IVES, Z., HALEVY, A., AND WELD, D. Updating XML. In *Proceedings of SIGMOD 2001* (Santa Barbara, CA, May 2001).
- [37] THE APACHE SOFTWARE FOUNDATION. Apache Xindice. <http://xml.apache.org/xindice>, 2002.
- [38] TUCHERMAN, L., FURTADO, A. L., AND CASANOVA, M. A. A pragmatic approach to structured database design. In *VLDB* (Florence, Italy, Oct 1983), pp. 219–231.
- [39] WANG, L., MULCHANDANI, M., AND RUNDENSTEINER, E. A. Updating XQuery Views Published over Relational Data: A Round-trip Case Study. In *Proc. of XML Database Symposium* (Berlin, Germany, Sep 2003).
- [40] WEIS, M. Development of an algorithm for generating an optimal XML Schema from a given relational schema. Bachelor Thesis, University of Berufsakademie Stuttgart, Germany, Sep 2003.