

# Designing Information-Preserving Mapping Schemes for XML

Denilson Barbosa  
University of Calgary  
denilson@scpsc.ucalgary.ca

Juliana Freire  
University of Utah  
juliana@cs.utah.edu

Alberto O. Mendelzon  
University of Toronto  
mendel@cs.toronto.edu

## Abstract

An XML-to-relational mapping scheme consists of a procedure for *shredding* XML documents into relational databases, a procedure for *publishing* databases back as documents, and a set of constraints the databases must satisfy. In previous work, we discussed two notions of information preservation for mapping schemes: *losslessness*, which guarantees the complete reconstruction of a document from a database; and *validation*, which guarantees that every update to a database corresponding to a valid document results in a database corresponding to another valid document. Also, we described one information preserving mapping scheme, called  $\text{Edge}^{++}$ , and showed that, under reasonable assumptions, lossless and validation are both undecidable. This leads to the question we study in this paper: how to design information-preserving mapping schemes. We propose to do it by starting with a scheme known to be information preserving (such as  $\text{Edge}^{++}$ ) and applying to it equivalence-preserving transformations written in weakly recursive ILOG. We study a particular incarnation of this framework, the LILO algorithm, and show that it provides significant performance improvements over  $\text{Edge}^{++}$  and that the constraints it introduces are efficiently enforced in practice.

## 1 Introduction

In order to use relational engines for managing XML data, we need a *mapping scheme* consisting of a procedure for shredding XML documents into relational databases, a procedure for publishing those databases as documents, and a set of constraints that those databases must satisfy. As with any other mapping strategy, it is important to study the information preservation properties of such a scheme in order to understand its suitability for a given application [25]. Although there is a rich literature on mapping schemes [5, 13, 16, 21, 29, 22, 31], little attention

---

*Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.*

**Proceedings of the 31st VLDB Conference,  
Trondheim, Norway, 2005**

has been given to their *correctness*; i.e., whether they preserve *enough* information. In previous work [3], we defined *lossless* mapping schemes as those that allow the reconstruction of the original documents, and *validating* mapping schemes as those in which all legal relational databases correspond to a valid XML document. We argued that while losslessness is enough for applications involving only queries over the documents, both losslessness *and* validation are required if the documents must conform to an XML schema and the application involves both queries and updates to the documents. We also described a mapping scheme,  $\text{Edge}^{++}$ , in which both losslessness and validation are guaranteed by constraints in its relational schema.

In this paper, we address the problem of designing information-preserving mapping schemes from an XML schema. Previous mapping design algorithms [5, 29] do not guarantee information preservation. Moreover, since both losslessness and validation are undecidable for a large class of mapping schemes that includes all those in the literature [3], arbitrary design procedures cannot guarantee information preservation in general. We propose a sound framework for designing information-preserving mapping schemes that can serve as the basis for design algorithms: start with an information preserving mapping scheme, and repeatedly apply *equivalence-preserving* transformations [25] to it. Following this procedure, information preservation is guaranteed *by construction*. Our framework is extensible and allows *any* transformation that can be written in weakly recursive ILOG with stratified negation [19] ( $\text{wrec-ILOG}^-$ ), which is powerful enough for expressing most of transformations proposed in the literature (e.g., [5, 28]). We also discuss an instance of our framework: the LILO (for Lossless Inlining, Lossless Outlining) algorithm, which uses  $\text{Edge}^{++}$  as starting point and defines several equivalence preserving transformations (some of which are extensions of transformations in the literature). Our experimental results show that LILO results in mapping schemes that outperform the previous  $\text{Edge}^{++}$  substantially.

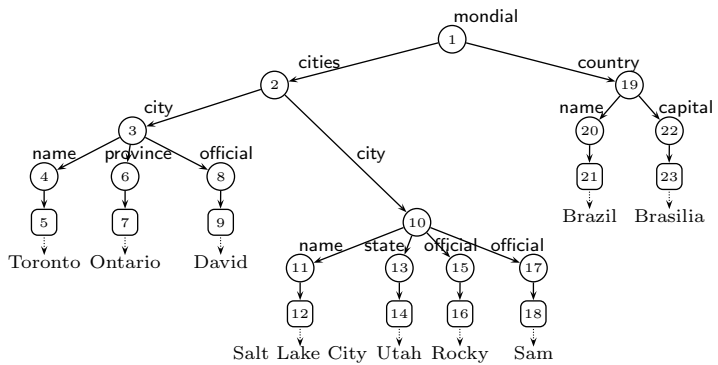
**Information Preservation in Mapping Schemes.** The example below illustrates inconsistencies that arise when updates are considered in a mapping scheme that is not validating.

$\Sigma = \{mondial, cities, country, city, province, state, name, official\}$

root = mondial

$R = \{mondial \leftarrow cities, country^*; cities \leftarrow city^*; city \leftarrow name, (province|state), official[1, 5]; country \leftarrow name, capital; name \leftarrow \#PCDATA; province \leftarrow \#PCDATA; state \leftarrow \#PCDATA; official \leftarrow \#PCDATA; capital \leftarrow \#PCDATA\}$

(a)



(b)

Figure 1: (a) XML document; elements are shown as circles and text nodes are shown as boxes. The labels of elements and the values of the text nodes are also shown. (b) a DTD for the document.

**Example 1.** Consider the schema in Figure 1(a), inspired by the Mondial Database<sup>1</sup>, describing cities and countries. For each city, the database contains its name; the name of the province or state the city is located; and from 1 to 5 government officials. Countries are described by their name and the name of their capitals. A typical relational schema derived from applying inlining techniques to the document schema above is as follows<sup>2</sup>:

```
city(cityId,name,provinceId,stateId)
official(officialId,cityId,name)
country(countryId,name,capital)
```

For the document in Figure 1(b), we would have the following database instance under this mapping:

```
city(1,'Toronto','Ontario',NULL)
city(4,'Salt Lake City',NULL,'Utah')
official(2,1,'David')
official(5,4,'Rocky')
official(6,4,'Sam')
capital(7,'Brazil','Brasilia')
```

This mapping scheme is lossless as the original document can be reconstructed from the database. However, by applying the following *legal* update:

```
UPDATE city SET province='Utah'
WHERE name='Salt Lake City'
```

we arrive at database that no longer corresponds to a valid document. This anomaly is due to the fact that the mapping scheme does not preserve the semantics of the choice (|) construct, i.e., that cities can have either a province or state, but not both. ■

In a lossless *and* validating mapping scheme, the relational schema is *equivalent* to the document schema. Thus, constraints in the relational schema prevent such anomalies by disallowing SQL updates that result in databases corresponding to invalid documents. The alternative to using an information-preserving mapping scheme would be re-validating the XML document re-

sulting from the update before changes are committed. However, note that testing whether even simple updates such as<sup>3</sup>:

```
update insert <official>John</official>
```

```
into /mondial/cities/city[name='Salt Lake City']
```

results in a valid document cannot be done statically (in this case, the legality of the update depends on the number of officials already associated with the city). While one could test whether an update is permissible by reconstructing the elements involved, applying the update to those elements, and using a *validator* to test whether the result conforms to the document schema, this approach has several disadvantages. First, validation is computationally expensive [26]; also, incremental validation techniques [4, 27] require and maintain auxiliary information that must be kept synchronized with the database; finally, this approach requires developing and maintaining a special-purpose application, and ignores the DBMS infrastructure for checking constraints.

**Outline and Contributions.** We argue in this paper that it is feasible to augment the relational schemas in mapping schemes with constraints for ensuring the *validity* of the elements in the XML documents with respect to the content models in an XML schema. We start in Section 2 discussing element and document validity, XML mapping schemes and information preservation. We propose a sound and extensible framework for designing information preserving mapping schemes (Section 3), which consists of applying equivalence preserving transformations to an information-preserving mapping scheme. We show how a mapping scheme and an arbitrary transformation written in  $wrec\text{-}ILOG^{\neg}$  can be rewritten as another mapping scheme (Section 3.3) in a mechanical way. We discuss several equivalence preserving transformations that result in mapping schemes defining simpler relational constraints for ensuring element validity, and introduce the LILO algorithm in Section 4. We show that

<sup>1</sup><http://dbis.informatik.uni-goettingen.de/Mondial>

<sup>2</sup>Primary keys are underlined and *nullable* columns are shown in italics

<sup>3</sup>Using the syntax of [23]

LILO provides significant performance improvements over Edge<sup>++</sup> and that the constraints it introduces can be efficiently enforced in practice (Section 5). We conclude with a discussion in Section 6.

## 2 Definitions and Terminology

XML documents are modeled by ordered labeled trees whose nodes represent either elements or textual content, thus capturing the essential data representation components of XML [7]. We use  $\tau$ ,  $\lambda$  and  $\nu$  to denote the the *type*, the *label* and the *value* of nodes in the tree, respectively. More precisely, let  $\mathcal{I}$ ,  $\mathcal{D}$  be two disjoint, countably infinite sets of node ids and values:

**Definition 1.** *An XML Document is a tuple  $\langle T, \lambda, \tau, \nu \rangle$ , where  $T$  is an ordered tree whose nodes are elements of  $\mathcal{I}$ ;  $\tau : \mathcal{I} \rightarrow \{\text{element}, \text{text}\}$  assigns types to nodes in  $T$ , such that  $\tau(r) = \text{element}$  if  $r$  is the root of  $T$ , and if  $\tau(e) = \text{text}$  then  $e$  is a leaf in  $T$ ;  $\lambda : \mathcal{I} \rightarrow \mathcal{D}$  assigns labels to nodes in  $T$  such that the label of all text nodes is  $\#PCDATA$ ; and  $\nu : \mathcal{I} \rightarrow \mathcal{D}$  assigns values to text nodes in  $T$ .*

For brevity, we will refer to subtrees rooted at *element* nodes simply as “elements”. We denote by  $\mathcal{X}$  the set of all XML documents.

The focus of this work is on preserving the validity constraint of elements, which is *identical* for both DTDs and XML Schemas: an element of a given *type* is said to be valid if its content is a word in a 1-unambiguous regular language [8] associated with that type. The only difference between these formalisms is the way that types are associated to elements: DTDs restrict all elements with the same label to have the same type, while XML Schema allows assigning different types to the same label, depending on the *context* in which the label appears. While our method handles XML Schemas and DTDs seamlessly, we restrict our discussion to DTDs for simplicity, and refer the reader to [3, 4] for handling type specialization.

**Definition 2.** *An XML DTD is a triple  $\langle \Sigma, r, \mathbf{R} \rangle$  where  $\Sigma$  is a set of element labels,  $r \in \Sigma$  is a distinguished label and  $\mathbf{R}$  is a mapping associating to each  $a \in \Sigma$  a content model expressed as a 1-unambiguous regular expression over  $\Sigma \cup \{\#PCDATA\}$ .*

Figure 1(a) shows a DTD; cardinality constraints of the form  $a[x, y]$  are shorthands for  $x$  copies of  $a$  followed by  $y$  copies of  $a$ ?

The validity of a document  $D$  with respect to a DTD  $X = \langle \Sigma, r, \mathbf{R} \rangle$  is done as follows. Let  $e$  be an element and  $c_1, \dots, c_n$  be its ordered children; the *content* of  $e$  is the string  $\lambda(c_1) \cdots \lambda(c_n)$ . We say that  $e$  is *valid* with respect to  $X$  if its content matches the regular expression associated with  $\lambda(e)$  in  $\mathbf{R}$ . A document  $D$  is valid with respect to  $X$ , written  $D \in L(X)$ , if all of its elements are valid with respect to  $X$  and the label of its root element is  $r$ . Using our notation,

the validation problem can be re-stated as: given  $D$  and  $X$ , is it the case that  $D \in L(X)$ ?

### 2.1 XML-to-relational Mapping Schemes

We model relational databases with separate domains for *surrogates* to element ids (e.g., `cityId` in Example 1) and other constants. Let  $\mathcal{J}$ ,  $\mathcal{D}$  be two disjoint countably infinite sets of surrogates and constants, respectively. A *relational* schema is a set of relation schemes and constraints; each relation scheme  $R$  has a set (possibly empty) of attributes with domain  $\mathcal{J}$ —called the *surrogate attributes* of  $R$ , and a set (possibly empty) of attributes with domain  $\mathcal{D}$ . Instances are defined as customary [1, 24]. A *constraint* is expressed as a boolean query and is said to be violated if that query evaluates to true. A relational database instance is *legal* if it does not violate any constraint in its schema.  $\mathcal{R}(S)$  denotes the set of legal instances of  $S$ .

No meaning is assigned to node ids the document tree, nor to the surrogates used for representing them; furthermore, no relationship between node ids and surrogates is assume either. That is, renaming node ids in Figure 1(b) does not yield a new document; similarly, renaming surrogates in the database in Example 1 does not create a new database. These properties are captured as follows:

**Definition 3.** *XML documents  $D_1 = \langle T_1, \lambda_1, \tau_1, \nu_1 \rangle$ , and  $D_2 = \langle T_2, \lambda_2, \tau_2, \nu_2 \rangle$ , are equivalent, denoted by  $D_1 \equiv_X D_2$ , if there exists an isomorphism  $\phi : \mathcal{I} \rightarrow \mathcal{I}$  between  $T_1$  and  $T_2$  such that  $\lambda_1(v) = \lambda_2(\phi(v))$ ,  $\tau_1(v) = \tau_2(\phi(v))$ , and  $\nu_1(v) = \nu_2(\phi(v))$ , for all  $v \in T_1$ .*

**Definition 4.** *Database instances  $I_1, I_2$  are equivalent, written  $I_1 \equiv_R I_2$  if there is a bijection on  $\mathcal{J} \cup \mathcal{D}$  that maps  $\mathcal{J}$  to  $\mathcal{J}$ , is the identity on  $\mathcal{D}$ , and transforms  $I_1$  into  $I_2$ .*

The notion of database equivalence above has been called *OID equivalence* in object databases [1].  $[D]$  denotes the equivalence class of document  $D$ ; that is  $[D] = \{D' \in \mathcal{X} \mid D' \equiv_X D\}$ ; similarly,  $[I]$  denotes the equivalence class of database  $I$ .

An XML-to-relational *mapping scheme* [3] as a triple  $\mu = (\sigma, \pi, S)$ , where  $S$  is a relational schema;  $\sigma$  is a *mapping function* that assigns instances of  $S$  to XML documents; and  $\pi$  is a *publishing function* that assigns XML documents to instances of  $S$ .

**Definition 5.** *An XML-to-relational mapping scheme is a triple  $\mu = (\sigma, \pi, S)$ , where  $\sigma : \mathcal{X} \rightarrow \mathcal{R}(S)$  is a partial function, and  $\pi : \mathcal{R}(S) \rightarrow \mathcal{X}$  is a total function, and the following hold: (1) for all  $D_1, D_2 \in \mathcal{X}$  we have that  $D_1 \equiv_X D_2$  implies  $\sigma(D_1) \equiv_R \sigma(D_2)$ ; and (2) for all  $I_1, I_2 \in \mathcal{R}(S)$  we have that  $I_1 \equiv_R I_2$  implies  $\pi(I_1) \equiv_X \pi(I_2)$ .*

Defining  $\sigma$  as a partial function accommodates mapping schemes customized for a specific DTD (e.g., [5, 29]); on the other hand, defining  $\pi$  to be total ensures

that any legal database *represents* (i.e., can be published as) a document. Conditions (1) and (2) ensure that both  $\sigma$  and  $\pi$  are *generic*: they map equivalent documents to equivalent databases and vice-versa.

## 2.2 The $\mathcal{XDS}$ Class of Mapping Schemes

A *class* of mapping schemes is defined in terms of the languages used for specifying  $\sigma$ ,  $S$ , and  $\pi$ . The power of these languages determines what kinds of mappings can be specified. The  $\mathcal{XDS}$  class of mapping schemes [3] is defined as follows.

**The Mapping Language.** The language consists of XQuery augmented with a clause `sql ... end` for specifying SQL insert statements, and to be used instead of the `return` clause in a FLOWR expression. The semantics of the mapping expressions is defined similarly to the usual semantics of FLOWR expressions: the `for`, `let`, `where`, `order by` clauses define a list of tuples which are passed, one at a time, to the `sql ... end` clause, and one SQL transaction is issued per such tuple.

**The Constraint Language.** Constraints are boolean programs in Datalog with stratified negation [1]. This language allows easy expression of standard relational constraints (e.g., functional dependencies and referential integrity), as well as graph connectivity—required for ensuring the database encodes a tree, and element validity (Section 3.2).

**The Publishing Language.** Publishing functions are arbitrary XQuery expressions over a “canonical” XML view of a relational database. That is, each relation is mapped into an element whose children represent the tuples in that relation in the standard way (i.e., one element per column). This is the approach taken by SilkRoute [15] and XPERANTO [9].

$\mathcal{XDS}$  is, by design, a powerful mapping tool. In fact, all mapping schemes that we are aware of in the literature can be expressed using it.

## 2.3 Information Preservation in $\mathcal{XDS}$

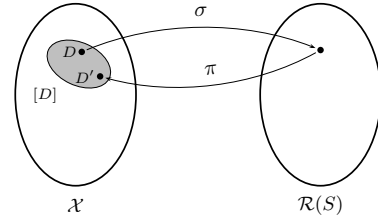
For completeness, we briefly revisit losslessness and validation here, and refer the reader to [3] for details. A mapping scheme is lossless if it allows the complete reconstruction of any (fragment of a) document from its database (see Figure 2(a)):

**Definition 6.** A mapping scheme  $\mu = (\sigma, \pi, S)$  is lossless if for all  $D \in \text{Dom}(\sigma)$ ,  $\pi(\sigma(D)) \equiv_X D$ .

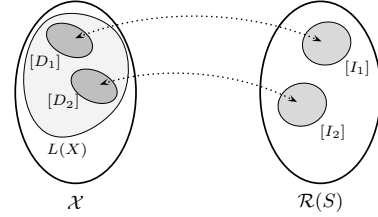
The following are easy to verify:

**Proposition 1.**  $\mu = (\sigma, \pi, S)$  is lossless if and only if  $\pi(\sigma(\cdot))$  is the identity on equivalence classes in  $\text{Dom}(\sigma)$ .

Validation is defined in terms of a DTD  $X$ . A validating mapping scheme is one in which every legal database instance corresponds to a document in  $L(X)$ :



(a) Lossless mapping scheme.



(b) Lossless and validating mapping scheme.

Figure 2: Information preservation in mapping schemes.

**Definition 7.** A mapping scheme  $\mu = (\sigma, \pi, S)$  is validating with respect to DTD  $X$  if  $\sigma$  is total on  $L(X)$ , and for all  $I \in \mathcal{R}(S)$ , there exists  $D \in L(X)$  such that  $I \equiv_R \sigma(D)$ .

Note that for a validating  $\mu = (\sigma, \pi, S)$ , every *successful* update on an instance of  $S$  results in a database that represents a valid document; it follows that only permissible updates over the original document can be effected on its corresponding relational database. As discussed in Section 1, when  $\mu$  is lossless only, testing if an update is permissible can be done by materializing  $\pi(\sigma(D))$ , effecting the update, and validating the resulting document, which can be prohibitively expensive in most cases.

As discussed in [3], losslessness and validation are orthogonal and one does not imply the other. Applications that define a DTD, and involve both querying and updating the XML documents require mapping schemes that have both properties. Note that when this is the case, it is guaranteed that all queries and updates over the documents can be done using SQL and that every database corresponds to a valid document:

**Proposition 2.**  $\mu = (\sigma, \pi, S)$  is lossless and validating with respect to DTD  $X$  if and only if  $\sigma$  and  $\pi$  are bijective, and  $\pi$  is the inverse of  $\sigma$  (up to equivalence).

*Proof.* It is easy to see that  $\mu$  is both lossless and validating if  $\sigma$  and  $\pi$  are as above. For the other implication, note that since  $\mu$  is validating w.r.t.  $X$  and both  $\sigma$  and  $\pi$  are generic (Definition 5), it follows that  $\sigma$  defines a bijection between equivalence classes of document in  $L(X)$  and database instances in  $\mathcal{R}(S)$ . A similar argument applies for  $\pi$ . Since  $\mu$  is also lossless,

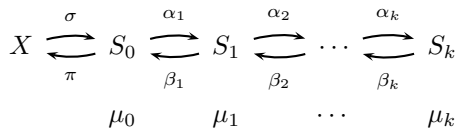


Figure 3: Designing of an information-preserving mapping scheme:  $X$  is a DTD,  $\mu_0$  is the  $\text{Edge}^{++}$  mapping scheme, and each  $\mu_i$  is the result of applying an information-preserving schema transformation.

Proposition 1 implies that  $\sigma$  and  $\pi$  are the inverse of each other up to equivalence.  $\square$

Put in other words, a lossless and validating mapping scheme defines a bijection among equivalence classes of valid documents and legal relational databases, as depicted in Figure 2(b). We call mapping schemes that are both lossless and validating (with respect to a DTD) *information-preserving*.

### 3 Designing Information-Preserving Mapping Schemes

We now discuss a general and sound framework for designing information-preserving mapping schemes based on applying *structural transformations* [18, 25] to an existing schema. In summary, we start with an initial mapping scheme  $\mu_0$  that is known to be information-preserving, and subsequently transform it until the desired one  $\mu_k$ , defined as

$$\mu_k = ((\alpha_k \circ \alpha_{k-1} \circ \dots \circ \alpha_1 \circ \sigma), (\pi \circ \beta_1 \circ \dots \circ \beta_k), S_k)$$

is found, as illustrated by Figure 3. To make this framework concrete we must specify the languages for expressing the (1) mappings between XML documents and database instances, as well as the (2) mappings  $\alpha_i, \beta_i$  between instances of different relational schemas. In our case, the mapping and publishing languages are those in  $\mathcal{XDS}$  (Section 2.2), and the mappings between relational instances are given as “weakly recursive” ILOG programs with stratified negation (wrec-ILOG<sup>-</sup>) [19]. Informally, this language corresponds to Datalog with stratified negation augmented with non-recursive “invention” rules that create new surrogates, and is powerful enough for the kinds of transformations we define.

**Soundness.** Information preservation in our framework is achieved as follows. First, we always start with a mapping scheme  $\mu_0 = (\sigma_0, \pi_0, S_0)$  that is both lossless and validating with respect to a DTD  $X$ ; as we discuss below, this is the case if and only if  $X$  and  $S_0$  are *equivalent* (Proposition 3). Second, *every* transformation applied to a mapping scheme  $\mu_i, 0 < i < k$  results in  $\mu_{i+1}$  whose relational schema is *equivalent* to that of  $\mu_i$ . Such transformations are called *equivalence preserving* [25]. It follows that the resulting mapping scheme  $\mu_k = (\sigma_k, \pi_k, S_k)$  is such that  $S_k$  is equivalent to  $X$  and thus information preserving.

In the remainder of this section, we relate classical notions of information preservation and relative capacity of schemas [18, 25] to our notions of lossless and

validation [3]. Then, we briefly describe the  $\text{Edge}^{++}$  mapping scheme, which is the starting point of LILO, discussed in Section 4. Finally, we give a procedure for incorporating arbitrary wrec-ILOG<sup>-</sup> programs into  $\mathcal{XDS}$  mapping and publishing functions (Section 3.3), thus making our framework extensible: *any* equivalence preserving transformation that can be specified in wrec-ILOG<sup>-</sup> can be used with our method.

#### 3.1 Dominance and Equivalence of Schemas

Equivalence of schemas  $S$  and  $T$  (within or across data models) is defined based on properties of the mappings between their instances [25]. Let  $\mathbf{I}(S)$  denote the set of all instances of  $S$  (e.g., valid documents if  $S$  is a DTD or legal database instances if  $S$  is a relational schema).

Let  $\alpha : \mathbf{I}(S) \rightarrow \mathbf{I}(T)$  be a mapping given in some appropriate language.  $\alpha$  is said to be *information-preserving* if it is *reversible*; that is, there exists  $\beta : \mathbf{I}(T) \rightarrow \mathbf{I}(S)$  such that  $\beta(\alpha(\cdot))$  is an equivalence relation in  $\mathbf{I}(S)$ . If this is the case, we say that  $T$  *dominates*  $S$  (via  $\alpha, \beta$ ), denoted  $S \preceq T$ . (We note that classical notions of schema dominance [18, 25] require  $\beta(\alpha(\cdot))$  to be the identity on  $\mathbf{I}(T)$ ; without loss of generality, we use an equivalence relation to account for renamings of element surrogates.)

If both  $\alpha$  and  $\beta$  are total and bijective, then  $S \preceq T$  via  $(\alpha, \beta)$  and  $T \preceq S$  (via  $(\beta, \alpha)$ ). In this case, we say that  $S$  and  $T$  are *equivalent*, denoted  $S \equiv T$ , and that  $\beta$  (resp.  $\alpha$ ) is an *equivalence preserving* transformation. Note that  $\preceq$  is a transitive relation, while  $\equiv$  is an equivalence relation.

Schema dominance and equivalence establish notions of “relative information capacity” between schemas. If  $S \preceq T$ , we say that  $T$  has *at least the information capacity of*  $S$ , since instances of  $S$  can represent any instance of  $T$ . Similarly, if  $S \equiv T$ , we say  $S$  and  $T$  have the *same information capacity*.

#### Equivalence of DTDs and Relational Schemas.

In a mapping scheme  $\mu = (\sigma, \pi, S)$ ,  $\sigma$  and  $\pi$  play the role of  $\alpha$  and  $\beta$  above, respectively:

**Proposition 3.** *If  $X$  is a DTD and  $\mu = (\sigma, \pi, S)$  is a mapping scheme such that  $L(X) \subseteq \text{Dom}(\sigma)$  then: (1)  $\mu$  is lossless if and only if  $X \preceq S$ ; (2)  $\mu$  is both lossless and validating with respect to  $X$  if and only if  $X \equiv S$ .*

*Proof.* (1) follows from Proposition 1, the fact that  $L(X) \subseteq \text{Dom}(\sigma)$ , and the definition of  $\preceq$ . (2) follows from Proposition 2 and the definition of  $\equiv$ .  $\square$

Since losslessness and validation are undecidable [3], it follows that:

**Corollary 1.** *Equivalence of DTDs and relational schemas is undecidable.*

#### 3.2 The $\text{Edge}^{++}$ Mapping Scheme

In what follows, we describe  $\text{Edge}^{++}$ , the initial mapping scheme in our framework and show it is information-preserving.  $\text{Edge}^{++}$  extends the  $\text{Edge}$



that all invention rules occur *after* all non-invention rules. We note that for every wrec-ILOG<sup>-</sup> program there is an equivalent one in normal form [19].

**Defining  $\sigma'$ .** Without loss of generality, we assume that each mapping expression in  $\sigma$  populates a single relation in  $S$ , and that each relation in  $S$  is populated by a single mapping expression in  $\sigma$ . (We can replace each procedure  $p$  populating relations  $r_1, \dots, r_n$  by  $n$  “copies” of  $p$  such that copy  $p_i$  populates relation  $r_i$  only. Also, we can replace  $n$  mapping expressions  $p_1, \dots, p_n$  where each populate relation  $r$  by a new procedure  $p'$  and  $n$  functions  $f_1, \dots, f_n$ , such that each  $f_i$  return the list of tuples inserted by  $p_i$  and  $p'$  inserts into  $r$  the “union” of the results of  $f_1, \dots, f_n$ .)

The first step of  $\sigma'$  materializes an “instance” of  $S$  (which we refer to as  $I_S$ ) as follows. Let  $R_1, \dots, R_n$  be the relation names in  $S$  and  $p_1, \dots, p_n$  be the mapping expressions in  $\sigma$  that populate them; we convert each  $p_i$  into a function  $f_i$  that returns the sequence of “tuples” that would be inserted in  $R_i$ . The result of each  $f_i$  is sorted lexicographically and duplicate tuples are removed. For the second step of  $\sigma'$ , we convert  $\alpha$  into a mapping expression that takes  $I_S$  and produces the desired instance of  $S'$ . Let  $r_1, \dots, r_i, \dots, r_n$  be the rules in  $\alpha$  stratified in a way that all invention rules are  $r_{i+1}, \dots, r_n$ . We define  $\alpha$  as two standard XQuery programs  $f_P$  and  $f_Q$ :  $f_P$  computes  $P = r_1, \dots, r_i$  on  $I_S$  using the conventional algorithm for fixpoint semantics for Datalog [1];  $f_Q$  is applied to the result of  $f_P$  and computes  $Q = r_{i+1}, \dots, r_n$ .

**Defining  $f_P$ .**  $f_P$  is a recursive function that takes a “database instance”  $d$  (represented as an XML document) as input, and computes a new “instance”  $d'$  (also represented as an XML document) that contains a copy of  $d$  and the results of the functions that compute the rules in  $P$ . If, after an iteration of  $f_P$ ,  $d$  and  $d'$  are not lexicographically identical, we iterate again using  $d'$  as input; otherwise (i.e., reach a fixpoint), we stop. The translation of the rules in  $P$  is as follows.

A rule  $r_i : A(\bar{x}) \leftarrow B_1(\bar{u}_1), \dots, B_n(\bar{u}_n)$  defining a conjunctive query where each  $B_j$  appears positively is implemented as a function  $f_i$  using nested loops to iterate over the cartesian product of relations  $B_1, \dots, B_n$ , thus computing all possible valuations to the variables in  $r_i$ , and returns the “tuples” that satisfy the conjunctive query in the body of the rule. Negation is dealt with in the usual way for Datalog with stratified negation:  $\neg R(\bar{x})$  holds if and only if  $\bar{x}$  is in the active domain but not in  $R$ . A set of rules

$$A(\bar{x}) \leftarrow \dots, \dots, A(\bar{x}) \leftarrow \dots$$

computing a union is implemented as a separate function for each rule (as above); the “union” is done by concatenating with duplicate elimination and lexicographic sorting the individual “relations” (i.e., sequences of tuples) returned by each function.

**Defining  $f_Q$ .** Since we do not allow recursive gen-

Edge <sub>0</sub>			FLC <sub>0</sub>		
<i>pid</i>	<i>eid</i>	<i>label</i>	<i>pid</i>	<i>first</i>	<i>last</i>
1	19	country	19	20	22
19	20	name			
19	22	capital			

ILS <sub>0</sub>		Value <sub>0</sub>	
<i>left</i>	<i>right</i>	<i>eid</i>	<i>value</i>
20	22	20	Brazil
		22	Brasilia

(a) Edge<sup>++</sup> mapping of the Mondial XML document.

Country <sub>1</sub>			Value <sub>1</sub>	
<i>country</i>	<i>name</i>	<i>capital</i>	<i>country</i>	<i>value</i>
19	20	22	20	Brazil
			22	Brasilia

(b) Inlining element ids.

Country <sub>2</sub>		
<i>country</i>	<i>name</i>	<i>capital</i>
19	Brazil	Brasilia

(c) Inlining element content.

Figure 5: Applying mapping scheme transformations.

eration of surrogates, we treat each invention rule  $A(*, \bar{x}) \leftarrow B_1(\bar{u}_1), \dots, B_n(\bar{u}_n)$  like a standard non-recursive rule as in  $f_P$ , except that an invented surrogate is added to each tuple in the result set, computed by a Skolem function  $f_A(\bar{x})$ . Note that the mapping expressions that populate an instance of  $S'$  from the result of  $f_Q$  are straightforward.

**Defining  $\pi'$ .** We obtain  $\pi'$  by converting  $\beta$  into XQuery functions as discussed above; note that, on a “canonical view” of an instance of  $S'$ , these functions compute a “canonical view” of the corresponding instance of  $S$ , to which we can apply  $\pi$  unchanged.

## 4 The LILO Algorithm

In this section we describe LILO (Lossless Inlining, Lossless Outlining), a mapping scheme design algorithm based on the framework discussed in Section 3. LILO uses Edge<sup>++</sup> as its initial mapping scheme and defines several equivalence-preserving transformations, some of which are similar to those defined in the literature. This section starts with an example that illustrates the application of two transformations to the Edge<sup>++</sup> mapping scheme, showing the Edge<sup>++</sup> mapping of the document in Figure 1(b), and the databases that result from each transformation. Next, we discuss in detail the transformations used in LILO and give an intuitive argument of why they are equivalence preserving. Finally, the LILO algorithm is discussed.

**Example 2.** Recall the fragment of the mondial

database shown in Figure 1(b) and its DTD in Figure 1(a). In this example, we focus on country elements and the constraints defined for them. For clarity, the figures omit relations not involved in the discussion. Note that subscripts are used to distinguish relations with the same name in different schemas (e.g.,  $\mathbf{Value}_1$  is a relation in  $S_1$ ).

**Edge<sup>++</sup> Mapping.** Figure 5(a) shows the Edge<sup>++</sup> mapping for the country element; we refer to its relational schema as  $S_0$  in the remainder of the example. In  $S_0$ , we have a recursive constraint to check that each element  $e$  labeled country has exactly two children  $c_1, c_2$ ;  $c_1$  is labeled name;  $c_2$  is labeled capital; and  $c_1$  precedes  $c_2$ . Moreover,  $S_0$  also has a structural constraint requiring both  $c_1, c_2$  to have a matching tuple in the  $\mathbf{Value}$  relation.

**Inlining Elements.** Consider a schema  $S_1$  that adds to  $S_0$  a relation  $\mathbf{Country}_1(\underline{country}, name, capital)$  for inlining country elements and their children, and constraints for enforcing the functional dependencies  $name \rightarrow country$  and  $capital \rightarrow country$ , for ensuring that each name and capital element has a single country element as its parent in the tree. The constraints for validating name and capital elements (Section 3.2) are also modified so that they access  $\mathbf{Country}_1$  as opposed to  $\mathbf{Edge}_1$ . Intuitively, a mapping scheme defined on  $S_1$  has one advantage over Edge<sup>++</sup>: the validation constraints for country elements can be specified directly using SQL [12], and thus enforced efficiently by relational database engines. The mapping  $\alpha_1 : \mathcal{R}(S_0) \rightarrow \mathcal{R}(S_1)$  is:

$$\begin{aligned}
& \mathbf{Diff}(e) :- \mathbf{Edge}_0(\_ , e, 'country') \\
& \mathbf{Diff}(e) :- \mathbf{Edge}_0(\_ , e, 'capital') \\
& \mathbf{Diff}(e) :- \mathbf{Edge}_0(\_ , c, 'country'), \mathbf{Edge}_0(c, e, 'name') \\
\mathbf{Country}_1(e, n, c) & :- \mathbf{Edge}_0(e, n, 'name'), \mathbf{Edge}_0(e, c, 'capital') \\
& \mathbf{Edge}_1(e, c, l) :- \mathbf{Edge}_0(e, c, l), \neg \mathbf{Diff}(e) \\
& \mathbf{FLC}_1(p, f, l) :- \mathbf{FLC}_0(p, f, l), \neg \mathbf{Diff}(p) \\
& \mathbf{ILS}_1(l, r) :- \mathbf{ILS}_0(l, r), \neg \mathbf{Diff}(l) \\
& \mathbf{Value}_1(e, v) :- \mathbf{Value}_0(e, v)
\end{aligned}$$

$\mathbf{Diff}$  computes all elements that are mapped into  $\mathbf{Country}_1$  instead of  $\mathbf{Edge}_1$ ,  $\mathbf{FLC}_1$  and  $\mathbf{ILS}_1$ . Figure 5(b) shows the resulting database instance. The instance of  $S_0$  is recovered by  $\beta_1 : \mathcal{R}(S_1) \rightarrow \mathcal{R}(S_0)$ :

$$\begin{aligned}
& \mathbf{Edge}_0(e, c, l) :- \mathbf{Edge}_1(e, c, l) \\
& \mathbf{Edge}_0(e, c, l) :- \mathbf{Edge}_1(e, \_ , 'country'), \mathbf{Country}(c, \_ , \_ ) \\
& \quad l = 'country' \\
& \mathbf{Edge}_0(e, c, l) :- \mathbf{Country}_1(e, c, \_ ), l = 'name' \\
& \mathbf{Edge}_0(e, c, l) :- \mathbf{Country}_1(e, \_ , c), l = 'capital' \\
& \mathbf{FLC}_0(p, f, l) :- \mathbf{FLC}_1(p, f, l) \\
& \mathbf{FLC}_0(p, f, l) :- \mathbf{Country}(p, f, l) \\
& \mathbf{ILS}_0(l, r) :- \mathbf{ILS}_1(l, r) \\
& \mathbf{ILS}_0(l, r) :- \mathbf{Country}_1(\_ , l, r) \\
& \mathbf{Value}_0(e, v) :- \mathbf{Value}_1(e, v)
\end{aligned}$$

It is easy to see that  $S_0 \equiv S_1$ :  $S_0 \preceq S_1$  via  $(\alpha_1, \beta_1)$  and that  $S_1 \preceq S_0$  via  $(\beta_1, \alpha_1)$

**Inlining Textual Content.** We can further modify  $S_1$  by storing the actual names and capitals of the countries instead of their surrogates; in this new schema  $S_2$ , there is no need to enforce the functional dependencies in  $S_1$ , nor the validating constraints for name and capital elements (note that the validation constraints for name elements that are not children of country are not affected by either transformation). Furthermore, no joins are required for finding names or capitals of the countries when processing queries. Instances of  $S_2$  are computed by  $\alpha_2 : \mathcal{R}(S_1) \rightarrow \mathcal{R}(S_2)$ :

$$\begin{aligned}
& \mathbf{Diff}(e) :- \mathbf{Country}_1(p, e, \_ ), \mathbf{Value}_1(e, \_ ) \\
& \mathbf{Diff}(e) :- \mathbf{Country}_1(p, \_ , e), \mathbf{Value}_1(e, \_ ) \\
& \mathbf{Edge}_2(e, c, l) :- \mathbf{Edge}_1(e, c, l) \\
& \mathbf{FLC}_2(p, f, l) :- \mathbf{FLC}_1(p, f, l) \\
& \mathbf{ILS}_2(l, r) :- \mathbf{ILS}_1(l, r) \\
\mathbf{Country}_2(e, n, c) & :- \mathbf{Country}_1(e, v_1, v_2), \\
& \quad \mathbf{Value}_1(v_1, n), \mathbf{Value}_1(v_2, c) \\
& \mathbf{Value}_2(e, v) :- \mathbf{Value}_1(e, v), \neg \mathbf{Diff}(e)
\end{aligned}$$

Figure 5(c) shows the resulting instance. To reconstruct the instances of  $S_1$ ,  $\beta_2 : \mathcal{R}(S_2) \rightarrow \mathcal{R}(S_1)$  uses invention rules (marked with  $*$ ) to replace the surrogates lost by  $\alpha_2$ :

$$\begin{aligned}
& \mathbf{Edge}_1(e, c, l) :- \mathbf{Edge}_2(e, c, l) \\
& \mathbf{FLC}_1(p, f, l) :- \mathbf{FLC}_2(p, f, l) \\
& \mathbf{ILS}_1(l, r) :- \mathbf{ILS}_2(l, r) \\
& \mathbf{PName}(*, e, n) :- \mathbf{Country}_2(e, n, \_ ) \quad (*) \\
& \mathbf{PCapital}(*, e, c) :- \mathbf{Country}_2(e, \_ , c) \quad (*) \\
& \mathbf{Country}_1(e, n, c) :- \mathbf{PName}(n, e, \_ ), \mathbf{PCapital}(c, e, \_ ) \\
& \mathbf{Value}_1(e, v) :- \mathbf{Value}_2(e, v) \\
& \mathbf{Value}_1(e, v) :- \mathbf{PName}(e, v, \_ ) \\
& \mathbf{Value}_1(e, v) :- \mathbf{PCapital}(e, \_ , v)
\end{aligned}$$

Again, note that  $S_2 \equiv S_1$ . Thus, the mapping scheme resulting from applying the two transformations above to Edge<sup>++</sup> is information-preserving. ■

#### 4.1 Equivalence Preserving Transformations

We now present the equivalence preserving transformations used in LILO. Our goal in defining them is twofold: we want to reduce the number of joins required for navigating the mapped document; and we want to replace the validation constraints defined by Edge<sup>++</sup>, which require recursive Datalog programs, into simpler ones that can be expressed in SQL. In the interest of space and for clarity of exposition, we omit the wrec-ILOG<sup>-</sup> programs in the discussion below; instead, we informally describe each transformation when arguing they are equivalence preserving.

**Equivalence Preservation.** We apply transformations to the relational schemas in the mappings because we are interested in eliminating some of the constraints defined in them. However, describing them as *changing the content models* in the DTD gives a more intuitive explanation of what they do. For example, suppose  $\mu = (\sigma, \pi, S)$  is the Edge<sup>++</sup> mapping scheme for the DTD in Figure 1(a); we can

view inlining the name of cities in  $S$  on as replacing  $city \leftarrow name, (province|state), official[1, 5]$  in the DTD by  $city^1 \leftarrow (province|state), official[1, 5]$ , modifying  $\sigma$  (resp.,  $\pi$ ) for storing (resp., fetching) the name elements in a separate relation. Some of our transformations result in *mixing* the content of different elements. For example, we can use a transformation for nesting (defined below) the content of cities elements as children of mondial, thus replacing  $mondial \leftarrow cities, country^*$  and  $cities \leftarrow city^*$  by  $mondial \leftarrow city^*, country^*$  (in this case,  $\sigma$  ignores cities elements and  $\pi$  is hard-coded to introduce them when processing a mondial element).

An important requirement for equivalence preservation in our transformations is that all content models they “introduce” are 1-unambiguous regular expressions, and if  $r = r_1, \dots, r_j, \dots, r_n$  is a 1-unambiguous regular expression and  $w \in L(r)$  then there exist *unique*  $w_1, \dots, w_n$  such that  $w_i \in L(r_i)$ ,  $1 \leq i \leq n$ , and  $w = w_1 \dots w_n$  [8]. (Note some  $w_i$  may be the empty string.) Thus, it follows that, even when we mix the contents of different elements, we can always distinguish, and thus reconstruct the original document.

**Notation.** In the remainder of this section,  $r, s$  denote regular expressions over  $\Sigma \cup \#PCDATA$  and  $a \in \Sigma$  denotes an element label.  $\mu = (\sigma, \pi, S)$  and  $\mu' = (\sigma', \pi', S')$  denote, respectively, the input and the output of each transformation. Recall that each rule  $t_i \leftarrow r_i$  in the DTD becomes a constraint in  $\text{Edge}^{++}$  involving the simulation of the appropriate DFA on the children of elements of type  $t_i$ ; we call such constraint an *edge* constraint and denote it by  $t_i \stackrel{R}{\leftarrow} r_i$  below.

The relational schema of the input mapping scheme is of the form  $S = \{E_1, \dots, E_n, M_1, \dots, M_k, \Gamma\}$ , where  $E_1, \dots, E_n$  are *edge* relations as defined by  $\text{Edge}^{++}$ ;  $M_1, \dots, M_k$  are *mapped* relations (resulting from the application of a transformation); and  $\Gamma = \Gamma^E \cup \Gamma^M$ , is the set of constraints in  $S$  where  $\Gamma^E$  contains all edge constraints, and  $\Gamma^M$  is a set of *mapped* constraints (resulting from the application of a transformation). Each transformation creates or modifies one or more mapped relations, and replaces one constraint in  $\Gamma^E$  by other constraints in  $\Gamma$  in a way that the resulting relational schema is equivalent to  $S$ .

#### 4.1.1 Inlining

Inlining [5, 29] consists of storing an element together with one or more of its children in the same relation.

**Inline\_Element**( $t, r_j, M$ ), where  $t \stackrel{R}{\leftarrow} r$  is an edge constraint,  $r = r_1, \dots, r_j, \dots, r_n$  and  $r_j$  is either  $a$  or  $a^?$ , results in using the mapped relation  $M$  to store pairs of surrogates of elements of type  $t$  and  $a$ , and applies only if  $s = r_1, \dots, r_{j-1}, r_{j+1}, \dots, r_n$  is 1-unambiguous. The schema for  $M$  is  $(\underline{t}, c_1, \dots, c_k, a)$  (the primary key is underlined), where  $c_1, \dots, c_k$ ,  $k \geq 0$ , are columns added by previous applications of inlining on type  $t$ . We replace the constraint  $t \stackrel{R}{\leftarrow} r$  by

$t \stackrel{R}{\leftarrow} s$ , and add a mapped constraint for checking that values of  $a$  in  $M$  are unique (that is, the FD  $a \rightarrow t$  holds). Since  $\mu'$  stores the surrogates for  $a$  elements in  $M$ , we also modify the edge constraint for validating  $a$  elements so that the datalog program refers to  $M$  as opposed to the original edge relations (recall the discussion in Example 2). Finally, if  $a$  elements are required (i.e.,  $r_j = a$ ), we define  $a$  to be not null in  $M$ .

Mapping instances of  $S$  into instances of  $S'$  is done by mapping the surrogates of inlined elements (of types  $c_1, \dots, c_k, a$ ) into  $M$ , while all other children of  $t$  elements are mapped into edge relations in  $S'$  (as illustrated in Example 2). The mapping in the other direction is done by copying the surrogates of the inlined elements into the edge relations in  $S$ . Since  $s = r_1, \dots, r_{j-1}, r_{j+1}, \dots, r_n$  is 1-unambiguous, we can recover the original element ordering in an instance of  $S$ : we know that each  $a$  element must occur *after* an element in  $r_1, \dots, r_{j-1}$  and *before* an element in  $r_{j+1}, \dots, r_n$ .

**Inline\_Union**( $t, r_j, M$ ) is a special case of inlining that applies when  $r_j = (s_1 | \dots | s_n)$ , and each  $s_i$  is either  $a_i$  or  $a_i^?$ . We add columns  $a_1, \dots, a_n$  to  $M$ , and define the following constraints. First, we need to enforce that at most one of these columns is not null in each tuple in  $M$ ; also, if no  $s_i$  is of the form  $a_i^?$ , then one of these columns must be not null. To ensure that at most one  $a_i$  is present at any time, we define a constraint comprising of a boolean formula of the form  $\phi_1 \vee \dots \vee \phi_k$  and each  $\phi_i$  checks  $a_i$  is not null while each  $a_j, i \neq j$  is. The constraint for ensuring that at least one of  $a_i$  is present is straightforward.

Note that the update anomaly discussed in Example 1 would be eliminated by defining the constraints discussed above in that relational schema.

**Inlining Textual Content.** When an element of content model  $\#PCDATA$  is inlined, we can use the mapped relation to store its content:

**Inline\_Cdata**( $a, M$ ), where  $a$  is the type of an inlined element in the mapped relation  $M$ , replaces the column for storing surrogates of  $a$  elements by a column for storing their textual content. As illustrated in Example 2, the mapping between instances of  $S$  and  $S'$  maps the tuples in **Value** corresponding to  $a$  elements into  $M$ . The reverse mapping is done by inventing new surrogates for each inlined element.

#### 4.1.2 Nesting Elements

This transformation eliminates some elements by nesting their contents within their parents:

**Nest**( $t, a$ ) applies when  $t \stackrel{R}{\leftarrow} r$  and  $a \stackrel{R}{\leftarrow} r'$  are edge constraints,  $r = r_1, \dots, a, \dots, r_n$ , and  $s = r_1, \dots, r', \dots, r_n$  is 1-unambiguous. As a result, we replace the constraints  $t \stackrel{R}{\leftarrow} r$  and  $a \stackrel{R}{\leftarrow} r'$  by  $t \stackrel{R}{\leftarrow} s$  in  $S'$ . The wrec-ILOG<sup>+</sup> programs for mapping instances of  $S$  and  $S'$  are straightforward, and since  $s$  is 1-unambiguous, we can distinguish the children of

**Input** : a DTD  $X = \langle \Sigma, r, \mathbf{R} \rangle$

**Output** : a mapping scheme  $\mu = (\sigma, \pi, S)$

1. Let  $\mu$  be the Edge<sup>++</sup> mapping scheme for  $X$
2. For each  $t \in \Sigma$ , create a mapped relation  $M_t$  in  $S$
3. Let  $t = r$
4. Let  $r = r_1 \dots, r_n$  be the content model associated with  $t$  in  $\mathbf{R}$
5. For each  $r_i, 1 \leq i \leq n$ , apply **Nest**( $t, r_i$ ), **Inline.Element**( $t, r_i, M_t$ ), **Inline.Cdata**( $r_i, M_t$ ), **Inline.Union**( $t, r_i, M_t$ ), **Outline**( $t, r_i$ ), **Outline.Union**( $t, r_i$ ) whenever possible and in this order
6. If a transformation creates a new type  $t'$ , add a mapped relation  $M_{t'}$  to  $S$
7. For each element type  $t_i$  occurring in  $r$ , let  $t = t_i$  and repeat from step 3 until no changes are made to  $S$
8. Remove from  $S$  all relations are not used by  $\sigma$  or  $\pi$
9. Return the resulting mapping scheme

Figure 6: LILO algorithm for designing information-preserving mapping schemes.

$a$  elements among the children of  $t$  elements in an instance of  $S'$ .

### 4.1.3 Outlining Elements

This transformation is the opposite of inlining and results in separate mapped relations for storing the content of elements of the same type. For instance, it could be used to separate city officials from other children of city elements; intuitively, this would replace  $city \leftarrow name, (province|state), official[1, 5]$  by  $city^1 \leftarrow name, (province|state)$  and  $city^2 \leftarrow official[1, 5]$  which could be further transformed independently.

**Outline**( $t, r_j$ ), applies when  $t \xleftarrow{R} r_1, \dots, r_j, \dots, r_m$  is an edge constraint, the participation of  $r_j$  is either  $*$  or  $+$ , and  $s = r_1, \dots, r_{j-1}, r_{j+1}, \dots, r_n$  is 1-unambiguous. In general, outline does not introduce new mapped relations; instead it replaces the original edge constraint by  $t \xleftarrow{R} s$  and  $t' \xleftarrow{R} r_j$ . The wrec-ILOG<sup>-</sup> programs for outlining map part of the content of each  $t$  element into the content of a new  $t'$  element (in the forward direction) and merge these pieces of content back (in the backward direction).

When  $r_j$  is either  $a*$  or  $a+$  we can use the mapped relation for capturing just parent-child relationship, ignoring labels (since they are all  $a$ ); thus, we can also drop the constraint  $t' \xleftarrow{R} r_j$  from the schema; the ordering of the  $a$  elements is still captured by the **FLC** and **ILS** relations.

A special case of outlining is when  $r_j = a[x, y]$  defining a cardinality constraint. In this case, we use a mapped relation for storing pairs of surrogates of  $t$  and  $a$  elements, and add a constraint that counts the number of  $a$  elements associated with each  $t$  element and checks this value is in the appropriate range. If the content model of  $a$  is #PCDATA, values of these elements can be inlined in the mapped relation.

**Outline.Union**( $t, r_j$ ) is a special case that applies

Operation	512KB	4MB	32MB	256MB	2GB
Edge <sup>++</sup> mapping scheme					
Insertion	7.72	15.87	17.70	20.48	20.45
Deletion	4.45	8.78	10.43	11.60	11.89
Q3	1.12	5.74	334.21	—	—
Q17	0.20	0.23	0.35	1.5	8
LILO mapping scheme					
Insertion	4.73	8.55	9.58	9.99	11.38
Deletion	3.53	6.41	7.48	8.16	8.76
Q3	0.09	0.12	0.18	1.47	50.41
Q17	0.03	0.03	0.03	0.09	0.59

Table 1: Experimental results; all times shown in milliseconds. For insertion and deletions, the time includes both modifying the database and checking the constraints.

when  $r_j = (u_1 | \dots | u_k)$  and results in replacing the original edge constraint  $t \xleftarrow{R} r_j$  by the following:

$$\begin{aligned}
 t_1 &\xleftarrow{R} r_1, \dots, r_{j-1}, u_1, r_{j+1}, \dots, r_n \\
 &\quad \vdots \\
 t_k &\xleftarrow{R} r_1, \dots, r_{j-1}, u_k, r_{j+1}, \dots, r_n
 \end{aligned}$$

It is easy to see that all regular expressions above are 1-unambiguous. One issue with this transformation is that the type of an element may change after its content is updated. Thus, whenever an update on element  $e$  of type  $t_i$  violates the edge constraint  $t_i \xleftarrow{R} r_1, \dots, r_{j-1}, u_i, r_{j+1}, \dots, r_n$ , we must check whether the new content of  $e$  satisfies the constraint of some other type  $t_j$ ; if so, we must also update the type of  $e$ .

## 4.2 The LILO Algorithm

The LILO algorithm is given in Figure 6. LILO visits each type  $t$  once, and after  $t$  is visited, its validating constraint is either mapped by some transformation or it is left unchanged. The order in which the transformations are attempted reflect the criteria discussed in the beginning of Section 4.1: reducing the number of joins required for navigating the document (achieved mostly by inlining), and simplifying the validation constraints introduced by Edge<sup>++</sup>.

**Theorem 1.** *Given a DTD  $X$ , LILO always terminates and produces a mapping scheme that is information-preserving with respect to  $X$ .*

*Proof.* It is easy to see that LILO always terminates: each rule in a the DTD is visited only once, and the total length (i.e., number of symbols) of all validation rules in  $\Gamma^E$  after each iteration of LILO is strictly smaller than at the beginning of that iteration. The resulting mapping scheme is information-preserving with respect to  $X$  since Edge<sup>++</sup> is information-preserving, and all transformation used by LILO are equivalence preserving (Section 3.1).  $\square$

## 5 Experimental Evaluation

Table 1 presents experimental results comparing the behavior of Edge<sup>++</sup> and LILO mapping schemes for

processing updates and queries, using XMark documents ranging from 512KB to 2GB. For processing updates, we must recompute the queries defining constraints in the relational schema; recomputing all such queries from scratch after each update is inefficient and, in most cases, unnecessary [17].

We used DB2 V8.1 on a Pentium-4 2.5GHz (1.5GB of RAM) running Linux 2.4; we limited DB2’s buffer to 45MB, to minimize the effects of buffering for larger databases. This is the only parameter we tune; furthermore, the only indices in each database are those created automatically for the primary keys of each relation. In our implementation of Edge<sup>++</sup>, we use horizontal partition in the **Edge** relation based on the type of the parent element, which eliminates some joins in the definition of the validation constraints; also, we adapt the incremental validation algorithm described in [4] to Edge<sup>++</sup>, thus avoiding the recomputation of the recursive constraints from scratch after every update. Finally, it is worth noting that, in this experiment, all transactions were executed at the user level, thus incurring overheads (e.g., query compilation and optimization) that can be avoided if these methods are implemented inside the database engine.

**Updates.** The update workload consists of 100 insertions and deletions of items for auctions in the North America region, each performed as a separate transaction. Each element inserted is valid and consists of an entire subtree of size comparable to those already in the document, and each deletion removes one of the “new” items inserted. Table 1 shows that while both mapping approaches scale very well with document size, LILO is on average 45% faster than Edge<sup>++</sup> for insertions and 26% faster for deletions.

**Queries.** While our focus in this paper is on information preservation, we compared Edge<sup>++</sup> and LILO for query processing as well. We used XMark queries whose focus was on navigating the documents, and here we show the results for two of them: Q3 and Q17; Q3 takes order into account and performs a join, while Q17 is among the simplest queries in the benchmark. For practical reasons, we set a timeout limit of 10 minutes for running each query. Table 1 shows that LILO is vastly superior to Edge<sup>++</sup> for query processing: Q3 times out on the 256MB Edge<sup>++</sup> mapping, while it takes about 5 seconds on the LILO mapping for the 2GB document. The improved performance is due mostly to the considerably fewer joins required for querying LILO mappings compared to Edge<sup>++</sup>.

## 6 Discussion and Related Work

We proposed a novel and sound framework for generating information-preserving XML-to-relational mapping schemes from an XML schema, which consists of applying equivalence-preserving transformations to schemes that are known to be information-preserving. As discussed, this process results in a mapping scheme

whose relational schema is *equivalent* to the original XML schema. Our framework is extensible and can handle arbitrary transformations that can be written in  $wrec\text{-}ILOG^{-}$ , which is powerful enough for the expressing the transformations of interest.

We discussed LILO: a mapping scheme design algorithm for  $\mathcal{XDS}$  that uses transformations to derive new mapping schemes that preserve the *validity* of XML documents. Using relational constraints to ensure validity has several advantages compared to the alternative of materializing and re-validating the portion of the document that is updated. Notably, our approach leverages the constraint-checking infrastructure already available in the DBMS, and does not require the development and maintenance of a separate validation tool. Moreover, our experiments have shown that even an implementation that incurs all the overheads of user-level transactions leads to acceptable performance.

While checking and incremental checking of constraints in the relational setting is well studied (see, e.g., [17] for a survey), researchers are only beginning to consider updating XML [6, 23, 30, 31], and the problems of validation and incremental validation after updates [4, 20, 27]. We studied two natural notions of information preservation for mapping schemes in previous work [3], including validation; we note that other authors have also identified the need for information preservation in mapping schemes and have informally discussed losslessness [14, 31].

While the literature on mapping schemes is vast, the focus of previous work has been mostly on performance of queries (see e.g., [21] for a survey) and updates [30, 31], rather than information preservation. Nevertheless, several existing methods (possibly with straightforward extensions) guarantee losslessness (only). For instance, numbering schemes that capture both element identity and ordering [31, 32], can be used to preserve the document structure [5, 14, 16, 29]. Edge<sup>++</sup> [3] captures the tree structure of the documents by representing the *successor* relation among elements (as opposed to explicit ordering); doing so results in more efficient update processing.

Some mapping schemes in the literature are oblivious to document schemas [16] while others exploit (and *require*) document schemas [5, 28, 29]. The latter approach has been shown to lead to better query performance, and is more closely related to our framework and LILO; a similar gain in query performance is noticed when comparing LILO to Edge<sup>++</sup>, whose relational schema is similar to the one in [16]. Some of the transformations we define can be viewed as extending those in [5, 28, 29] to guarantee information preservation. LegoDB [5] uses a cost-based model for designing mapping schemes whose goal is minimizing the *estimated* cost of executing an input query workload on an input XML document. Such estimates are

obtained by using features in modern RDBMSs that allow the query optimizer to use statistical information describing a *hypothetical* database instance. While it would be interesting to extend LILO with a cost-based model, strong empirical evidence suggests that, even for simple queries, there can be discrepancies of several orders of magnitude between the *estimated* costs and the execution costs observed on that actual instance [11]. Thus, other cost models may have to be considered.

Several techniques have been proposed for translating specific XML Schema constraints into relational ones in mapping schemes; for instance, techniques for translating keys [13], foreign-keys [10], cardinality constraints [5, 22], ID/IDREF attributes [4], and type specialization [3] have been proposed. However, to the best of our knowledge, no work has addressed the problem of mapping the element validity constraint, requiring the content of valid elements to be words in regular languages defined in the document schemas [7], which is achieved by LILO.

Finally, an interesting observation made in [2] is that some of the strategies defined in the literature are orthogonal, and new mapping schemes can be derived by mixing them. The framework we propose is extensible and can accommodate *any* transformation that can be described in  $wrec\text{-}ILOG^\neg$ . Thus, while we considered transformations for preserving element validity only, our framework is not tied to specific kinds of constraints: different transformations that preserve additional constraints can be easily incorporated.

## References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison Wesley, 1995.
- [2] S. Amer-Yahia, F. Du, and J. Freire. A Comprehensive Solution to the XML-to-Relational Mapping Problem. In *WIDM*, 2004.
- [3] D. Barbosa, J. Freire, and A. O. Mendelzon. Information Preservation in XML-to-Relational Mappings. In *XSym*, 2004.
- [4] D. Barbosa, A. O. Mendelzon, L. Libkin, L. Mignet, and M. Arenas. Efficient Incremental Validation of XML Documents. In *ICDE*, 2004.
- [5] P. Bohannon, J. Freire, P. Roy, and J. Siméon. From XML Schema to Relations: A Cost-based Approach to XML Storage. In *ICDE*, 2002.
- [6] V. P. Braganholo, S. B. Davidson, and C. A. Heuser. From XML View Updates to Relational View Updates: old solutions to a new problem. In *VLDB*, 2004.
- [7] T. Bray, J. Paoli, C. M. Sperberg-McQueen, and E. M. (Editors). *Extensible Markup Language (XML) 1.0*. World Wide Web Consortium, third edition, 2004.
- [8] A. Brüggemann-Klein and D. Wood. One-Unambiguous Regular Languages. *Information and Computation*, 142, 1998.
- [9] M. J. Carey, J. Kiernan, J. Shanmugasundaram, E. J. Shekita, and S. N. Subramanian. XPERANTO: Middleware for Publishing Object-Relational Data as XML Documents. In *VLDB*, 2000.
- [10] Y. Chen, S. Davidson, C. S. Hara, and Y. Zheng. RRRF: Redundancy reducing XML storage in relations. In *VLDB*, 2003.
- [11] M. Consens, D. Barbosa, A. Teisanu, and L. Mignet. Goals and Benchmarks for Autonomic Configuration Recommenders. In *SIGMOD*, 2005. To appear.
- [12] C. J. Date and H. Darwen. *A Guide to the SQL Standard*. Addison Wesley, 4th edition, 1997.
- [13] S. Davidson, W. Fan, C. Hara, and J. Qin. Propagating XML Constraints to Relations. In *ICDE*, 2003.
- [14] A. Deutsch, M. Fernández, and D. Suciu. Storing Semistructured Data with STORED. In *SIGMOD*, 1999.
- [15] M. Fernández, Y. Kadiyska, D. Suciu, A. Morishima, and W.-C. Tan. SilkRoute: A Framework for Publishing Relational Data in XML. *TODS*, 27(4), 2002.
- [16] D. Florescu and D. Kossmann. Storing and Querying XML Data Using an RDBMS. *IEEE Data Engineering Bulletin*, 22(3), 1999.
- [17] A. Gupta and I. S. Mumick, editors. *Materialized Views - Techniques, Implementations and Applications*. MIT Press, 1998.
- [18] R. Hull. Relative Information Capacity of Simple Relational Database Schemata. *SIAM Journal of Computing*, 15(3), 1986.
- [19] R. Hull and M. Yoshikawa. ILOG: Declarative Creation and Manipulation of Object Identifiers. In *VLDB*, 1990.
- [20] B. Kane, H. Su, and E. A. Rundensteiner. Consistently Updating XML Documents Using Incremental Constraint Check Queries. In *CIKM*, 2002.
- [21] R. Krishnamurthy, R. Kaushik, and J. F. Naughton. XML-SQL Query Translation Literature: the State of the Art and Open Problems. In *XSym*, 2003.
- [22] D. Lee and W. W. Chu. Constraints-Preserving Transformation from XML Document Type Definition to Relational Schema. In *ER*, 2000.
- [23] P. Lehti. Design and implementation of a data manipulation processor for an xml query language. Master's thesis, Universität Darmstadt, 2001.
- [24] D. Maier. *The Theory of Relational Databases*. Computer Science Press, 1983.
- [25] R. Miller, Y. Ioannidis, and R. Ramakrishnan. The Use of Information Capacity in Schema Integration and Translation. In *VLDB*, 1993.
- [26] M. Nicola and J. John. XML Parsing: a Threat to Database Performance. In *CIKM*, 2003.
- [27] Y. Papakonstantinou and V. Vianu. Incremental Validation of XML Documents. In *ICDT*, 2003.
- [28] M. Ramanath, J. Freire, J. R. Haritsa, and P. Roy. Searching for Efficient XML-to-Relational Mappings. In *XSym*, 2003.
- [29] J. Shanmugasundaram, K. Tufte, C. Zhang, G. He, D. J. DeWitt, and J. F. Naughton. Relational Databases for Querying XML Documents: Limitations and Opportunities. In *VLDB*, 1999.
- [30] I. Tatarinov, Z. G. Ives, A. Y. Halevy, and D. S. Weld. Updating XML. In *SIGMOD*, 2001.
- [31] I. Tatarinov, S. D. Viglas, K. Beyer, J. Shanmugasundaram, E. Shekita, and C. Zhang. Storing and Querying Ordered XML Using a Relational Database System. In *SIGMOD*, 2002.
- [32] C. Zhang, J. F. Naughton, D. J. DeWitt, Q. Luo, and G. M. Lohman. On Supporting Containment Queries in Relational Database Management Systems. In *SIGMOD*, 2001.