

Updating Recursive XML Views of Relations

Byron Choi^{1,4} (蔡冠球), Gao Cong^{2,4} (丛 高), Wenfei Fan^{3,4} (樊文飞), and Stratis D. Viglas⁴

¹*Division of Information System, School of Computer Engineering, Nanyang Technological University, 639798, Singapore*

²*Microsoft Research Asia, Beijing 100080, China*

³*Bell Laboratories, Murray Hill, NJ07974-0636, U.S.A.*

⁴*University of Edinburgh, Edinburgh EH8 9LE, Scotland, U.K.*

E-mail: kkchoi@ntu.edu.sg; gaocong@microsoft.com; {wenfei, sviglas}@inf.ed.ac.uk

Received September 3, 2007; revised March 7, 2008.

Abstract This paper investigates the view update problem for XML views published from relational data. We consider XML views defined in terms of mappings directed by possibly recursive DTDs compressed into DAGs and stored in relations. We provide new techniques to efficiently support XML view updates specified in terms of XPath expressions with recursion and complex filters. The interaction between XPath recursion and DAG compression of XML views makes the analysis of the XML view update problem rather intriguing. Furthermore, many issues are still open even for relational view updates, and need to be explored. In response to these, on the XML side, we revise the notion of side effects and update semantics based on the semantics of XML views, and present efficient algorithms to translate XML updates to relational view updates. On the relational side, we propose a mild condition on SPJ views, and show that under this condition the analysis of deletions on relational views becomes PTIME while the insertion analysis is NP-complete. We develop an efficient algorithm to process relational view deletions, and a heuristic algorithm to handle view insertions. Finally, we present an experimental study to verify the effectiveness of our techniques.

Keywords XML, XML publishing, XML views, view update

1 Introduction

As a classical technical problem, view updates have been studied for relational databases for decades (see, e.g., [1–4]), and the techniques developed in that area have been introduced into commercial DBMSs^[5–7]. Recently, a number of systems have been developed for publishing relational data to XML^[5–10]. The published XML documents can be seen as *XML views* of the relational data. For all the reasons that updating data through its relational views is needed, it is also important to update relational databases through their XML views.

In this paper we study the *XML view update problem*, which can be stated as follows. Given an XML view of a relational database, we want to propagate updates of the XML view to the relational tables, without compromising the integrity of neither the XML nor the relational data. Formally put, given an XML view defined as a mapping $\sigma : \mathcal{R} \rightarrow D$ from relations of a schema \mathcal{R} to XML documents (trees) of a DTD D , a relational instance I of \mathcal{R} , the XML view $T = \sigma(I)$, and

updates Δ_X on the XML view T , we want to compute relational updates Δ_R such that $\Delta_X(T) = \sigma(\Delta_R(I))$. That is, the relational updates Δ_R , when propagated to XML via the mapping σ , yield the desired XML updates Δ_X on the view T .

While several commercial systems^[5–7] allow users to define XML views of relations, their support for XML view updates is either very restricted or not yet available. Previous work on XML view updates^[11] has addressed the problem by translating XML view updates to relational view updates and delegating the problem to the relational DBMS; however, most commercial DBMSs only have limited view-update capability^[5–7]. The state of the art in XML view update research^[12,13] solves the problem by explicitly focusing on *non-*recursively defined XML views and XML updates defined *without* recursive XPath queries. Though it is a complete solution, the restrictions posed in [13] are unfortunate since the recent proposals on XML update languages^[14,15] employ recursive XPath queries while DTDs (and thus XML view definitions) found in practice are often recursive^[16].

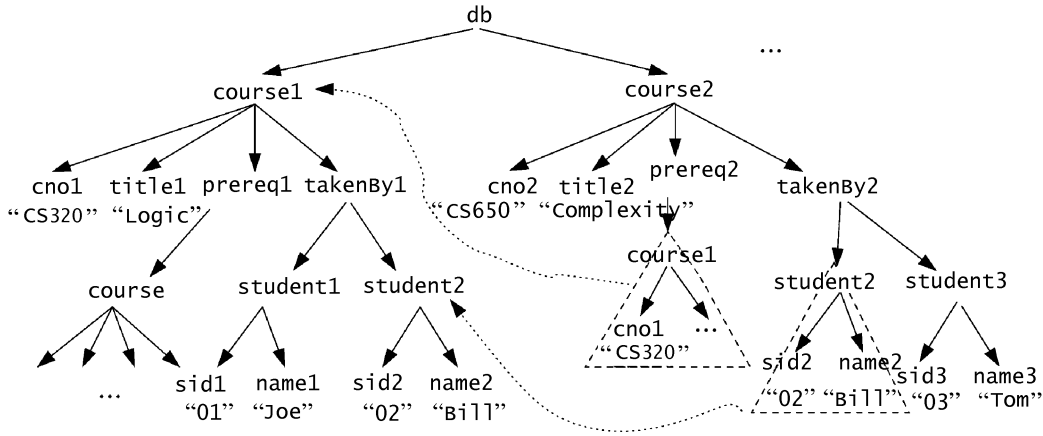


Fig.1. Example XML view.

In accordance with these requirements, we consider more general XML views and updates: possibly recursive XML view definitions and XML updates specified in terms of XPath expressions with recursion (descendant-or-self “//”) and complex filters, as illustrated by an example below.

Example 1. Consider a registrar database I_0 , which is specified by the relational schema R_0 (with keys underlined):

```

course(cno, title, dept)
project(cno, title, dept)
student(ssn, name)
enroll(ssn, cno)
prereq(cno1, cno2),

```

where a tuple $(c1, c2)$ in *prereq* indicates that $c2$ is a prerequisite of $c1$.

As depicted in Fig.1 (the dotted lines will be explained shortly), an XML view T_0 of the relational database is published for the CS department by extracting CS course-registration data from I_0 . The view is required to conform to the DTD D_0 below (the definition of elements whose type is PCDATA is omitted):

```

<!ELEMENT db (course*)
<!ELEMENT course (cno, title, prereq, takenBy)
<!ELEMENT prereq (course*)
<!ELEMENT takenBy (student*)
<!ELEMENT student (ssn, name)

```

Note that the view is defined recursively since the DTD D_0 is recursive (*course* is defined indirectly in terms of itself via *prereq*). Now consider an XML update $\Delta_X = \text{insert } T' \text{ into } P_0$ posed on the XML view T_0 , where P_0 is the (recursive) XPath query $\text{course}[cno=CS650]//\text{course}[cno=CS320]/\text{prereq}$, and T' is the subtree representing the course CS240. It is to find all the CS320 nodes below CS650 in T_0 and for each

CS320 node v , insert T' as a prerequisite of v . To carry out Δ_X , we need to find updates Δ_R on the underlying database I_0 such that $\Delta_X(T_0) = \sigma_0(\Delta_R(I_0))$.

Already a hard problem for relational views, the view update problem for XML views introduces several new challenges, which previous work^[11–13] on XML view updates cannot handle.

First, update semantics should be revised in the context of XML views of relations. Referring to the example above, the operation asks for inserting CS240 as a *prereq* of *only* those CS320 nodes below CS650; however, CS320 nodes also occur elsewhere below the root. As the XML view is published from the same relational database, all the courses, and therefore CS320, have unique *prereq* hierarchies. Such an insertion on selected paths of the hierarchy will result in *side effects* that need to be detected. In such a case, the users need to be consulted and, if they insist on carrying on updating, the semantics of insertion is revised such that the insertion will be performed at *every* CS320 node. Thus the insertion can accommodate side effects while being consistent with the semantics of the XML view. Note that such side effects are orthogonal to both the publishing middleware used and the storage scheme of the XML views. The details of side effects on deletions are even more subtle and call for a new semantics (see Section 2).

Second, the XML view may be *compressed* by storing each subtree shared by multiple nodes in the tree *only once*, as indicated in Fig.1 (replacing the subtrees in the dotted triangles by dotted edges). The need for this is evident: the compressed view becomes a directed acyclic graph (DAG), which is often significantly (at times even exponentially) smaller than the original tree. Furthermore, one may want to store the view (DAG) in *relations* itself. This raises another question: how

should one define relational views that characterize the compressed XML view? If one is to reduce the XML view update problem to its relational counterpart, this question has to be answered. However, this is nontrivial: the XML view is recursively defined, and a naïve relational encoding may require *infinitely many* relational views.

Third, to locate where the updates take place, one has to evaluate recursive XPath queries on DAGs instead of XML trees. Added to the complication, we need to detect if the update will result in side effects. As observed in [17], it is nontrivial to translate (recursive) XPath queries (resp. updates) over recursive XML views (stored in relations) to SQL queries (resp. updates). To our knowledge, no efficient algorithm has been proposed for evaluating XPath queries with *complex filters* on DAGs stored in relations and none of the previous work considers to detect side effects in the evaluation of XPath.

While these are issues beyond what we have encountered in the relational realm, automated processing of relational view updates is already intricate, even under various restrictions on the views^[1–3]. In fact, even the updatability problem, i.e., the problem of determining whether a relational view is updatable, w.r.t. given updates, is mostly unsolved and few complexity results are known^[1,18]. This tells us that it is unrealistic to reduce the XML view update problem to its relational counterpart and then rely on the DBMSs to do the rest.

Contributions. This paper is the full version of [19] with all proofs and additional algorithms. The paper consists of new techniques for updating *compressed* and possibly *recursively* defined XML views via *schema-directed XML publishing*, in particular ATGs^[8]. (Our techniques are applicable to XML views published from relations via other systems (e.g., SilkRoute, XPERANTO) as long as they represent the XML views in terms of SPJ queries.) Given XML updates of an XML view which is compressed into a DAG and stored in relations, we do the followings: (a) define relational views that characterize the compressed XML view, such that the *number* of relational views is bounded by the size of the XML view, even if the XML view is recursively defined; (b) translate single updates at the XML level to group updates of the relational view representation; (c) translate updates over the relational views to updates over the underlying (published) relational database. More specifically, we make the following contributions.

- *On the XML Side.* (a) We refine the update semantics for XML views of relations to accommodate

XML side effects, based on the semantics of XML views. (b) We develop an algorithm to translate (*recursive*) updates on a (*possibly recursively defined*) XML view to updates on the relational representation of the XML view. (c) To detect XML side effects and translate the updates, we present an efficient algorithm for evaluating XPath queries with *complex filters* on DAGs, based on a new indexing structure to handle recursion and a new technique for handling filters. (d) We also develop efficient algorithms to incrementally maintain the indexing structure.

On the Relational Side. (a) We identify a *key-preservation* condition on SPJ views, which is less restrictive than the conditions imposed by previous work^[1–3]. This condition does not reduce the expressive power of ATGs. (b) We establish complexity results for the updatability problem. We show that under key-preservation on SPJ views, while the problem for tuple insertions is NP-complete, it becomes *tractable* for *group* deletions (which is NP-complete without key preservation). (c) We propose a PTIME algorithm for processing group deletions on SPJ views. (d) To process group insertions we give an efficient heuristic algorithm.

Experimental Study. Our experimental results verify the effectiveness and efficiency of our techniques.

These techniques are the first for processing XML updates with *recursion and complex filters* on *compressed and possibly recursively defined* XML views, without relying on the high-end and mostly unavailable view-update functionality of the underlying relational DBMS. They provide the capability of supporting XML view updates within the immediate reach of most XML publishing systems. On the relational side, our complexity results and algorithms are a useful addition to the study of relational view updates.

Organization. Section 2 defines XML updates and reviews a tool for publishing relational data, namely ATGs. Section 3 develops algorithms for translating XML updates to relational view updates, and Section 4 presents our complexity results and algorithms for handling relational view updates. An experimental study is presented in Section 5, followed by related work in Section 6. We conclude in Section 7.

2 View Updates Revisited in the XML Setting

In this section we define the syntax and semantics of XML updates. We review how XML views may be generated from a relational database and outline our approach to processing the updates over DAG compression of relationally stored XML views.

2.1 XML View Updates: Side Effects and Semantics

Syntax. Following [14,15], we specify XML updates in terms of XPath expressions: (a) insert (A, t) into p , (b) delete p . Here, A is an element type, and t is an instantiation of the semantic attribute $\$A$ of A . Given the instantiation we can uniquely identify the root of a subtree of type A (see Subsection 2.3). We define p as an XPath expression:

$$\begin{aligned} p &::= \epsilon \mid A \mid * \mid // \mid p/p \mid p[q], \\ q &::= p \mid p = \text{“}s\text{”} \mid \text{label}() = A \mid q \wedge q \mid q \vee q \mid \neg q, \end{aligned}$$

where ϵ , A , $*$ and $/$ denote the *self-axis*, a label (tag), a wildcard and the *child-axis*, and $//$ stands for */descendant-or-self::node()/*, respectively; q in $p[q]$ is called a *filter*, in which s is a constant (string value), and \wedge , \vee and \neg denote conjunction, disjunction and negation, respectively. For $//$, we abbreviate $p_1//$ as $p_1//$ and $//p_2$ as $//p_2$.

Side Effects. We next study the side effects of XML view updates. On detecting side effects, users can choose either to abort the update, or to carry on under the semantics we provide. Detection of side effects will be further elaborated in Subsection 3.2.

Recall the update Δ_X from Example 1. The update is to change the subtrees (prerequisite hierarchy) of only those CS320 nodes below CS650. This update will result in side effects since CS320 also appears elsewhere below the root. The subtree property of XML publishing tells us that the subtree of a CS320 node is *uniquely determined* by the value of its semantic attribute $\$course$, which is in turn determined by the set of relational records for *all* CS320 nodes. In other words, changes incurred to the subtree of any CS320 node must also be reflected to *all* CS320 nodes, rather than only to those below CS650.

The side effect issue is more subtle for deletions. Consider `delete course[cno=CS650]/prereq/course[cno=CS320]` on the XML tree of Fig.1. The deletion aims to remove course CS320 from the prerequisites of course CS650. The subtree property instructs that we should remove all CS320 nodes, not only the CS320 node under the CS650 node. On the other hand, this cannot be simply performed by physically removing all CS320 nodes as in previous work on XML view updates^[11–13]: CS320 is itself, an independent CS course and, moreover, may be a prerequisite of other courses. For a correct deletion we first need to find all the *parents* of the nodes to be removed, i.e., those *prereq* nodes below CS650 nodes, and then remove CS320 from the *children* list of only those *parent* nodes.

Semantics of XML View Updates. It is obvious that a new semantics should be developed to cope with *side effects* like the ones mentioned. This semantics needs to respect the hierarchical nature of XML views. Note that this semantics is *different* from the semantics of updates on XML data^[14,15]. Given an XML view T with root r , an insert operation: (a) it finds the set of all *elements* reachable from r via p in T , denoted by $r[[p]]$; (b) for each element v in $r[[p]]$, it adds the new subtree $ST(A, t)$ as the rightmost child of v ; and moreover, (c) for each element u that has the same type and semantic attribute value as v , it adds also $ST(A, t)$ as the rightmost child of u as required by the semantics of XML views.

A delete operation on XML views (a) computes $r[[p]]$; (b) for each node $v \in r[[p]]$, removes the subtree $ST(A, t)$ from the children list of the parent node u of v , where A is the type of v and t is the value of $\$A$ at v ; and (c) for any node u' that has the *same type and semantic attribute value* as the *parent* u of v , removes $ST(A, t)$ from the children list of u' .

Compared to previous work^[11–13], we support XML view updates that (a) are defined with much richer XPath expressions with *recursion and complex filters*, (b) operate on (possibly) recursively defined XML views, and (c) possess a new semantics that captures *side effects*, if any, of XML view updates. We also provide techniques to *detect* whether there are side effects and, in those cases, allow the users to cancel the update; otherwise, the operation will carry on with the semantics described earlier.

2.2 Attribute Translation Grammars (ATGs)

In this subsection, we review XML publishing with Attribute Translation Grammars (ATGs). It should be remarked that the proposed techniques in this paper are applicable to other XML publishing frameworks, e.g., SILKROUTE and XPERANTO.

An ATG can be understood as a mapping $\sigma : \mathcal{R} \rightarrow \mathcal{D}$, where \mathcal{R} is a relational schema and \mathcal{D} is a predefined (possibly recursive) DTD. Given an instance I of \mathcal{R} , σ produces an XML view T , denoted as $\sigma(I) = T$, that conforms to \mathcal{D} . A DTD \mathcal{D} is a triplet (E, P, r) , where E is a finite set of (*element*) *types*; $r \in E$ is called the *root type*; P defines the element types: a *production*, $A \rightarrow \alpha$, is associated with each A in E , where α is an expression of the form:

$$\alpha ::= p\text{cdata} \mid \epsilon \mid B_1, \dots, B_n \mid B_1 + \dots + B_n \mid B^*,$$

where ϵ is the empty word, B is a type in E (a *child* type of A), and $,$, $+$, and $*$ denote concatenation,

alternation, and the Kleene star, respectively.^① A DTD is *recursive* if a type is defined (directly or indirectly) in terms of itself. It is shown in [16] that DTDs found in practice are often recursive.

Example 2. Consider a relational schema R_0 shown in Example 1. We define an ATG σ_0 , shown in Fig.2, for publishing the registrar's database into an XML view. The bold text highlights the DTD D_0 embedded in σ_0 . Note that D_0 is recursive as **course** is indirectly defined with prerequisite **courses**. The XML view generated shall conform to D_0 . A possible XML view T_0 is illustrated with Fig.1. Consider the course CS320 subtree. It appears at different places in the XML view. It is more efficient to keep a single copy of the CS320 subtree and use references (the dotted arrows) to represent multiple occurrences of the subtree. We shall discuss the details of this representation shortly.

```

db → course*
  $course λ Qdb_course
  Qdb_course: select c.cno, c.title
             from course c
             where c.dept = 'CS'
course → cno, title, prereq, takenBy
  $cno = $course.cno, $title = $course.title,
  $prereq = $course.cno, $takenBy = $course.cno
prereq → course*
  $course λ Qprereq_course($prereq)
  Qprereq_course(c1): select c.cno, c.title
                    from prereq p, course c
                    where p.cno1=c1 and p.cno2=c.cno
takenBy → student*
  $student λ QtakenBy_student($takenBy)
  QtakenBy_student(c): select s.ssn, s.name
                    from enroll e, student s
                    where e.cno=c and e.ssn=s.ssn

```

Fig.2. Example ATG σ_0 .

Given an instance I of \mathcal{R} , the ATG σ systematically extracts portions of I into an XML view as follows. (a) For each element type A in \mathcal{D} , σ defines a semantic attribute, or simply called *tuple*, $\$A$, with fixed arity and type; intuitively, $\$A$ governs the generation of an A -subtree, and is passed to the production of A 's children as the view is generated. (b) For each production $p = A \rightarrow \alpha$ in \mathcal{D} and each type B in α , σ specifies an SPJ query, $rule(p)$, which extracts data from a relational database; using the data and $\$A$, it generates the B children of an A node and the tuple for $\$B$. For example, for the production $prereq \rightarrow course^*$, the SPJ

query can be specified as $Q_{prereq_course}(\$prereq)$. In all, σ generates the XML view top-down with reference to \mathcal{D} .

Example 3. Consider a **prereq** node v with the tuple $\$prereq$. $\$prereq$ is used as a constant in the query Q_{prereq_course} to extract data from the source database I . For each tuple t returned by $Q_{prereq_course}(\$prereq)$, a **course** child node v_c is generated and t is associated with v_c . Then the production of **course** is invoked with v_c and t in a similar manner.

2.3 Relational Coding of Recursively Defined XML Views

Consider an ATG $\sigma : \mathcal{R} \rightarrow D$ that defines XML views of relational databases \mathcal{R} . To reduce the update problem for XML views defined by σ to its relational counterpart, we define relational views \mathcal{V}_σ to characterize σ . This is nontrivial: (a) σ is possibly recursively defined; on such views the encoding methods of previous work (e.g., [11]) may lead to *infinitely* many relational views; (b) we consider DAG compressions of XML views, i.e., a DAG representation of $\sigma(I)$ where I is an instance of \mathcal{R} as opposed to trees assumed in previous work. To this end we define \mathcal{V}_σ by means of the edge relations in $\sigma(I)$ as follows.

(a) We assume a compact, unique value associated with each tuple value of semantic attribute $\$A$ in $\sigma(I)$. We abstract away the implementation of this identity value by assuming, w.l.o.g., the existence of a Skolem function gen_id that, given the tuple value of $\$A$, computes id_A that is unique among all identities associated with *all* semantic attributes. We use gen_A to denote the set of the identities of all $\$A$ tuples, which is computed once.

(b) We encode an XML view definition σ in terms of \mathcal{V}_σ as a set of SPJ queries $Q_{edge_A_B}$ materializing the edge relations of σ . More specifically, for each production $A \rightarrow P(A)$ in the DTD of σ , and for each child type B in $P(A)$, we create a relation $edge_A_B$ with two columns, id_A and id_B . Consider productions of the form $A \rightarrow B^*$, where $\$B \leftarrow Q(\$A)$ is the associated query in σ . Then $edge_A_B$ is the set of pairs (ia, ib) such that $ia = gen_id(a)$, $ib = gen_id(b)$, where $a \in gen_id(a)$, $b \in Q(a)$. The definition of $Q_{edge_A_B}$ is similar for productions of other forms. One example of an edge-relation query derived from the σ_0 ATG of Fig.2 is $Q_{edge_prereq_course}$:

```

select gen_id(gp), gen_id(c.cno, c.title)
from gen_prereq gp, prereq p, course c

```

^①An arbitrary DTD can be normalized into a DTD in the form defined by introducing additional element types in linear time. A post-publishing processing then transforms the XML view into one that conforms to the DTD in $O(|T|)^{[20]}$.

where $p.cno1 = gp.cno$ and $p.cno2 = c.cno$

Observe the following about \mathcal{V}_σ . 1) \mathcal{V}_σ encodes the DAG *compression* of XML view $\sigma(I)$. Indeed, for any subtree $ST(A, \$A)$ in $\sigma(I)$, each edge (ia, ib) in $ST(A, \$A)$ is stored *only once* in a relation $edge_{A_B}$ no matter how many times $ST(A, \$A)$ (and thus the edge) appears in $\sigma(I)$. 2) Each $Q_{edge_{A_B}}$ in \mathcal{V}_σ is defined by an SPJ query. Thus \mathcal{V}_σ consists of only SPJ *views*. 3) \mathcal{V}_σ consists of a *bounded* number of *relational views* even if σ is *recursively* defined.

It should be remarked that there have been a few alternative encoding schemes for XML views (possibly with compression). For example, inlining techniques^[21] were proposed to encode recursive XML in a finite number of relations. For presentation brevity, we propose edge-based relations from an ATG but skip the details of the applications of other particular XML encoding schemes. Furthermore, as we shall see soon, our DAG compression supports efficient view updates in XML settings, e.g., side-effect detection and XPath evaluation. Any alternative encodings employed for XML view updates must address these issues.

Updates on Relational Views. Given an update Δ_X on a DAG compressed XML view $\sigma(I)$, we convert it to updates Δ_V on the relational view $V = \mathcal{V}_\sigma(I)$. The relational view updates Δ_V consist of edge tuples of the form $t = (ia, ib)$ to be inserted into or deleted from an edge relation $edge_{A_B}$.

Note that a shared tree cannot be simply removed. Consider again the deletion of Subsection 2.1 on the XML view of Fig.1. We cannot remove the subtree of CS320 completely even if all CS320 nodes are in the *prereq* subtree of some CS650 nodes. This is because some subtrees inside CS320 (i.e., certain *students*) may be shared and referenced by other nodes outside of the subtree.

In response to this, we compute the relational view updates Δ_V such that (a) a newly inserted subtree is only stored once in V no matter how many times it appears in the updated view, and (b) a deleted subtree is not physically removed: only the tuple (ia, ib) in V representing the corresponding parent-child edge is deleted

from its edge relation $edge_{A_B}$. More specifically, the tuple corresponding to ia is not removed from gen_A because ia is a parent node in $r[[p]]$ and needs to be kept in the XML view. To cope with subtree sharing, ib is not removed from gen_B when the edge (ia, ib) is removed from $edge_{A_B}$; instead, upon the completion of processing Δ_V , our incremental maintenance algorithm runs in the *background* to remove tuples from gen_B 's that are no longer linked to any node; it is at the completion of Δ_V when gen_B 's are updated (similarly for insertions). Note that gen_B 's are not defined as a view; they are derived from V (i.e., the edge relations \mathcal{V}_σ) and maintained in the background.

2.4 Processing XML View Updates

We propose a framework for processing XML view updates, as shown in Fig.3. For each ATG (XML view definition) $\sigma : \mathcal{R} \rightarrow D$, we maintain a relational database I of \mathcal{R} , and the relational views V that encode the DAG compression of $T = \sigma(I)$. The users pose updates on (the virtual view) T . Given a single XML update Δ_X on T as input, we are to generate a group update Δ_R on I such that $\Delta_X(T) = \sigma(\Delta_R(I))$ if such Δ_R exists; otherwise *reject* Δ_X as early as possible. Specifically, the framework processes an XML update Δ_X on T in three phases, namely, DTD *validation*, *translation from Δ_X to Δ_V* (Section 3), and *translation from Δ_V to Δ_R* (Section 4). If our algorithm detects a side effect, we report it to the user. After the relational update Δ_R is computed, we update the underlying database I using Δ_R , update the relational views V using Δ_V , and finally, *in the background*, invoke our incremental algorithm to maintain the indexing structures and to remove from gen_A those node ids that are no longer reachable from the root of the XML view T .

Before we end this section, we discuss DTD validation. The other steps in processing the XML view updates will be discussed in subsequent sections. Given XML updates Δ_X , we first perform static optimization by validating the predefined DTD D with respect to

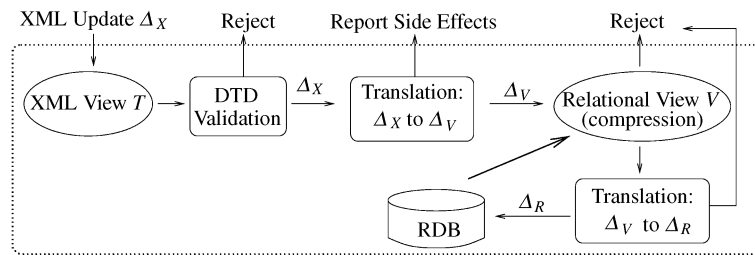


Fig.3. Overview of XML view updates.

Δ_X , and reject the updates if $\Delta_X(T)$ does not conform to D as required by the schema-directed definition of σ .

The validation is conducted at the schema level by leveraging the DTD normalization given in Subsection 2.2, as follows. Let Δ_X be defined in terms of an XPath query p . We first “evaluate” p on the DTD D to find the types of the elements reached via p . We then check whether the insertion or deletion of subtrees of these elements (types) violates their productions in the DTD D . Note that an insertion (resp. deletion) of a B child under an A element does not violate D only if the production of A is of the form $A \rightarrow B^*$. Thus updates of other forms can be immediately rejected. This can be checked in $O(|p| |D|^2)$ time, where $|p|$ and $|D|$ are the sizes of the XPath query p and the DTD D , respectively.

3 Mapping XML View Updates to Relations

In this section we present a technique for translating XML updates on an XML view to updates on relational views representing the DAG compression of the XML view, derived from the ATG Subsection 2.3. The technique consists of four parts: (a) indexing structures for checking ancestor-descendant relationships in a DAG (Subsection 3.1), (b) an efficient algorithm for evaluating XPath queries on DAGs and detecting side effects (Subsection 3.2), (c) algorithms to translate updates on the XML view to updates on its relational representation (Subsection 3.3), based on the indexing structures and the evaluation algorithm, and (d) incremental algorithms for maintaining the indexing structures (Subsection 3.4).

3.1 Auxiliary Structures

To efficiently process “//” and filters on a DAG, we introduce two auxiliary structures: a *topological order* and a *reachability matrix*. The reachability matrix can efficiently support “//” in XPath queries on a DAG while the topological order is crucial in evaluating XPath filters as well as in computing and maintaining the reachability matrix.

Topological Order. Recall from Section 2 the function *gen_id()*, which generates a unique id for each node based on the value of its semantic attribute. Given a representation of a DAG V , we create a list \mathbf{L} consisting of all the distinct node identities in V topologically sorted such that u precedes v in \mathbf{L} only if u is not an ancestor of v in the DAG, i.e., there is no path from u to v in the DAG.

The list \mathbf{L} can be computed in $O(|V|)$ time (see, e.g., [22]), where $|V|$ is the size of the relational views. Its size, $|\mathbf{L}|$, is the number of *distinct nodes* in the DAG, denoted by n . Note that \mathbf{L} is computed once when V is created and it is maintained incrementally.

Reachability Matrix. To evaluate the ancestor-descendant relationship between a pair of nodes in a DAG, we use an $n \times n$ *reachability matrix* \mathcal{M} : a cell in \mathcal{M} is a bit. Given a row i denoting node n_i and a column j indicating node n_j , if cell \mathcal{M}_{ij} is set, n_i is an ancestor of n_j in the XML view (resp. n_j is a descendant of n_i).

To store \mathcal{M} , we conceptually need as many bits as n^2 . The cost for that is prohibitive. To overcome this, we store only information about the set bits of the reachability matrix. That is, \mathcal{M} is physically stored as a relation $\mathbf{M}(anc, desc)$, where *anc* denotes an ancestor node, and *desc* a descendant. We use $desc(a)$ (resp. $anc(a)$) to denote the descendants (resp. ancestors) of node a retrieved from \mathbf{M} .

Input: the relational view V and topological order \mathbf{L} .
Output: reachability matrix \mathbf{M} .

1. $\mathbf{M} := \emptyset$;
2. **for** ($k := |\mathbf{L}|$; $k > 0$; $k - -$) /* backward topological order */
3. $d := \mathbf{L}[k]$;
4. $A_d := \{a_2 \mid a_2 \in anc(a_1), a_1 \in parent(d)\}$;
5. **insert** (a, d) **into** \mathbf{M} **for each** $a \in A_d$;
6. **return** \mathbf{M}

Fig.4. Algorithm Reach.

Relation \mathbf{M} can be computed in $O(|V|^2 \log |V|)$ time from V (see, e.g., [22]). Capitalizing on the topological order \mathbf{L} we give Algorithm Reach, shown in Fig.4, that computes \mathbf{M} in $O(n |V|)$ time. It is based on dynamic programming: it ensures that for a node d the ancestors of the nodes in the set of parents of d , denoted by $parent(d)$, are already known before we compute ancestors A_d , such that we can compute A_d by using those previously computed ancestors (lines 4~5). This can be achieved by processing the nodes in the order of \mathbf{L} from right to left (line 2). Note that $parent(d)$ can be computed from the edge relations in V .

To see that Algorithm Reach runs in $O(n |V|)$ time, observe the followings: (a) for each node in \mathbf{L} we visit its parents once and thus any node v is visited $in(v)$ times, where $in(v)$ is the in-degree of v , i.e., the number of incoming edges to v in the DAG; (b) the sum of $in(v)$'s for all v is $|V|$; and (c) each visit takes at most $O(n)$ time. In practice, $|\mathbf{M}| \ll n^2 \ll |V|^2$, where V is typically much smaller than the XML tree T , even up to an exponential factor.

3.2 Evaluating XPath Queries on DAGs

To translate updates Δ_X on XML views to updates Δ_R on relational views and detect whether the update will yield side effects, we have to evaluate the XPath expression embedded in Δ_X . The DAG compression of XML views introduces new challenges: previous work on XPath evaluation has mostly focused on trees rather than DAGs. While the evaluation algorithms were developed for path queries on DAGs^[23,24], they cannot be applied in our setting because (a) they do not deal with complex filters which, as will be seen shortly, require a separate pass of the input DAG, and (b) they do not address maintenance of the indexing structures they employ, which is necessary when the DAG is updated. Path-query evaluation algorithms were also developed for semi-structured data (general graphs). However, these algorithms neither treat DAGs differently from cyclic graphs (and thus may not be efficient when dealing with DAGs), nor consider XPath queries used in XML view updates.

To this end we outline an efficient algorithm for evaluating an XPath query p on an XML tree T that is (a) compressed as a DAG, and (b) stored in edge relations V . The algorithm takes as input an XPath query p over T , the relational views V , and the reachability matrix M . It computes (a) a set $r[[p]]$ consisting of, for each node reached by p , a pair (B, v) , where v is the id and B the type of the node, respectively; (b) a set $E_p(r)$ consisting of, for each v reached by p , tuples of the form $((C, u), v)$, where u is the id of a parent of v in the DAG (i.e., there is an edge from u to v) such that p reaches v through u , and C is the type of u ; we shall see that the set $E_p(r)$ is needed for handling deletions; and (c) the set of nodes S in T which are affected by the update but are not reachable via p . If the set S is not empty, the update will generate XML side effects. Note that for each v there are possibly multiple (C, u) pairs, since we are dealing with a DAG (in which a node may have multiple parents) rather than a tree.

For XML data stored as a tree T , [25] developed an algorithm that evaluates an XPath query p in two passes (linear scans) of T . The basic idea of [25] is to first convert T to a binary-tree representation (before the two-pass process is invoked), and then run a bottom-up tree automaton on the binary tree to evaluate filters, followed by a run of a top-down tree automaton to identify nodes reached by p . It has linear-time complexity, the “optimal” one can expect^[25]. We next show that a *comparable complexity* can be achieved when evaluating XPath queries on a DAG stored in relations.

Our evaluation algorithm uses the following vari-

ables. (a) A list Q of filters including all the sub-expressions of filters in p , topologically sorted such that for any q_i, q_j in Q , q_i precedes q_j if q_i is a sub-expression of q_j . (b) For each q in Q and each node v in L , two Boolean variables $\text{val}(q, v)$ and $\text{desc}(q, v)$ to denote whether or not the filter q holds at v and at any descendant u of v , respectively.

Using these variables, we present a two-pass algorithm to evaluate p on V : a bottom-up phase that evaluates *filters* in p and computes the Boolean variables associated with each node v in L , followed by a top-down phase that computes $r[[p]]$ and $E_p(r)$ using the filters computed. We next outline the algorithm below.

Bottom-Up. The key idea is based on dynamic programming. For each node v in the topological order L , and for each sub-filter q in the topological order Q , we compute the values of $\text{val}(q, v)$ and $\text{desc}(q, v)$. This can be done by structural induction on the form of q . For example, when q is $\text{label}() = A$, $\text{val}(q, v)$ is true if and only if v is in gen_A . When q is $q_1 \vee q_2$, $\text{val}(q, v) := \text{val}(q_1, v) \vee \text{val}(q_2, v)$. When q is a path expression p , p can be rewritten into a “normal form” $\eta_1 / \dots / \eta_n$, where each η_i is either (a) $\epsilon[q_i]$, (b) a label A , (c) wildcard “*”, or (d) “//”. The normal form can be obtained in $O(|p|)$ time by capitalizing on the following rewriting rules: $p[q] \equiv p/\epsilon[q]$, and $\epsilon[q_1] \dots [q_n] \equiv \epsilon[q_1 \wedge \dots \wedge q_n]$. For example, if q is rewritten as $//\eta_2 / \dots / \eta_n$ with $\eta_1 = //$, $\text{val}(q, v)$ is true if either $\text{val}(\eta_2 / \dots / \eta_n, v)$ or $\text{desc}(\eta_2 / \dots / \eta_n, u)$ is true for some child u of v ; correspondingly, $\text{desc}(q, v)$ is true if either $\text{val}(q, v)$ or $\text{desc}(q, u)$ holds. Note that the children of v can be efficiently identified by using the indexes on V . In addition, the algorithm proceeds in the topological orders L and Q . Therefore, the truth values of $\text{val}(\eta_2 / \dots / \eta_n, v)$ and $\text{desc}(\eta_2 / \dots / \eta_n, u)$ are already available before assigning a value for $\text{val}(q, v)$ and $\text{desc}(q, v)$. Similarly $\text{val}(q, v)$ can be computed for all other possible rewrites of q .

Top-Down. Upon completion of the bottom-up phase, we compute $r[[p]]$, $E_p(r)$ and S as follows. As mentioned earlier p can be normalized in the form of $\eta_1 / \dots / \eta_n$, in which all the filters have already been evaluated to be a truth value at each node. Starting from the root r , we find nodes C_i reached after each step η_i and meanwhile maintain a set of nodes S in T that are not reachable via p but will be affected by the update. When η_i is “/” (resp. “//”), S is extended with the parent (resp. ancestor) nodes of C_i that are not reached via p . These nodes can be easily found by using indexes on the edge relations V when η_i is A or “*”, and by means of the reachability matrix M when η_i is “//”. The nodes reached by the last step η_n are put

in $r[p]$, along with their types. The parents through which they are reached via p are put in $E_p(r)$ along with their types during the traversal. One can verify that there is a side effect iff S is not empty. As remarked in Section 2, users may either abort the update or carry out the update using our update semantics.

Example 4. Consider the XML update $\Delta_{X_1} = \text{delete } //course [cno=CS320]//student[sid=S02]$ on the XML tree in Fig.1, which is to delete student S02 from the subtree of course CS320's subtree. Consider course₁ and desc(cno=CS320, course₁). The recurrence relation tells us that desc(cno=CS320, course₁) would be true if desc(cno=CS320, cno₁) is true. In a bottom-up pass, val(cno=CS320, cno₁) and desc(cno=CS320, cno₁) have been evaluated to be true before desc(cno=CS320, course₁). Similarly, the two variables val and desc of all nodes in the DAG can be computed in dynamic programming fashion in one bottom-up pass.

In a top-down pass, we can efficiently evaluate $//course$ by using the index. course₁ is selected since val(cno=CS320, cno₁) = true and course₁ is reachable from the root. course₂ is not selected because val(cno=CS320, cno₂) is false. Similarly, we select student₂ for $//student$. We compute $E_p(r)$ by checking the parent of student₂ that is on a path satisfying $//course[cno=CS320]//student[sid=S02]$. We obtain ((takenBy, takenBy₁), student₂).

Complexity. In the bottom-up phase, each node v is visited at most $\text{in}(v)$ times, where $\text{in}(v)$ is the in-degree of v . In the top-down phase, each node is visited only once except the final step when a node u may be included in $E_p(r)$ at most $\text{out}(u)$ times, where $\text{out}(u)$ is the out-degree of u . Putting these together, the complexity of the algorithm is $O(|p| |V|)$.

Compared to the algorithm of [25], observe the followings. (a) While our algorithm operates on the DAG representation, it visits the nodes of the corresponding (uncompressed) tree at most twice, i.e., it has the same complexity as that of [25]. When dealing with DAGs that do not have a tree structure, it is necessary to visit all the edges in the DAGs in the worst case and thus our algorithm is asymptotically optimal. (b) In contrast to [25], our algorithm does not require the conversion of the input data into binary trees and the construction of tree automata, which are potentially very large. (c) Our algorithm works on DAGs (including trees) while [25] cannot work on DAGs.

3.3 Translating Updates from XML to Relations

On account of the relational representation (DAG) of XML views, a single XML update may be mapped

to multiple relational updates (a group update) over the edge tables V (see Subsection 2.3). We next give two algorithms, Xinsert and Xdelete, for translating XML view insertions and deletions to relational view updates Δ_V , respectively.

Insertion. Algorithm Xinsert is presented in Fig.5. Given $\Delta_X = \text{insert}(A, t)$ into p on the XML view T , the objective is to return the group of insertions Δ_V over V (which will then be tested for acceptance). The first step is to find the set of edges in the newly inserted subtree $ST(A, t)$ with the root r_A , which is computed by the algorithm of [8] and the function $gen_id()$ (lines 2~3). We then generate the relational view updates: for each edge (ui, vi) in the newly inserted subtree, we add (ui, vi) to Δ_V (lines 4~5); moreover, for each $(B, ui) \in r[p]$, we add (ui, r_A) as a new edge to Δ_V (lines 6~7). The set $r[p]$ of nodes (pairs (B, ui) of node ids along with their types) reached by XPath p from the root of T (line 6) is computed using the evaluation algorithm of Subsection 3.2.

<p>Input: an insertion of the form $\Delta_X = \text{insert}(A, t)$ into p over T, and the relational view V.</p> <p>Output: a group insertion Δ_V over V.</p> <ol style="list-style-type: none"> 1. $\Delta_V := \emptyset$; 2. $E_A := \{(B, gen_id(\\$u)), (C, gen_id(\\$v)) \mid (u, v) \text{ is an edge in } ST(A, t), u, v \text{ with type } B, C \text{ resp.}\}$; 3. $r_A :=$ the id of $ST(A, t)$'s root as generated by $gen_id(t)$; 4. for each $((B, ui), (C, vi)) \in E_A$ 5. $\Delta_V := \Delta_V \cup \{\text{insert}(ui, vi) \text{ into } edge_B_C\}$; 6. for each $(B, ui) \in r[p]$ 7. $\Delta_V := \Delta_V \cup \{\text{insert}(ui, r_A) \text{ into } edge_B_A\}$; 8. return Δ_V;
--

Fig.5. Algorithm Xinsert.

<p>Input: a deletion $\Delta_X = \text{delete } p$ over T and the rel. view V.</p> <p>Output: a group deletion Δ_V over V.</p> <ol style="list-style-type: none"> 1. $\Delta_V := \emptyset$; 2. for each $((C, ui), vi) \in E_p(r)$, where $(B, vi) \in r[p]$ 3. $\Delta_V := \Delta_V \cup \{\text{delete}(ui, vi) \text{ from } edge_C_B\}$; 4. return Δ_V;
--

Fig.6. Algorithm Xdelete.

Deletion. Algorithm Xdelete is shown in Fig.6. Given $\Delta_X = \text{delete } p$, Algorithm Xdelete returns the group of relation view deletions Δ_V over V , which will be passed to subsequent steps to test for acceptance (Subsection 4.2). For each node vi in $r[p]$ and each parent ui of vi in $E_p(r)$, Xdelete removes the edge (ui, vi) from V (lines 2~3). The parent-child relation is computed by using the set $E_p(r)$, whose computation is coupled with that of $r[p]$ (See Subsection 3.2).

Observe that these algorithms implement *the new semantics* of XML view updates given in Section 2. This is achieved by leveraging the characterization of the XML view T in terms of relational views V . Indeed, for two edges (u, v) , (u', v) in T , if two parents u and u' of the same node v have the same element type A and the same value of the semantic attribute $\$A$, the two edges are represented by a *single* tuple in some edge relation $edge_A.B$. Thus there is no need to search V to find different nodes sharing (A, t) , i.e., XML side effects described in Section 2 do not incur extra cost. Furthermore, the set semantics of V ensures that a newly inserted subtree is stored *only once*. In addition, Algorithm Xdelete does *not* physically remove a deleted subtree; instead, only the corresponding parent-child edge is removed. These naturally comply with the requirements of DAG update semantics given in Section 2.

Example 5. Reconsider the XML update $\Delta_{X_1} = \text{delete } //course [cno=CS320]//student[sid=S02]$ on the XML tree in Fig.1. Given this as input, Algorithm Xdelete yields $\Delta_{V_1} = \{(takenBy_1, student_2)\}$. As another example, given $\Delta_{X_2} = \text{delete } //student[sid=S02]$, we get $\Delta_{V_2} = \{(takenBy_1, student_2), (takenBy_2, student_2)\}$.

Complexity. Algorithm Xinsert takes $O(|E_A| + |r[p]|)$ time at most, which is the cost of inserting the “inner” connections of $ST(A, t)$ into V and connecting $ST(A, t)$ to the rest of V , where $|E_A|$ is the number of edges in $ST(A, t)$. Algorithm Xdelete takes $O(|E_p(r)|)$ time. Together with the complexity $O(|p| |V|)$ of evaluating p , this is the cost of generating Δ_V from Δ_X .

3.4 Maintaining Auxiliary Structures

We next outline how to maintain the reachability matrix M and the topological order L in response to updates over V . We should remark that the maintenance of M and L is performed in the *background* in parallel with the processing of relational updates Δ_R ; as a result, in our framework (Fig.3), maintenance does not slow down the process of carrying out XML view updates. The maintenance can be cumbersome, as illustrated by the next example.

Example 6. Recall the XML update Δ_{X_1} from Example 5. This entails that all reachability information to S02 be deleted from the root of the CS320 subtree and from *all nodes* on the path to S02. Moreover, this course may be a prerequisite of other courses, e.g., CS650; since CS320’s subtree is shared, the reachability information from CS650 to S02 should be updated.

Recomputing M from the updated V bears a prohibitive cost. What we ideally would like is

to *incrementally* update M . Existing incremental techniques^[26,27] for updating reachability information are not applicable since they rely on special auxiliary structures which are themselves expensive to construct and maintain (e.g., [26] requires the computation of a spanning tree, taking $O(n |V|)$ time for each node insertion). On the other hand, incremental algorithms of updating topologically ordered lists (e.g., [28]) takes $O(|V|)$ time per edge insertion. Given these high individual complexities we follow a hybrid approach by maintaining both auxiliary structures at once.

We next give two algorithms, $\Delta_{(M, L)}\text{insert}$ and $\Delta_{(M, L)}\text{delete}$, for maintaining auxiliary structures M and L in response to XML view insertions and deletions, respectively.

Insertion. Algorithm $\Delta_{(M, L)}\text{insert}$ is shown in Fig.7. Given $\Delta_X = \text{insert}(A, t)$ into p , it finds the Δ_M over M to maintain the reachability information, and moreover, updates the topological order L in response to the insertion of $ST(A, t)$.

<p>Input: an insertion of the form $\Delta_X = \text{insert}(A, t)$ into p over T, the rel. view V, reachability matrix M and topological order L.</p> <p>Output: insertions Δ_M over M, and updated list L.</p> <ol style="list-style-type: none"> 1. compute N_A and r_A, as lines 2~4 in Algorithm Xinsert; 2. $L_A :=$ the topological order of nodes in $ST(A, t)$; 3. $\Delta_M :=$ reachability matrix for $ST(A, t)$; /*using Algorithm Reach*/ 4. for each $a \in \text{anc}(r[p])$ and each $d \in N_A$ /* computing Δ_M */ 5. $\Delta_M := \Delta_M \cup \{\text{insert}(a, d) \text{ into } M\}$; 6. $N_C :=$ the set of common nodes in lists L and L_A; /*update L^**/ 7. $L_{N_C} :=$ the topological order of nodes in N_C; 8. for $(k = L_{N_C} ; k > 1; k - -)$ /*align L_A and L with L_{N_C}*/ 9. $u := L_{N_C}[k]; v := L_{N_C}[k - 1]$; 10. if $\text{ord}_{L_A}(u) < \text{ord}_{L_A}(v)$ then $\text{swap}(L_A, u, v)$; 11. if $\text{ord}_L(u) < \text{ord}_L(v)$ then $\text{swap}(L, u, v)$; 12. if $r_A \in L$ then for each u in $r[p]$ 13. if $\text{ord}_L(u) < \text{ord}_L(r_A)$ then $\text{swap}(L, u, r_A)$; 14. $L :=$ merge L_A into L; 15. return (Δ_M, L);
--

Fig.7. Maintenance algorithm $\Delta_{(M, L)}\text{insert}$ for insertions.

It is simple to compute Δ_M , which consists of two parts: (a) the reachability matrix for the newly inserted DAG $ST(A, t)$ is computed by invoking Algorithm Reach (line 3); (b) for each $a \in \text{anc}(r[p])$ (ancestors of nodes in $r[p]$) and each $d \in ST(A, t)$, we add (a, d) to Δ_M (lines 4~5).

Maintaining L is a bit cumbersome. As will be shown, M is useful in maintaining L . Before con-

sidering inserting a DAG ($ST(A, t)$), we first consider how to maintain \mathbf{L} when one edge is inserted. For an edge insertion (u, v) , if v is already in front of u in \mathbf{L} , \mathbf{L} remains valid without any change; otherwise, special care is needed to update node positions in \mathbf{L} . We illustrate this by an example. Consider part of \mathbf{L} : $\langle \dots, d_u, u, a_{u_1}, a_1, d_{v_1}, a_{u_2}, v, \dots \rangle$, where a_{u_1} and a_{u_2} are ancestors of u , d_{v_1} is a descendant of v , d_u is a descendant of u , and a_1 is neither an ancestor of u nor a descendant of v . After (u, v) is inserted, we can obtain a correct topological order by moving v and its descendants (d_{v_1}) between u and v such that they precede u . This yields $\langle \dots, d_u, d_{v_1}, v, u, a_{u_1}, a_1, a_{u_2}, \dots \rangle$. Note that d_{v_1} must be neither an ancestor of u (otherwise there is a cycle) nor an ancestor of a_1 . To formalize this, we denote the nodes between u and v in \mathbf{L} as $\mathbf{L}[u : v]$. Given an edge insertion (u, v) , the correct topological order can be obtained by moving the nodes in $\mathbf{L}[u : v] \cap \text{desc}(v)$ immediately in front of u in \mathbf{L} . The procedure of changing \mathbf{L} for reflecting the insertion (u, v) is denoted as $\text{swap}(\mathbf{L}, u, v)$, where u precedes v in \mathbf{L} before the move.

We next explain the algorithm for updating \mathbf{L} when inserting $ST(A, t)$ (lines 6~14). Let \mathbf{L}_A be the topological order for $ST(A, t)$ (line 2) and N_C be the set of common nodes in \mathbf{L} and \mathbf{L}_A . The basic idea of the algorithm is to make the relative orders of nodes in N_C consistent in lists \mathbf{L} and \mathbf{L}_A before we merge \mathbf{L} and \mathbf{L}_A to obtain the updated \mathbf{L} . To do this, we compute the topological orders \mathbf{L}_{N_C} for nodes in N_C by considering the edges that connect nodes of N_C in either T or $ST(A, t)$ (line 7), and then align \mathbf{L} and \mathbf{L}_A with \mathbf{L}_{N_C} to make their positions consistent with \mathbf{L}_{N_C} (lines 8~11). One subtlety is worth mentioning: when performing the alignment we follow the order of \mathbf{L}_{N_C} from the right to the left. This processing order ensures that the position of aligned nodes will not be changed by subsequent alignment. To be specific, the aligned nodes are not descendants of nodes to be aligned and thus will not be moved any more when $\text{swap}(\mathbf{L}, u, v)$ is called in subsequent alignment (they are not descendants of v). Furthermore, if the root of $ST(A, t)$ is already in T , we may need to change the order of \mathbf{L} in response to the inserted edge (u, r_A) , where $u \in r[[p]]$ ($u \notin L_A$) (lines 12~13). After we obtain two consistent lists \mathbf{L} and \mathbf{L}_A , we can merge \mathbf{L}_A into \mathbf{L} to generate the updated \mathbf{L} (line 14). This can be done by regarding the nodes in N_C as “pivots” and inserting the new nodes (i.e., $\mathbf{L}_A \setminus N_C$) into \mathbf{L} before their respective “pivots”.

Deletion. Maintenance of auxiliary structures in response to XML view deletions takes place in the form of Algorithm $\Delta_{(M, L)}\text{delete}$, shown in Fig.8. The al-

gorithm efficiently produces the followings by scanning the elements of an XML deletion Δ_X : (a) deletions Δ_M over M , (b) an updated \mathbf{L} , and (c) as a bonus, the set of edges Δ'_V in the deleted subtree that are no longer connected to any nodes in the DAG and are to be passed to the garbage collector for *background* processing (see Section 2). The set Δ'_V is a direct consequence of deletions Δ_V computed by Algorithm Xdelete. The need arises when a node $d \in \Delta_V$ is to be completely removed from the subtree. This happens when either all its incoming edges are in $E_p(r)$ (described in Subsection 3.2), or all its parent nodes are deleted.

<p>Input: a deletion of the form $\Delta_X = \text{delete } p$ over T, the rel. view V, reachability matrix M and topological order \mathbf{L}.</p> <p>Output: deletions Δ'_V over V, Δ_M over M, and updated list \mathbf{L}.</p> <ol style="list-style-type: none"> 1. $\Delta'_V := \emptyset$; $\Delta_M := \emptyset$; 2. $\mathbf{L}_R :=$ the sorted list $\text{desc}(r[[p]])$ according to topological order \mathbf{L}; 3. $\text{keep}(d) := \text{true}$ for each $d \in T$; /*initialize state*/ 4. for each d in a backward traversal of \mathbf{L}_R 5. $P_d := \emptyset$; 6. for each $a \in \text{parent}(d)$ 7. if $((C, a), d) \notin E_p(r)$ and $\text{keep}(a) = \text{true}$ 8. then $P_d := P_d \cup \{a\}$; 9. $A_d := \{a_2 \mid a_2 \in \text{anc}(a_1), a_1 \in P_d\}$; 10. for each $a \in \text{anc}(d) \setminus A_d$ 11. $\Delta_M := \Delta_M \cup \{\text{delete } (a, d) \text{ from } M\}$; 12. if $P_d = \emptyset$ /*compute Δ'_V and update \mathbf{L}^*/ 13. then $\text{keep}(d) := \text{false}$; 14. delete d from list \mathbf{L}; 15. for any child d' (of type H) of d (of type G) 16. $\Delta'_V := \Delta'_V \cup \{\text{delete } (d, d') \text{ from } \text{edge_G-H}\}$; 17. return $(\Delta'_V, \Delta_M, \mathbf{L})$

Fig.8. Maintenance algorithm $\Delta_{(M, L)}\text{delete}$ for deletions.

The algorithm progresses by populating deletions Δ_M while, at the same time and whenever applicable, removing elements from \mathbf{L} and populating Δ'_V . The first step is arranging all nodes in all deleted subtrees in a list \mathbf{L}_R (line 2). To do so, we compute $\text{desc}(r[[p]])$, i.e., the descendants of all nodes in $r[[p]]$; we then sort \mathbf{L}_R according to \mathbf{L} ; this is always possible since $\mathbf{L}_R \subseteq \mathbf{L}$. For each node d in T we associate a state $\text{keep}(d)$, initialized to true , and keeping track of whether or not the node should be deleted in the end (line 3). \mathbf{L}_R is then traversed backwards (line 4); this processing order of \mathbf{L}_R ensures that each d in \mathbf{L}_R is processed after its ancestors. Thus it guarantees correct deletion semantics. For each d in \mathbf{L}_R we compute its undeleted parents (lines 6~8) P_d (i.e., any node a in its parent set for which $\text{keep}(a)$ is true) and then its *new* ancestors A_d (line 9). If there is a node in d 's current ancestors

$\text{anc}(d)$ that is not in A_d , it should be removed from \mathbf{M} (lines 10,11). If d does not have any parents (i.e., $P_d = \emptyset$) we set its keep state to false and delete it from \mathbf{L} (lines 13,14). Observe that according to the semantics of \mathbf{L} , an element removal does not affect the topological order of the rest of its elements. In addition, all outgoing edges from a deleted node d are deleted from V (lines 15,16); children d' of d can be readily identified from d 's type.

Example 7. Recall Δ_{X_1} from Example 5. Given Δ_{X_1} , Algorithm $\Delta_{\mathbf{M},\mathbf{L}}\text{delete}$ returns 1) $\Delta'_{V_1} = \emptyset$, 2) unchanged \mathbf{L} , and 3) $\Delta_{M_1} = \{(\text{prereq}_2, \text{student}_2), (\text{prereq}_2, \text{sid}_2), (\text{prereq}_2, \text{name}_2), \dots\}$, i.e., the reachability information from nodes prereq_2 , course_1 and takenBy_1 to nodes in the S02 subtree (student_2 , sid_2 and name_2). Note that $\{(\text{takenBy}_2, \text{student}_2), (\text{takenBy}_2, \text{sid}_2), (\text{takenBy}_2, \text{name}_2), \dots\}$, i.e., the connection between node takenBy_2 (and thus course_2) and the S02 subtree still holds and is not included in Δ_{M_1} .

Complexity. The worst-case time complexity of Algorithm $\Delta_{(\mathbf{M},\mathbf{L})}\text{insert}$ is $O(|E_A| + |E_{N_C}| + (|N_C| + |r[p]|)n + |N_A||E_A| + |N_A|n)$, where (a) $|N_A|$ is the number of distinct nodes, and $|E_A|$ is the number of edges in the inserted subtree $ST(A, t)$, (b) $|N_C|$ is the number of common nodes in \mathbf{L} and \mathbf{L}_A , $|E_{N_C}|$ is the number of those edges that connect nodes of N_C in either T or $ST(A, t)$, and (c) n is the number of distinct nodes in T . In practice $|N_C| < |N_A| < |E_A| \ll n \ll |V|$. The first and second factors are the cost of computing \mathbf{L}_A and \mathbf{L}_{N_C} , respectively, and the third factor is the cost of maintaining \mathbf{L} , where $\text{swap}()$ is called at most $2|N_C| + |r[p]|$ times and each takes at most $O(n)$ time. Note that $\text{swap}(\mathbf{L}, u, v)$ is in $O(|\mathbf{L}[u : v]|)$ time, which is usually much smaller than n . The fourth factor is the cost of computing the reachability matrix for $ST(A, t)$, while the last factor is the cost of maintaining the reachability between nodes in $ST(A, t)$ and the nodes in T . The worst-case time complexity of Algorithm $\Delta_{(\mathbf{M},\mathbf{L})}\text{delete}$ is $O(n|V|)$, which is the cost of computing new ancestors for nodes in \mathbf{L}_R . For each node in \mathbf{L}_R we visit its parents once, which in total takes $O(|V|)$ time in the worst-case (in practice it is often much smaller than $|V|$); at each visit, the algorithm takes $O(n)$ time.

We make the following observations on the analysis. (a) The analysis given above is the worst-case complexity. While it seems no better than the complexity of re-computing \mathbf{M} and \mathbf{L} from scratch, in practice the updated XML view $\Delta_X(T)$ typically differs slightly from the old view T , and $|r[p]|$ and $|\text{anc}(r[p])|$ are of-

ten far smaller than n . (b) \mathbf{L}_A and \mathbf{L}_R are typically much smaller than \mathbf{L} ; this makes the fourth factor of the complexity of $\Delta_{(\mathbf{M},\mathbf{L})}\text{insert}$ and the complexity of $\Delta_{(\mathbf{M},\mathbf{L})}\text{delete}$ much smaller than $n|V|$ in practice. (c) As mentioned earlier, the computation of $\Delta_{\mathbf{M}}$ and updating of \mathbf{L} is in fact conducted in the background. (d) Our experimental study verified that the incremental approach is far more efficient than the batch counterpart.

4 Updating Relational Views

In this section, we extend the study of relational view updates by providing complexity results and techniques for processing SPJ view updates under key preservation. These results are not only important for updating XML views, defined in terms of ATGs, but are also useful for studying relational view updates.

4.1 Key Preservation and Relational View Updates

Foremost, we propose a mild condition on SPJ views. Then we show that this condition simplifies the analysis of relational view updates.

Key Preservation. Consider an SPJ query $Q(R_1, \dots, R_k)$ that takes base relations R_1, \dots, R_k of \mathcal{R} as input, and returns tuples of the schema $R(\mathbf{a})$. We say that Q is *key preserving* if for each R_i the primary key of R_i is included in \mathbf{a} (with possible renaming). That is, the primary keys of all the base relations involved in Q are included in the projection fields of (the SPJ query) Q .

Next, we make a couple remarks on key preservation. First, key preservation is far less restrictive than other conditions proposed in earlier work for handling relational view updates (e.g., [2, 3]; see Section 6). Second, every SPJ query in the definition of an ATG view σ can be made key-preserving by extending its projection-attribute list to include the primary keys. The extension does not affect the expressive power of ATGs. For example, Q_3 in σ_0 of Fig.2 can be made key-preserving by adding *e.cno* to its select clause. Third, key preservation is a property of a view. This property does not assume how the base relations are defined or specified. The key attributes on the view would be useful for translating the view update efficiently. Thus, in the sequel we assume w.l.o.g. that all the queries in ATGs are key-preserving.

Analysis. We consider the following decision problem:

PROBLEM:	SPJ View Updatability Problem
INPUT:	A collection of views \mathcal{V} defined as SPJ queries <i>under key preservation</i> , a relational database I of schema \mathcal{R} , and a group view update Δ_V .
QUESTION:	Is there a group update Δ_R on the database I such that $\Delta_V(\mathcal{V}(I)) = \mathcal{V}(\Delta_R(I))$?

Here Δ_V consists of either only tuple deletions or only tuple insertions, as produced by the translation algorithm of the last section. These deletions and insertions in Δ_V are translated to deletions and insertions in Δ_R , respectively. We use V to denote the view $\mathcal{V}(I)$.

It is known^[18] that without key preservation, the updatability problem is already NP-hard for a single deletion and a single PJ view, i.e., when Δ_V consists of a single deletion and \mathcal{V} is a view defined with projection and join operators only. In contrast, we show that key preservation simplifies the updatability analysis for a collection of SPJ views and group deletions.

Theorem 1. *For group view deletions Δ_V , the SPJ view updatability problem is in PTIME.*

In Subsection 4.2 we present a PTIME algorithm for computing database deletions Δ_R from view deletions Δ_V which suffices to prove Theorem 2.

However, the problem is intractable for insertions under key preservation; the lower bound can be verified by reduction from the non-tautology problem, which is NP-complete (cf. [29]).

Theorem 2. *The SPJ view updatability problem is NP-complete when Δ_V has a single insertion and \mathcal{V} has a single view.*

Proof. An NP algorithm for checking CQ view updatability works as follows: it first guesses a group insertion Δ_R and then checks whether $\mathcal{V}(\Delta_R(I)) = \Delta_V(V)$, which can be done in PTIME (data complexity).

We next show the problem is NP-hard, by reduction from the non-tautology problem. Consider an instance of the problem: $\phi = C_1 \vee \dots \vee C_n$, where all the variables in ϕ are x_1, \dots, x_k , C_j is of the form $l_{j_1} \wedge l_{j_2} \wedge l_{j_3}$, and l_{j_i} is either x_s or \bar{x}_s , $s \in [1, k]$. The problem is to determine whether there is a truth assignment such that ϕ is false, i.e., ϕ is not valid. This problem is known to be NP-complete.

Given ϕ , we define a relational database I , a single CQ view \mathcal{V} under key preservation, and a single view insert Δ_V on $V = \mathcal{V}(I)$, such that ϕ is not valid iff there exists Δ_R and $\mathcal{V}(\Delta_R(I)) = \Delta_V(V)$.

Relational Database I. The database consists of three base relations, R , R_ϕ and R_E , defined as follows.

- $R(A, B)$, where A is the key of the relation and B

is a Boolean. Intuitively, A is to hold a number in $[1, k]$ encoding a variable, and B is a truth value (T or F). That is, $R(A, B)$ is a truth assignment for ϕ . Initially $R(A, B)$ consists of a single special tuple $(0, T)$.

- $R_\phi(j, j_1, X_1, j_2, X_2, j_3, X_3)$, where j is the key of the relation. Initially, for each $C_j = l_{j_1} \wedge l_{j_2} \wedge l_{j_3}$, there is a tuple $(j, l_{j_1}, X_1, l_{j_2}, X_2, l_{j_3}, X_3)$ in R_ϕ such that l_{j_i} is s if $l_{j_i} = x_s$ or \bar{x}_s , X_i is T if $l_{j_i} = x_s$, and X_i is F if $l_{j_i} = \bar{x}_s$. Intuitively, each of these tuples in R_ϕ codes a clause in ϕ . A special tuple $(0, 0, T, 0, T, 0, T)$ is also in R_ϕ .

- $R_E(e_1, e_2, \dots, e_k)$, where e_1, \dots, e_k are the key. Intuitively e_i is to code i in $[1, k]$. Initially, R_E consists of a single special tuple $(0, \dots, 0)$.

View. We define a single view $\mathcal{V} = V_1 \times V_2$ in terms of conjunctive queries and under key-preservation as follows.

- $V_1 = \pi_{j, j_1, j_2, j_3} \sigma_C(R_1 \times R_2 \times R_3 \times R_\phi)$, where R_1, R_2, R_3 are renaming of R , and C is a Boolean condition $c_1 \wedge c_2 \wedge c_3$, in which c_i is $R_i(A) = R_\phi(j_i) \wedge R_i(B) = R_\phi(X_i)$ ($i = 1, 2, 3$). Intuitively, C holds if and only if one of the C_j 's is true.

- $V_2 = \pi_{e_1, e_2, \dots, e_k} \sigma_D(R_E \times R_1 \times R_2 \times \dots \times R_k)$, where R_1, R_2, \dots, R_k are renamings of R , and D is a Boolean condition $\bigwedge_{i=1}^k R_i(A) = R_E(e_i)$.

Initially $V = \mathcal{V}(I)$ has a single tuple $(0, \dots, 0)$ ($k + 40$'s).

View Insert. We define Δ_V to insert a single tuple $(0, 0, 0, 0, 1, \dots, k)$ into V .

We next verify that Δ_V is side-effect free iff ϕ is not a tautology. Indeed, if ϕ is not a tautology, then there is a truth assignment μ such that ϕ is false, and thus C_j is false w.r.t. μ . We define Δ_R based on μ as follows: insert tuples to $R(A, B)$ such that (i, T) is inserted into $R(A, B)$ iff $\mu(x_i) = T$, and (i, F) is inserted into $R(A, B)$ iff $\mu(x_i) = F$; furthermore, insert $(1, \dots, k)$ into R_E . Then obviously Δ_V is side-effect free. Conversely, suppose that there is Δ_V that is side-effect free. Then $(1, \dots, k)$ needs to be inserted into R_E , and a unique tuple of the form (i, X) needs to be inserted into the base relation R for each $i \in [1, k]$ due to the key constraint on R , such that Δ_V is indeed an update on the view V . Here X is either T or F , and thus after the insertion of Δ_V , $R(A, B)$ contains a valid truth assignment for ϕ . Since Δ_V is side-effect free, V_1 will remain $(0, 0, 0, 0)$ after Δ_V is performed. That is, C_j remains false. Thus ϕ is not a tautology.

4.2 Processing Group Deletions

We give a PTIME algorithm for computing database tuple deletions Δ_R from a group of view deletions Δ_V . Let V_Q be the view $Q(I)$, and consider a tuple t in Δ_V

that is to be deleted from V_Q . The key preservation condition allows us to identify, for each S_j , a *unique* tuple t_j via its key in t , such that t_1, \dots, t_l produce t via Q . Let us use $Sr(Q, t)$ to denote the set consisting of all the pairs (S_j, t_j) , referred to as the *deletable source* of t in V_Q . Observe the followings. (a) Deleting any t_j from S_j suffices to remove t from V_Q . (b) Deletion of a source tuple t_j from V_Q is *side-effect free* if and only if (S_j, t_j) is not in the deletable source of any tuple $t' \in \mathcal{V}(I) \setminus \Delta_V$ that is to remain in the view after Δ_V is carried out. From these one can see that t can be deleted from V_Q if and only if there exists $(S_j, t_j) \in Sr(Q, t)$ such that for all $Q' \in \mathcal{V}$ and all t' that are in $Q'(I)$ but not in Δ_V , (S_j, t_j) is not in $Sr(Q', t')$. Note that as far as *the updatability problem* is concerned, deleting any of such t_j suffices, i.e., one can choose an arbitrary t_j from $Sr(Q, t)$ satisfying the condition (b) given above, if there exists any.

<p>Input: a view definition \mathcal{V}, a relational database I, the view $V_Q = Q(I)$ for each $Q \in \mathcal{V}$, and a group deletion Δ_V.</p> <p>Output: a group update Δ_R on I if it exists.</p> <ol style="list-style-type: none"> 1. $\Delta_R := \emptyset$; 2. for each (Q, t) in Δ_V 3. compute $Sr(Q, t)$, the deletable source of t in V_Q; 4. for each Q' in \mathcal{V} and each t in $V_{Q'}$ but not in Δ_V 5. compute $Sr(Q', t')$; 6. for each (Q, t) in Δ_V 7. if there exists (S_j, t_j) in $Sr(Q, t)$ such that (S_j, t_j) is not in $Sr(Q', t')$ for any Q' in \mathcal{V} and any t' in $V_{Q'}$ but not in Δ_V 8. then $\Delta_R := \Delta_R \cup \{(S_j, t_j)\}$; 9. else reject Δ_V and exit; 10. return Δ_R

Fig.9. Algorithm delete.

On this basis, we give Algorithm delete in Fig.9, which is self-explanatory. The worst-case complexity of Algorithm delete is in $O(|\Delta_V|(|\mathcal{V}(I)| - |\Delta_V|))$ time.

Minimal Deletions. Algorithm delete does not address which Δ_R to select if multiple valid Δ_R 's exist. In the presence of multiple Δ_R 's it is natural for one to choose the *smallest* set Δ_R of tuples to delete, i.e., a set Δ_R such that $|\Delta_R|$ is the smallest. The *minimal view deletion problem* is thus to find, given a collection \mathcal{V} of view definitions, a database I and view deletions Δ_V , the smallest set of tuple deletions Δ_R such that $\Delta_V(\mathcal{V}(I)) = \mathcal{V}(\Delta_R(I))$. However desirable, the minimal view deletion problem is intractable, even under the key preservation condition. The lower bound can be verified by reduction from the minimal set cover problem, which is known to be NP-complete (cf. [29]).

Theorem 3. *For SPJ views under key preservation, the minimal view deletion problem is NP-complete.*

Proof. We show the problem is NP-hard by reduction from the minimal set cover problem. An instance of the minimal set cover problem consists of a collection C of subsets of a finite set S ; it is to find a subset $C' \subseteq C$ such that every element in S belongs to at least one member of C' and moreover, $|C'|$ is minimal.

Given S and C , we define an instance of the minimal view deletion problem. Let $S = \{x_i \mid i \in [1, n]\}$. We construct $|C|$ many base tables, n CQ views and a group view deletion, as follows.

1) For each $S_j \in C$, we define a base relation R_j consisting of a single column.

Let I_j , the instance of R_j , be $\{j\}$, and let the database instance I be the collection of all I_j 's defined above.

2) For each x_i , let T_i be the collection of all the subsets in C that contain x_i . Enumerate the elements of T_i as $(S_{i^1}, \dots, S_{i^{n_i}})$. Define $V_i = R_{i^1} \times \dots \times R_{i^{n_i}}$. Note that $V_i(I) = (i^1, \dots, i^{n_i})$. Let \mathcal{V} be the collection of V_i 's for $i \in [1, n]$.

Obviously, the views defined as above are key-preserving.

3) The group deletion Δ_V is to remove all tuples from all the views.

Note that the tuple is removed from V_i without side effect if and only if the tuple from any R_{i^j} is removed.

The minimum view deletion problem is to find a smallest set of the base relations $R_1, \dots, R_{|C|}$ from which tuples are removed, while ensuring that the view tuples from V_i for $i \in [1, n]$ are deleted without side effect.

We next verify that the construction above is indeed a reduction from the minimum set cover problem. First suppose that C' is a minimal cover of S . We define Δ_R such that it consists of deletion of the tuples from each base relation in $\{R_j \mid S_j \in C'\}$. Clearly, $\mathcal{V}(\Delta_R(I)) = \Delta_V(\mathcal{V}(I)) = \emptyset$ since C' is a cover of S . Furthermore, Δ_R is minimal since C' is minimal. Conversely, suppose that Δ_R is a solution to the minimal view deletion problem. Then let C' be the subset of C such that an element S_j of C is in C' if and only if Δ_R involves deletion of the tuple from the corresponding relation R_j . To see that C' is a cover of S , note that $\mathcal{V}(\Delta_R(I)) = \Delta_V(\mathcal{V}(I)) = \emptyset$, and thus for each $i \in [1, n]$, some set R_{i^j} is in C' . Moreover, C' is minimal since Δ_R is minimal. \square

4.3 Processing Group Insertions

Theorem 2 shows that any practical algorithm for handling group view insertions is necessarily heuristic. We approach this by reducing the SPJ view insertion problem to SAT, one of the most studied NP-complete

problems. This allows us to leverage a well-developed SAT solver^[30] to efficiently compute Δ_R if it exists.

An instance of SAT (cf. [29]) is $\phi = \bigwedge_{i \in [1, n]} C_i$, where C_i is a disjunction of literals, i.e., propositional variables or their negation. It is to find a truth assignment μ that satisfies ϕ , if such a μ exists.

Below we outline our heuristic algorithm, referred to as Algorithm *insert*. The algorithm takes the same input as that of Algorithm *delete* given in Fig.9, namely, $\mathcal{V}, I, V_Q(I)$ for each $Q \in \mathcal{V}$, and Δ_V , except that tuples in Δ_V are to be inserted into the views. It either finds a set of insertions D_R such that $\Delta_V(\mathcal{V}(I)) = \mathcal{V}(\Delta_R(I))$, or it rejects Δ_V . The major steps of the algorithm can be described as follows.

- Compute a propositional logic formula ϕ (i.e., a SAT instance) from $\mathcal{V}, I, V_Q(I)$'s, and Δ_V , such that ϕ is satisfiable if and only if there exists D_R such that $\Delta_V(\mathcal{V}(I)) = \mathcal{V}(\Delta_R(I))$.
- Utilize an existing heuristic tool^[30] for SAT to process ϕ .
- If the tool returns a truth assignment μ that satisfies ϕ , compute Δ_R from μ ; otherwise reject the view updates Δ_V as well as Δ_X .

We next illustrate each of the three steps in detail.

Deriving ϕ . The encoding is a little involved. It takes four steps.

First, we derive tuples that have to be present in base relations so that Δ_V can be computed through queries in \mathcal{V} . Consider (Q, t) in Δ_V , which indicates that tuple t is to be inserted into the view $Q(I)$, as illustrated in Subsection 4.2. For each t and each relation R_i involved in Q , we derive an R_i tuple template $t_i = (\mathbf{a}_i, \mathbf{b}_i, \mathbf{z}_i)$ from t and Q , where \mathbf{a}_i corresponds to the (primary) key of R_i , \mathbf{b}_i to the other columns of R_i whose values can be determined from t , and \mathbf{z}_i to variables whose values are unknown. Note that \mathbf{a}_i is known due to the key preservation condition. If there is no tuple t' in the instance I_i of R_i with the key \mathbf{a}_i , we add t_i to a set X_i . Note that there are no more than $|Q| |\Delta_V|$ many tuple templates in these X_i 's.

Example 8. Consider two relations R_1, R_2 and an SPJ view Q given below, where keys are underlined:

$$R_1 = (\underline{A} : \text{int}, B : \text{bool}), \quad R_2 = (\underline{C} : \text{int}, D : \text{bool}), \\ Q = \pi_{A,C} (\sigma_{B=D} (R_1 \times R_2)).$$

Suppose that tuples (a, c) and (a, c') are to be inserted into $Q(I)$. Then X_1 contains a tuple template (a, x_1) and X_2 contains (c, x_2) and (c', x_3) , if no tuple bearing the key a is already in I_1 and no c, c' tuples are in I_2 . For $(a, c), (a, c')$ to be inserted into the view, it is necessary that (a, x_1) is inserted into I_1 after x_1 is

instantiated to a truth value, and that $(c, x_2), (c', x_3)$ are added to I_2 .

Second, we “evaluate” each view query Q on the database I incremented by adding X_i to I_i . For succinctness of presentation, we present the details of the evaluation in Appendix A. In the evaluation we “instantiate” variables in the tuple templates, as well as the selection (conjunctive) condition in Q . In Example 8, for instance, the evaluation yields view tuples (a, c) with condition $x_1 = x_2$, and (a, c') with condition $x_1 = x_3$. We then inspect the result of Q to determine whether or not tuple templates may yield side effects. Specifically, for each tuple t in the result, if it is in neither the view nor Δ_V , we consider the following cases.

- (a) If t is not associated with any condition, i.e., it certainly has some side effects, then we *reject* the view updates Δ_V and Δ_X immediately.
- (b) If t has a condition in which at least one variable represents an attribute with an infinite domain, we can always pick a distinct value for the variable that makes the condition false. This eliminates t from the result and thus t does not yield a side effect.
- (c) If t has a condition ϕ_t in which all variables correspond to attributes with a finite domain, we add the negation $\neg\phi_t$ as a conjunct to the logic formula ϕ that we are constructing.

Furthermore, for each t that is in Δ_V , we also add its associated condition ϕ_t as a conjunct to ϕ . Observe that these conjuncts are bounded by $|\Delta_V|$, and those in case (c) involve only attributes with a finite domain (with a fixed cardinality, a *constant*).

Example 9. Referring to Example 8, the conjuncts added to ϕ in the second step are $x_1 = x_2$ and $x_1 = x_3$.

Third, to complete the construction of ϕ , for each variable x bounded to a finite domain, we add the following formula to ϕ as a conjunct: $x = c_1 \vee \dots \vee x = c_k$, where c_1, \dots, c_k are all the values in that domain. In Example 8, for instance, we add $x_i = \text{true} \vee x_i = \text{false}$ for $i \in [1, 3]$.

Finally, we convert ϕ to a propositional formula (i.e., a SAT instance). We use propositional variables and their negation to code variables introduced in the encoding: p for $x = c$ and \bar{p} for $y \neq c$. We also add conjuncts $(\bar{p} \vee \bar{p}')$ to ensure that p and p' cannot be both true if, e.g., p codes for $x = c$, p' for $x = c'$, and $c \neq c'$.

The correctness of the reduction is ensured by the following.

Theorem 4. *If Δ_V is not rejected during the coding, then ϕ is satisfiable iff there is Δ_R such that $\Delta_V(Q(I)) = Q(\Delta_R(I))$.*

Proof. We verify that if Δ_V is not rejected during

the coding of an instance Q , Δ_V and I of the CQ view insertion problem, then there exists a truth assignment μ that satisfies ϕ_Q if and only if there exists Δ_R such that $\Delta_V(Q(I)) = Q(\Delta_R(I))$.

Assume that there exists a truth assignment μ that satisfies ϕ_Q . Then we define Δ_R as follows. For each X_j and each tuple template t in X_j , we assign a value to each variable z in t based on μ . If z is bounded in ϕ_Q by $(z = c)$ for some constant c and $(z = c) \leftrightarrow x$, then we let $z = c$ if $\mu(x)$ is true; after this process if z is not assigned any value, z must be a free variable that ranges over an infinite domain τ_i and thus we can always pick a value c' for z without violating ϕ . Indeed, our coding distinguishes (bounded) variables with a finite domain from those (free) variables with an infinite domain, and encodes possible value selections of those variables having a finite domain in terms of additional clauses; the coding ensures that the value of z can be picked without causing side effects. For each relation I_i , let Δ_R^i consist of all these instantiated tuple templates from all X_j 's that are a renaming of R_i . Let Δ_R be the collection of Δ_R^i 's for $i \in k$. Then $\Delta_V(Q(I)) = Q(\Delta_R(I))$. Indeed, these newly inserted tuples do not produce view tuples that have a key of R_i that is not already in Δ_V , since otherwise this had been caught in the coding process and Δ_V would have been rejected. Furthermore, these newly insertions do not yield tuples that are not in Δ_V but share keys of Δ_v , as ensured by the coding ϕ_Q . Finally, all the tuples in Δ_V are coded in ϕ_Q and are guaranteed to be produced by $\Delta_R(I)$. Thus Δ_R carries out the desired view insertions without side effects.

Conversely, assume that there exists a group update Δ_R to I such that $\Delta_V(Q(I)) = Q(\Delta_R(I))$. Then by reversing the derivation of Δ_R given above we can define a truth assignment μ to propositional variables in ϕ_Q ; indeed, we let $\mu(x)$ be true iff $(z = c)$ and $(z = c) \leftrightarrow x$ are in ϕ_Q , if z has the value c in Δ_R . It is easy to verify that μ satisfies the formula ϕ_Q . \square

Processing ϕ . We invoke Walksat^[30] with ϕ as the input. Walksat, an extension of GSAT, employs an efficient approximation algorithm to solve the maximum satisfiability problem. If ϕ is satisfiable, it finds a truth assignment μ for ϕ above a certain percentage.

Computing Δ_R . If μ is found, we derive Δ_R , i.e., the set of tuples to be inserted into each I_i , by instantiating variables in the tuple templates in X_i 's based on μ and the interpretation of propositional variables given above. More specifically, for each tuple template t in X_i , we assign a value to each variable z in t based on μ : if z is bounded in ϕ by $(z = c)$ for some constant c and $(z = c) \leftrightarrow x$, then we let $z = c$ if $\mu(x)$ is true. After this process if z is not assigned any value, then

either (a) z ranges over an infinite domain and thus we can always pick a value c' for z that is not in the active domain of the database, or (b) the value of z does not have any impact on the satisfaction of ϕ ; in both cases we can find a value for z without violating ϕ . Then Δ_R consists of query templates instantiated by these values.

If μ is not found, we reject Δ_V and Δ_X . Note that Walksat^[30] may not find a truth assignment for ϕ even if ϕ is satisfiable, since SAT is intractable and so is the view insertion updatability problem (Theorem 2). However, this only happens within a certain percentage given the excellent performance of Walksat^[31].

Complexity. From the construction of ϕ one can see that its size $|\phi|$ depends on $|\Delta_V|$, \mathcal{R} and $|Q|$ only, whereas the size of the database I is irrelevant. Our algorithm has a low (data) complexity, and is effective in practice as verified by our experimental study.

5 Experimental Study

We conducted an experimental study of our proposed view update mechanism in order to verify its effectiveness. The reported numbers are warm numbers and are the average of five runs per query. The standard deviation of the reported numbers is no greater than 5%.

All experiments were conducted on a synthetic dataset. It allows us to produce highly nested XML views with diverse structure and to have more control over the experimental settings (e.g., data size). (We have not found any real highly-recursive relational dataset to perform our experiment.) The dataset consists of four base relations: $C(\underline{c}_1, \dots, c_{16})$, $F(\underline{f}_1, \dots, f_{16})$, $H(\underline{h}_1, \underline{h}_2)$ and $C_U(\underline{c}'_1, \dots, c'_{16})$, where underlined attributes indicate keys. The domain of f_1 is equal to the domain of c_1 and c'_1 . The remaining C and F attributes were used to control how many joining C and F tuples were filtered out. The domains of h_1 and h_2 are the same as the domain of c_1 . The generator ensured that 1) for each $c \in C \cup C_U$ there would be on average three tuples $h \in H$, where $c_1 = h_1$, and 2) $h_1 < h_2$, where $(h_1, h_2) \in H$. The universe of C , namely C_U , consisting of 100M C -tuples, ensured that whenever h_2 joined with c_1 it always yielded a C -tuple. The sizes of F and H were proportional to the size of C , which we use for reporting the size of the synthetic database; more specifically, the size we report is $|C|$, which ranges from 1000 to 1000000 tuples, while $|F| = |C|$ and $|H| \simeq 3|C|$. We defined an ATG view of the relations C , F and H ; as indicated in Fig.10(a), the C nodes in the view were recursively defined, and a

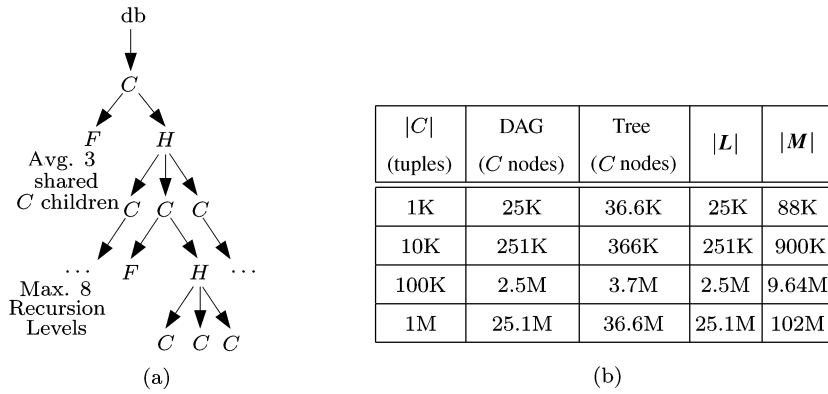


Fig.10. Description of the datasets. (a) XML view. (b) Statistics of the datasets.

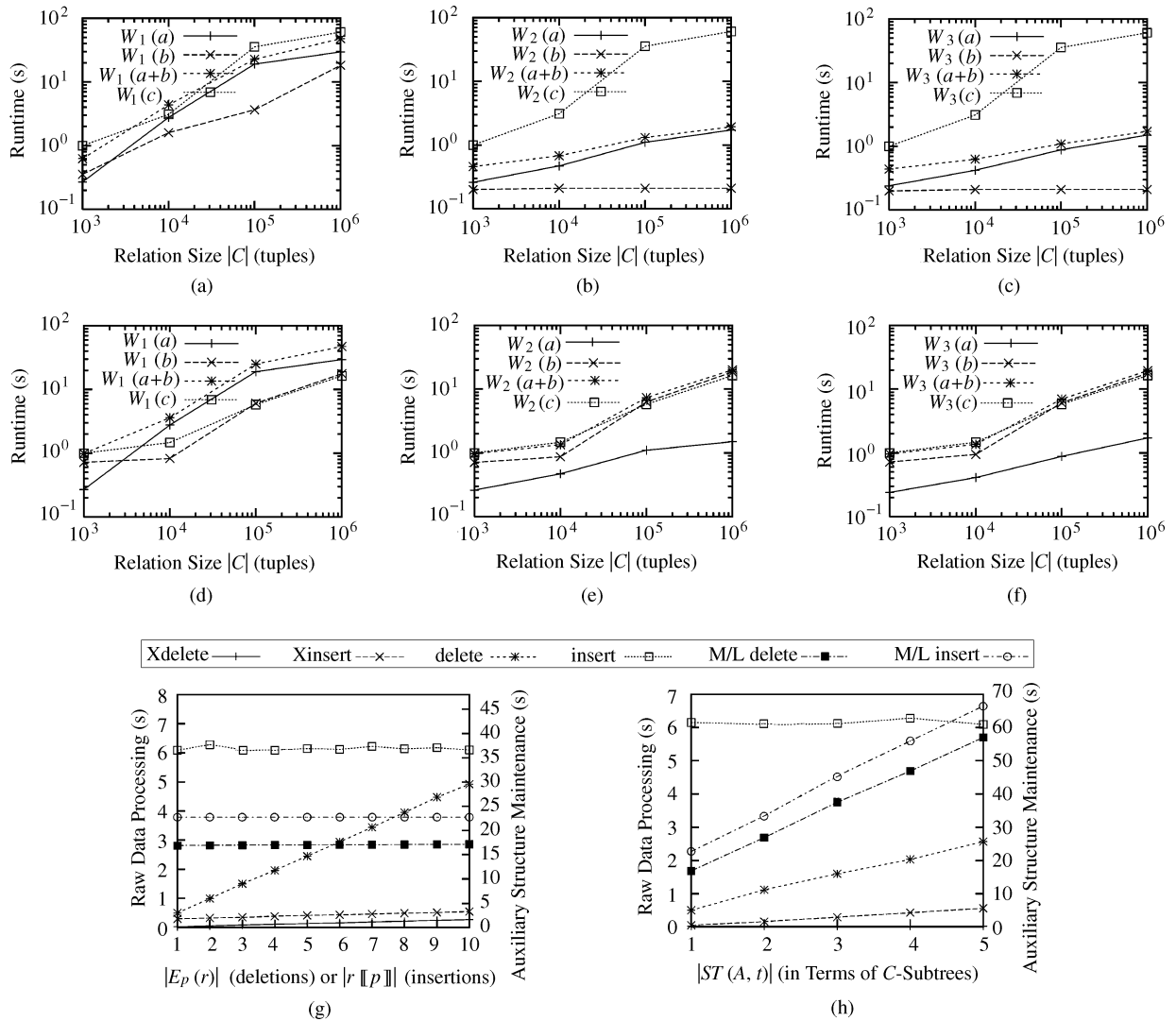


Fig.11. Update performance as a function of the size of the underlying relational database and the view update size. (a) W_1 deletion. (b) W_2 deletion. (c) W_3 deletion. (d) W_1 insertion. (e) W_2 insertion. (f) W_3 insertion. (g) Varying $|r[p]|$ or $|E_P(r)|$. (h) Varying $|ST(A, t)|$.

recursion of C in the view can be understood as

$$\pi_{c_1, f_1, h_1, h_2}(\sigma_{c_1=f_1 \wedge f_1=h_1 \wedge h_2=c'_1 \wedge c_2=f_2 \wedge c_3=f_3 \wedge c_4=f_4} (C \times F \times H \times C_U)).$$

Recall that [11,13] cannot handle recursions of C in the view. Compression was achieved by sharing C subtrees; in our dataset subtree sharing accounted for 31.4% of C instances. Fig.10(b) lists some statistics on the number of published C subtrees and their compressed DAGs, and the corresponding sizes of the reachability matrix \mathbf{M} and topological order \mathbf{L} .

Varying Database Size. We generated two random update workloads over the XML view, one for insertions, and one for deletions; each workload consisted of three update classes, each class including ten operations. The classes were characterized by the XPath queries used for defining the updates. Specifically, class W_1 involved XPath queries using “//” and value-based filters; XPath queries in W_2 used “/” and value-based filters; finally, W_3 contained XPath queries with “/”, and both structural and value filters. The times we report include the followings: (a) the time to evaluate XPath queries (Subsection 3.2); (b) the time to translate Δ_X to Δ_V (Algorithms Xinsert and Xdelete) and subsequently Δ_V to Δ_R (Section 5), and the time to execute the update; and (c) the time to maintain the auxiliary structures (Algorithms $\Delta_{(\mathbf{M}, \mathbf{L})}$ insert and $\Delta_{(\mathbf{M}, \mathbf{L})}$ delete). Note that (c) is executed in the *background*.

Figs.11(a), 11(b) and 11(c) show the performance of the deletion algorithms for W_1 , W_2 and W_3 , respectively. We plot the runtime of performing the updates broken into their (a), (b) and (c) above constituents for various relational database sizes. Note that both x - and y -axes use a logarithmic scale. As shown, the algorithms scale linearly with the size of the relational database. It is evident that deletion time is dominated by XPath evaluation. Observe that although the cost for (c) is relatively high, it is performed in the background. $W_1(b)$ is the highest reported time among the three workloads since its XPath queries generate more edges (i.e., $E_p(r)$), which are then examined by Algorithm delete.

Similar results are reported for insertions, as shown in Figs.11(d), 11(e) and 11(f) for W_1 , W_2 and W_3 , respectively (again, using logarithmic scales). The size of the inserted subtree was fixed. The SAT solver^[30] we used returned a truth assignment in 78% of the cases and we only report the time for insertions where the SAT solver successfully returned a truth assignment. As in the case of deletions, our insertion algorithms also scale linearly with the size of the database.

Varying Update Size. For these experiments, we fixed $|C|$ to be 100K tuples. Fig.11(g) shows the performance of each algorithm as we varied $|E_p(r)|$ (see Subsection 4.2) for deletions and $|r[p]|$ for insertions, while keeping $ST(A, t)$ a constant single C -subtree. The runtimes for Algorithms Xinsert, Xdelete, delete and insert are measured on the left y -axis, while the runtimes for algorithms $\Delta_{(\mathbf{M}, \mathbf{L})}$ insert and $\Delta_{(\mathbf{M}, \mathbf{L})}$ delete are measured on the right y -axis. The translation time from Δ_X to Δ_V for Algorithm Xinsert (resp. Algorithm Xdelete) increases slightly as $|r[p]|$ (resp. $|E_p(r)|$) increases, as expected. The slope of the curve for Algorithm delete is large, as the increase of $|E_p(r)|$ involves more database queries to determine the source tuples to be deleted. The performance of Algorithm insert, which models the translation of Δ_V to Δ_R for insertion workloads, is dominated by the coding time. As $|C|$ is far larger than $|ST(A, t)|$ and $|r[p]|$, and the number of database queries required remained fixed, the coding time remains roughly constant, though the size of the resulting coding increases; however, that only results in a non-observable increase in the SAT solver’s runtime keeping the curve relatively flat. The performance of Algorithm $\Delta_{(\mathbf{M}, \mathbf{L})}$ insert and Algorithm $\Delta_{(\mathbf{M}, \mathbf{L})}$ delete is almost unaffected by $|r[p]|$ (resp. $|E_p(r)|$) since $|ST(A, t)|$ is fixed.

Similar results are shown in Fig.11(h) where we varied the size of $|ST(A, t)|$ while fixing $|E_p(r)| = 1$ and $|r[p]| = 1$. The performance of Algorithm Xdelete remains unchanged and its runtime is negligible as it nearly overlaps with the x -axis for a fixed $|E_p(r)|$. Algorithm Xinsert scales linearly with the update size $|ST(A, t)|$ as it needs to process $ST(A, t)$ to generate Δ_V . Algorithms $\Delta_{(\mathbf{M}, \mathbf{L})}$ insert and $\Delta_{(\mathbf{M}, \mathbf{L})}$ delete evidently scale linearly w.r.t. the update size for reasons similar to the ones outlined above.

Table 1. Incremental Maintenance of \mathbf{L} and \mathbf{M} vs. Recomputation

Sizes	Incremental (s)		Recomputation (s)	
	Insertion	Deletion	\mathbf{L}	\mathbf{M}
1K	1.0	1.0	6.3	9.8
10K	4.6	3.1	86	288
100K	22.7	16.9	631	3 600
1M	84.2	61.5	8611	14 000

Effectiveness of Incremental Maintenance. The cost of incrementally maintaining the reachability matrix \mathbf{M} and the topological order \mathbf{L} as opposed to recomputing them is shown in Table 1. The first column presents the size of the relational datasets. The total time needed for incrementally maintaining both auxiliary structures is given in the second column for Algorithm $\Delta_{(\mathbf{M}, \mathbf{L})}$ insert

and in the third column for Algorithm $\Delta_{(M, L)}\text{delete}$. The time for recomputing each structure is shown in the last two columns. As expected, the advantages of incremental maintenance become more prominent as the size of the data increases.

6 Related Work

Commercial database systems^[5–7] provide support for defining XML views of relations and restricted view updates. For example, IBM DB2 XML Extender^[5] supports limited view maintenance. It supports only propagation of updates from relations to simple XML views but does not support updates through XML views. Oracle XML DB^[6] provides XMLType views to wrap relational tables in XML views using SQL statements. It does not support recursive XPath queries and update operations. It does not allow updates on XML (XMLType) views. In SQL Server^[7], updates of XML views generated by an annotated schema are represented in an *updategram*, a data structure for users to express changes in XML data, by specifying the difference of the images of the data before and after a change. Then, the system generates the SQL update statements that correspond to the updategram. However, the views supported are very restricted: only key-foreign key joins are allowed; neither recursive views nor updates defined in terms of recursive XPath expressions are supported.

There have been recent studies on updating XML views published from relational data^[11,13]. In [11], XML views are defined as *query trees* and are mapped to relational views. XML view updates are translated to updates of relations only if XML views are well-nested (i.e., joins are through keys and foreign keys), and if the query tree is restricted to avoid duplication. Existing technique on relational view updates is reused for the update translation. [32] studies a *round-trip* mapping that shreds XML data into relations in order to ensure that XML views are always updatable. More general XML views where duplication is allowed is considered in [33]. A detailed analysis on deciding whether or not an update on XML views is translatable to relational updates and the decision algorithms are presented in [13]. A framework for [13] is demonstrated in [12]. The limitations of previous work, e.g., [11–13], have been discussed in Section 1.

There has been a host of work^[1–7] on relational view updates. [2] provides algorithms for translating restricted view updates to base-table updates without side effects in the presence of certain functional dependencies. The algorithm in [3] handles translation

(which may allow side effects) for a restricted class of SPJ view: base tables may only be joined on keys and must satisfy foreign keys; a join view corresponds to a single tree where each node refers to a relation; join attributes must be preserved; and comparisons between two attributes are not allowed in selection conditions. Clearly, our key preservation condition is less restrictive than those considered in [2, 3]. There has also been work^[1,4] on relational view complements. An update of a view can be correctly translated into updates of base relations if and only if there exists at least one complement that is not changed by the view update, i.e., a constant complement exists. Obviously, it is easier to decide the translatability of a view update with a small view complement. However, finding a minimal view complement is NP-complete^[1]. Furthermore, the problem of constructing an update translator given a complement view remains largely unexplored.

An algorithm for deletion translation using data lineage is given in [34], which is very different from Algorithm *delete* of Fig.9. The algorithm runs in exponential time in the worst case. However, if the view is key-preserving, the computation of data lineage is simplified and the algorithm can determine a side-effect free deletion in PTIME.

Commercial DBMSs^[5–7] allow updates on very restricted relational views (while users may specify updates manually with INSTEAD OF triggers). For example, for views to be deletable, IBM DB2^[5] restricts the FROM clause to reference only one base table.

Few complexity bounds are known for (relational) view updates. The complexity of deletion on views is given in [18]. To the best of our knowledge, this work and the work on annotation propagation in [35] are the only work that establishes complexity bounds for both deletion and insertion on views under key preservation.

A number of XPath evaluation algorithms have been proposed (e.g., [25, 36]) for trees and cannot answer XPath queries on DAGs. Path query evaluation has been studied in [23, 24] for DAGs. However, they cannot be directly used in the context of XML view updates as discussed in Subsection 3.2.

7 Conclusions

We have proposed new techniques for updating XML views published from relational data. The novelty of our technique consists of (a) the ability to handle XML updates defined with *recursive* XPath queries over (possibly) *recursively defined* XML views; (b) the first method to rewrite XML updates into group updates on relational views that represent a DAG *com-*

pression of an XML view, capturing XML view-update side effects; (c) a key-preservation condition on SPJ views that is less restrictive than constraints imposed by previous work but simplifies the analysis of relational view updates; and (d) efficient (heuristic) algorithms for handling *relational* SPJ view updates under key preservation, along with complexity results. Our results contribute to the study of view updates in both an XML and a relational setting. On the XML side, these yield an effective approach to dealing with XML view updates without relying on the limited view-update support of relational DBMSs. On the relational side, our complexity results and algorithms extend the line of research for processing relational view updates.

References

- [1] Stavros S Cosmadakis, Christos H Papadimitriou. Updates of relational views. *Journal of ACM*, 1984, 31(4): 742–760.
- [2] Umeshwar Dayal, Philip A Bernstein. On the correct translation of update operations on relational views. *ACM Transactions on Database Systems (TODS)*, 1982, 7(3): 381–416.
- [3] Arthur Keller. Algorithms for translating view updates to database updates for views involving selections, projections, and joins. In *Proc. the fourth ACM SIGACT-SIGMOD Symposium on Principles of Database Systems (PODS)*, Portland, Oregon, USA, 1985, pp.154–163.
- [4] Jens Lechtenborger, Gottfried Vossen. On the computation of relational view complements. *ACM Transactions on Database Systems (TODS)*, 2003, 28(2): 175–208.
- [5] IBM DB2 universal database SQL reference. IBM. www-306.ibm.com/software/data/db2/.
- [6] SQL reference. Oracle. www.oracle.com/technology/documentation/.
- [7] SQL server. MSDN Library. msdn.microsoft.com/en-us/sqlserver/.
- [8] Philip Bohannon, Byron Choi, Wenfei Fan. Incremental evaluation of schema-directed XML publishing. In *Proc. the 2004 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, Paris, France, 2004, pp.503–514.
- [9] Michael J Carey, Jerry Kiernan, Jayavel Shanmugasundaram, Eugene J Shekita, Subbu N Subramanian. XPERANTO: Middleware for publishing object-relational data as XML documents. In *Proc. the 26th International Conference on Very Large Data Bases (VLDB)*, Cairo, Egypt, 2000, pp.646–648.
- [10] Mary F Fernandez, Atsuyuki Morishima, Dan Suciu. Efficient evaluation of XML middleware queries. In *Proc. the 2001 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, Santa Barbara, CA, USA, 2001, pp.103–114.
- [11] Vanessa P Braganholo, Susan B Davidson, Carlos A Heuser. From XML view updates to relational view updates: Old solutions to a new problem. In *Proc. the Thirtieth International Conference on Very Large Data Bases (VLDB)*, Toronto, Canada, 2004, pp.276–287.
- [12] L Wang, E A Rundensteiner, Murali Mani. UFilter: A lightweight XML view update checker. In *Proc. the 22nd International Conference on Data Engineering (ICDE)*, Atlanta, USA, 2006, p.126.
- [13] L Wang, E A Rundensteiner, Murali Mani. Updating XML view published over relational databases: Towards the existence of a correct update mapping. *Data and Knowledge Engineering (DKE)*, 2006, 58(3): 263–298.
- [14] Laux A, Martin L. XUpdate — XML Update Language. 2000. <http://www.xmldb.org/xupdate/xupdate-wd.html>.
- [15] Sur G, Hammer J, Siméon J. An XQuery-based language for processing updates in XML. In *Proc. Programming Language Technologies for XML (PLAN-X)*, Venice, Italy, 2004.
- [16] Byron Choi. What are real DTDs like. In *Proc. the Fifth International Workshop on the Web and Databases (Webdb)*, Madison, Wisconsin, USA, 2002, pp.43–48.
- [17] Rajasekar Krishnamurthy, Raghav Kaushik, Jeffrey Naughton. XML-SQL query translation literature: The state of the art and open problems. In *Proc. Database and XML Technologies, First International XML Database Symposium (XSym)*, Berlin, Germany, 2003, pp.1–18.
- [18] Buneman P, Khanna S, Tan W. On propagation of deletions and annotations through views. In *Proc. the Twenty-First ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*, Wisconsin, USA, 2002, pp.150–158.
- [19] Byron Choi, Gao Cong, Wenfei Fan, Stratis Vigiias. Updating recursive XML views of relations. In *Proc. the 23rd International Conference on Data Engineering (ICDE)*, Istanbul, Turkey, 2007, pp.766–775.
- [20] Michael Benedikt, Chee Yong Chan, Weifei Fan, Rajeev Rastogi, Shihui Zheng, Aoying Zhou. DTD-directed publishing with attribute translation grammars. In *Proc. the 28th International Conference on Very Large Data Bases (VLDB)*, Hong Kong, China, 2002, pp.838–849.
- [21] Jayavel Shanmugasundaram, Kristin Tuftte, Chun Zhang, Gang He, David J DeWitt, Jeffrey F Naughton. Relational databases for querying XML documents: Limitations and opportunities. In *Proc. 25th International Conference on Very Large Data Bases (VLDB)*, Edinburgh, Scotland, UK, 1999, pp.302–314.
- [22] Cormen T H, Leiserson C E, Rivest R L, Stein C. Introduction to Algorithms. McGraw-Hill, 2001.
- [23] Li Chen, Amarnath Gupta, M Erdem Kurul. Stack-based algorithms for pattern matching on DAGs. In *Proc. the 31st International Conference on Very Large Data Bases (VLDB)*, Trondheim, Norway, 2005, pp.493–504.
- [24] Ralf Schenkel, Anja Theobald, Gerhard Weikum. Efficient creation and incremental maintenance of the HOPI index for complex XML document collections. In *Proc. the 21st International Conference on Data Engineering (ICDE)*, Tokyo, Japan, 2005, pp.360–371.
- [25] Christoph Koch. Efficient processing of expressive node-selecting queries on XML data in secondary storage: A tree automata-based approach. In *Proc. the 29th International Conference on Very Large Data Bases (VLDB)*, Berlin, Germany, 2003, pp.249–260.
- [26] Italiano G F. Finding paths and deleting edges in directed acyclic graphs. *Inf. Process. Lett.*, 1988, 28.
- [27] King V, Sagert G. A fully dynamic algorithm for maintaining the transitive closure. In *Proc. ACM Symposium on Theory of Computing*, 1999.
- [28] Alberto Marchetti-Spaccamela, Umberto Nanni, Hans Rohnert. Maintaining a topological order under edge insertions. *Information Processing Letters*, 1996, 59(1): 53–58.
- [29] Michael Garey, David Johnson. Computers and Intractability: A Guide to the Theory of NP-Completeness. 1979.
- [30] Bart Selman, Henry Kautz. Walksat Home Page. 2004. <http://www.cs.washington.edu/homes/kautz/walksat/>.
- [31] Elias Koutsoupias, Christos H Papadimitriou. On the greedy algorithm for satisfiability. *Information Processing Letters*, 1992, 43(1): 53–55.

- [32] Wang L, Mulchandani M, Rundensteiner E. Updating XQuery views published over relational data: A round-trip case study. In *Proc. XML Database Symposium*, 2003, pp.223–237.
- [33] Wang L, Rundensteiner E A. Updating XML views published over relational databases: Towards the existence of a correct update mapping. Technical Report, Worcester Polytechnic Institute, 2004.
- [34] Yingwei Cui, Jennifer Widom. Run-time translation of view tuple deletions using data lineage. Technical Report, Stanford University, 2001.
- [35] Gao Cong, Wenfei Fan, Floris Geerts. Annotation propagation revisited for key preserving views. In *Proc. the 15th ACM International Conference on Information and Knowledge Management*, Arlington, Virginia, USA, 2006, pp.632–641.
- [36] Cohen E, Kaplan H, Milo T. Labeling dynamic XML tree. In *Proc. the Twenty-First ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*, Madison, Wisconsin, USA, 2002, pp.271–281.



Byron Choi received his B.Eng. degree in computer engineering from the Hong Kong University of Science and Technology (HKUST) in 1999 and his M.SE. and Ph.D. degrees in computer and information science from the University of Pennsylvania in 2002 and 2006, respectively. He has been an assistant professor in the School of Computer Engineering,

Nanyang Technology University, Singapore, in 2005~2008. He joins Hong Kong Baptist University as an assistant professor in 2008. His research interests include XML, query processing and optimization.



Gao Cong is an assistant professor in Aalborg University, Denmark. His research interests include text and data mining, web data management, and database. He has published papers in conferences, such as SIGMOD, VLDB, KDD, ICDE, SIGIR, ACL, AAAI, etc. He also served on the programme committee of some of the aforementioned conferences.



Wenfei Fan is a professor (Chair of Web Data Management) in the School of Informatics, University of Edinburgh, and a Research Scientist at Bell Laboratories, Alcatel-Lucent. He received his Ph.D. degree from the University of Pennsylvania, and his M.S. and B.S. degrees from Peking University. He is the recipient of the Roger Needham Award in

2008, the Chang Jiang Scholar Award in 2007, the Outstanding Overseas Young Scholar Award in 2003, the Career Award in 2001, the ICDE Best Paper Award in 2007,

and the Best Paper of the Year Award from Computer Networks in 2002. His current research interests include data quality, data integration, integrity constraints, distributed query processing, Web services and XML.



Stratis D. Viglas is a lecturer in database systems at the School of Informatics of the University of Edinburgh. he holds a Bachelor's and a Master's degree in computer science from the University of Athens, and a Ph.D. degree in computer science from the University of Wisconsin-Madison. He joined the School of Informatics at the University of Edinburgh in 2003. His research interests lie in the areas of database systems and, in particular, on query optimization and evaluation, and data storage and management for both centralized and distributed systems.

His research interests lie in the areas of database systems and, in particular, on query optimization and evaluation, and data storage and management for both centralized and distributed systems.

Appendix Query Evaluation on Database with Variables

Given the original database I_i ($i = 1, \dots, n$), the set of relational tuples to be inserted X_i ($i = 1, \dots, n$) and the conjunctive query $Q = \pi_P(\sigma_C(T_1, \dots, T_n))$, where C is a conjunction of equalities and P is a set of projected attributes, the problem is how to evaluate query Q on database I_i incremented by X_i that contains variables to capture whether insertions X_i will yield side effects. The challenge here is that the selection conditions of Q cannot be evaluated on tuples with variables and thus SQL queries cannot work directly on tuples enriched with variables.

Before analyzing how side-effects are generated and discussing how to evaluate Q to capture side-effects, we will do some preprocessing in order to 1) guarantee that Δ_R can be generated from the conjunctive query (view) on $I_i \cup X_i$ for any instantiation of the variables in X_i ; and 2) reduce the number of variables. The preprocessing consists of several steps: 1) if there is a selection condition such that $z_{ik} = z_{jl}$, $z_{ik} \in \mathbf{z}_i$, $z_{jl} \in \mathbf{z}_j$, we use one variable to rename z_{ik} and z_{jl} ; 2) if a variable is not involved in selection conditions, it can be filled with a dummy value because the instantiation of the variable is not relevant to side-effects; and (3) If there already exists a base tuple r' sharing key with r in X_i , we fill the missing values in r according to r' .

We observe that there are only two types of side-effects.

1) A view tuple is a side effect if it contains at least one key from $I_i \setminus X_i$ and at least one key from $X_j \setminus I_j$.

2) A view tuple is a side effect if it is generated from X_i ($i = 1, \dots, n$), but is not a tuple in

$\Delta_R \cup Q(I_1 \cap X_1, \dots, I_n \cap X_n)$.

The above two kinds of side effects cover all possible side effects raised by the insertion of Δ_R while other possibility, such as $Q(I_1, \dots, I_n)$, will not generate any side effect tuples. For convenience of presentation, we divide $I_i \cup X_i$ into three non-overlapping subsets for each $i \in [1, n]$:

- $U_i = X_i \setminus I_i, i \in [1, n]$,
- $A_i = I_i \setminus X_i, i \in [1, n]$,
- $B_i = X_i \cap I_i, i \in [1, n]$.

To capture the first kind of side effect, for all possibilities of T_1, \dots, T_n , where $T_i \in \{U_i, A_i, B_i\}$, such that there exist an $i, j \in [1, n]$, $T_i = U_i$ and $T_j = A_i$, we rewrite Q to accommodate the variables in U_i and thus to capture side effects. More specifically, we rewrite the selection conditions and projected attributes. We illustrate the rewriting using an example: given $Q := \pi_P(\sigma_C(R_1, R_2, R_3))$ and one combination (U_1, U_2, A_3) , to capture the side effects from the combination we rewrite the Q into $Q' = \pi_{P_1}(\sigma_{C_2}(U_1, U_2, A_3))$. The selection conditions C in Q are decomposed into C_1 and C_2 , where C_1 only contains equality conditions involving variables (must be in U_1 and U_2 in this example) while C_2 contains the other selection conditions. P_1 contains only the attributes contained in C_2 . Observe that 1) the selection conditions in C_2 that do not contain variable can be imposed on Q' , and 2) the projection on P_1 ensures that any two of generated side effect tuples produce different encoding. The second kind of side effect is captured similarly.

The algorithm is given in Fig.A. Its input consists of 1) a set of base relations $\{I_1, \dots, I_n\}$, 2) a view V defined in terms of conjunctive query $V = \pi_P(\sigma_C(R_1 \times \dots \times R_n))$, and 3) a group insertion $\Delta_R = \{t_1, \dots, t_k\}$ against V . The first kind of side-effect is encoded in lines 7~16. If a returned tuple does not contain any variable, it is a side-effect tuple (line 13); if it contains some variables, we need to instantiate the variables such that the selection conditions in C_1 are not satisfied in order to avoid side-effect. More specifically, for each return tuple t_k containing variable, we construct for

each condition c_j in C_1 one inequality $x_{k_j} \neq z_{k_j}$, where x_{k_j} is a variable and z_{k_j} can be either a constant or a variable. Side-effect tuple t_k can be avoided only if at least one of the above inequalities holds. Similarly, we encode the second kind of side effect (lines 17~25).

<p>Input: relations I_1, \dots, I_n, view V, a group insertion Δ_R, the view definition $\pi_P(\sigma_C(R_1 \times \dots \times R_n))$, Output: side-effect encode or reject (exception)</p> <ol style="list-style-type: none"> 1. Compute X_i from Δ_R w.r.t. R_i, for $i \in [1, n]$; 2. Preprocess X_i; 3. $\Theta := \emptyset$ /* SAT instance */ 4. $U_i := X_i \setminus I_i, i \in [1, n]$ 5. $A_i := I_i \setminus X_i, i \in [1, n]$ 6. $B_i := X_i \cap I_i, i \in [1, n]$ <p>/* detect the first type of side-effect */</p> <ol style="list-style-type: none"> 7. for each combination of T_1, \dots, T_n, s.t., $\exists i \exists j [T_i = U_i \wedge T_j = A_j], \wedge \forall k [(k \neq i \wedge k \neq j) \rightarrow (T_k = U_k \vee T_k = R_k)]$ 8. $C_1 :=$ selection conditions involving variables in T_i 9. $C_2 := C \setminus C_1$ 10. $P_1 :=$ attributes involved in conditions in C_1 11. $\Delta V_1 := \pi_{P_1}(\sigma_{C_2}(T_1, \dots, T_k))$ 12. for each $t' \in \Delta V_1$ 13. if t' does not contain variable then reject Δ_R 14. return 15. else $\Theta := \Theta \wedge (\bigvee_{c_j \in C_1} ((x_{k_j} \neq z_{k_j})))$ 16. endfor <p>/* detect the second type of side-effect */</p> <ol style="list-style-type: none"> 17. for each combination of T_1, \dots, T_n, s.t., $\exists i [T_i = U_i] \wedge \forall k [(k \neq i) \rightarrow (T_k = X_k)]$ 18. $C_1 :=$ selection conditions involving variables in T_i 19. $C_2 := C \setminus C_1$ 20. $\Delta V_2 := \sigma_{C_2}(T_1, \dots, T_k)$ 21. for each $t' \in \Delta V_2 \wedge t' \notin U$ 22. if t' does not contain variable then reject Δ_R 23. return 24. else $\Theta := \Theta \wedge (\bigvee_{c_j \in C_1} (x_{k_j} \neq z_{k_j}))$ 25. endfor 26. return Θ
--

Fig.A. Algorithm insert.