

Information Preserving XML Schema Embedding

Philip Bohannon[‡] Wenfei Fan^{† *} Michael Flaster[‡] P. P. S. Narayan[‡]

[‡] Bell Laboratories, Lucent Technologies {bohannon,mflaster,ppsn}@research.bell-labs.com

[†]University of Edinburgh & Bell Laboratories, wenfei@research.bell-labs.com

Abstract

A fundamental concern of information integration in an XML context is the ability to *embed* one or more source documents in a target document so that (a) the target document conforms to a target schema and (b) the information in the source document(s) is *preserved*. In this paper, information preservation for XML is formally studied, and the results of this study guide the definition of a novel notion of *schema embedding* between two XML DTD schemas represented as graphs. Schema embedding generalizes the conventional notion of graph similarity by allowing an edge in a source DTD schema to be mapped to a path in the target DTD. Instance-level embeddings can be defined from the schema embedding in a straightforward manner, such that conformance to a target schema and information preservation are guaranteed. We show that it is NP-complete to find an embedding between two DTD schemas. We also provide efficient heuristic algorithms to find candidate embeddings, along with experimental results to evaluate and compare the algorithms. These yield the first systematic and effective approach to finding information preserving XML mappings.

1 Introduction

A central technical issue for the exchange, migration and integration of XML data is to find mappings from documents of a source XML (DTD) schema to documents of a target schema. While one can certainly define XML mappings in a query language such as XQuery or XSLT, such queries may be large and complex, and in practice it is often needed that XML mappings (1) guarantee *type-safety* and (2) *preserve information*.

* Supported in part by EPSRC GR/S63205/01, EPSRC GR/T27433/01 and NSFC 60228006.

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

It is clearly desirable that the document produced by an XML mapping conforms to a target schema, guaranteeing *type safety*. But this may be difficult to check for mappings defined in XQuery or XSLT [4]. Further, since in many applications one does not want to lose the original information of the source data, a mapping should also preserve information. Criteria for *information preservation* include: (1) *invertibility* [16]: can one recover the source document from the target? and (2) *query preservation*: for a particular XML query language, can all queries on source documents in that language be answered on target documents? We now illustrate these concepts with an example.

Example 1.1: Consider two source DTDs S_0, S_1 and a target DTD S represented as graphs in Fig. 1 (we omit the `str-PCDATA-child` under `cno`, `credit`, `title`, `year`, `term`, `instructor`, `gpa` in Fig. 1(c)). A document of S_0 contains information of *classes* currently being taught at a school, and a document of S_1 contains *student* data of the school. The user wants to map the document of S_0 and the document of S_1 to a single instance of S , which is to collect data about *courses* and *students* of the school in the last five years. Here we use edges of different types to represent different constructs of a DTD, namely, *solid edges* for a concatenation type (a unique occurrence of each child), *dashed edges* for disjunction (one and only one child), and *star edges* (edge labeled ‘*’) for Kleene star (zero or more child). □

In this example, invertibility asks for the ability to reconstruct the original *class* and *student* documents from an integrated *school* document, while query preservation requires the ability to answer XML queries posed on *class* and *student* documents using the *school* document. Two natural questions are: (a) can one determine whether an XML mapping is information preserving? (b) is there an efficient method to find information-preserving XML mappings?

While type safety and information preservation are clearly desirable, an additional feature is the ability to map documents of DTDs that have *different structures*. A given source DTD may differ in structure from a desired target DTD. This is typical in data integration, where the target DTD needs to accommodate data from *multiple sources* and thus cannot be similar to any of the sources; see, e.g., the *class*, *student* DTDs and the *school* DTD in Fig. 1.

Background. While information preservation has been studied for traditional database transformations [3, 16, 27, 28], to our knowledge, no previous work has considered it for XML mappings. In fact, a variety of tools and models

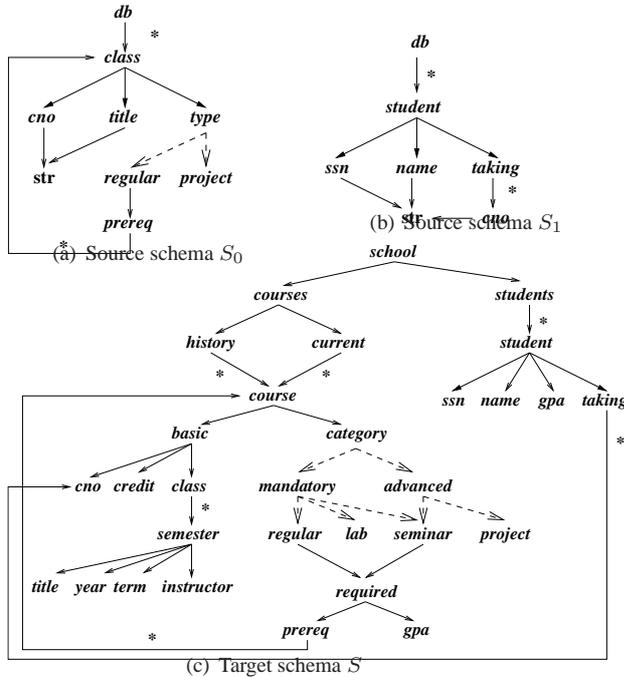


Figure 1: Example: source and target schemas

have been proposed for finding XML mappings at schema- or instance-level [13, 22, 24, 25, 26, 29]; however, none has addressed invertibility and query preservation for XML. Most tools either focus on *highly similar* structures, or adopt a strict graph similarity model like bisimulation (see., e.g., [1]) to match structures, which is incapable of mapping DTDs with *different structures* such as those shown in Fig. 1, and can ensure neither invertibility nor query preservation w.r.t. XML query languages. Another issue is that it is unclear that mappings found by some of these tools guarantee type safety when it comes to complex XML DTDs.

Contribution. To this end we study information preserving XML mappings, and make the following contributions.

First, as criteria for information preservation we revisit the notions of invertibility and query preservation [3, 16, 27, 28] for XML mappings (Section 2). While the two notions coincide for relational mappings w.r.t. relational calculus [16], we show that they are in general different for XML mappings w.r.t. XML query languages. Furthermore, we show that it is undecidable to determine whether or not an XML mapping defined in a simple fragment of XQuery (or XSLT) is information preserving (Section 3).

Second, to cope with the undecidability result, we introduce an XML mapping framework based on a novel notion of schema embeddings. A *schema embedding* is a natural extension of graph similarity in which an edge in a source DTD schema may be mapped to a *path*, rather than a single edge, in a target DTD. For example, the source DTDs S_0 and S_1 of Fig 1 can both be embedded in S , while there is no sensible mapping from them to S based on graph similarity. From a schema embedding, an instance-level XML mapping can be directly produced that has all the properties mentioned above. In particular, such mappings are invertible, query preserving w.r.t. regular XPath (an exten-

sion of XPath introduced in [23]), and ensure type safety. As with schema-mapping techniques for other data models, by automatically producing this mapping the user is saved from writing and type-checking a complex mapping query. Moreover, we show that the *inverse* and *query rewriting functions* for the mapping are efficient (Section 4).

Third, we provide algorithms to compute schema embeddings. We show that it is NP-complete to find an embedding between two DTDs, even when the DTDs are nonrecursive. Thus algorithms for finding embeddings are necessarily heuristic. A building block of our algorithms is an efficient algorithm to find a *local embedding* for individual productions in the source schema. Based on this, we develop three heuristic algorithms to compute embeddings. The first two algorithms repeatedly attempt to assemble local embeddings into a schema embedding (using a random or quality-specific order of the local embeddings, respectively), and when conflicts arise, attempt to generate new, non-conflicting local embeddings. The third algorithm generates a candidate pool of local embeddings, and then uses a heuristic solution to Maximum-Independent-Set to assemble a valid schema embedding (Section 5).

Finally, we have implemented our algorithms and conducted an experimental study based on mapping schemas taken from real-life and benchmark sources to copies of these schemas with varying amounts of introduced noise. These experiments verify the accuracy and efficiency of our heuristics on schemas up to a few hundred nodes in size (Section 6), and suggest that schema embeddings will lead to a promising tool for automatically computing information preserving XML mappings. We discuss related work in Section 7. Proofs are in the full version [8] of this paper.

To the best of our knowledge, this work is the first to study information preservation in the XML context, and it yields a systematic and effective approach to defining and finding information preserving XML mappings.

2 DTDs, XPath, Information Preservation

In this section we review DTDs and (regular) XPath, and revisit information preservation [16, 28] for XML.

2.1 XPath and Regular XPath

We consider a class of *regular* XPath queries proposed and studied in [23], denoted by \mathcal{X}_R and defined as follows:

$$\begin{aligned}
 p &::= \epsilon \mid A \mid p/text() \mid p/p \mid p \cup p \mid p^* \mid p[q], \\
 q &::= p \mid p/text() = 'c' \mid position() = k \\
 &\mid \neg q \mid q \wedge q \mid q \vee q.
 \end{aligned}$$

where ϵ is the empty path (*self*), A is a label (element type), ' \cup ' is the *union* operator, ' $/$ ' is the *child-axis*, and $*$ is the Kleene star; p is an \mathcal{X}_R expressions, k is a natural number, c is a string constant, and \neg, \wedge, \vee are the Boolean negation, conjunction and disjunction operators, respectively.

An *XPath fragment* of \mathcal{X}_R , denoted by \mathcal{X} , is defined by replacing p^* with $p//p$ in the definition above, where $//$ is the *descendant-or-self axis*.

A (regular) XPath query p is evaluated at a *context node* v in an XML tree T , and its result is the set of nodes (ids) of T reachable via p from v , denoted by $v[[p]]$.

2.2 DTDs

We consider DTDs of the form (Ele, P, r) , where Ele is a finite set of *element types*; r is a distinguished type in Ele , called the *root type*; P defines the element types: for each A in Ele , $P(A)$ is a regular expression of the form:

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

where str denotes PCDATA, ϵ is the empty word, B is a type in Ele (referred to as a *child of A*), and ‘+’, ‘,’ and ‘*’ denote *disjunction* (with $n > 1$), *concatenation* and the *Kleene star*, respectively. We refer to $A \rightarrow P(A)$ as the *production of A*. Note that this form of DTDs does not lose generality since any DTDs S can be converted to S' of this form (in linear time) by introducing new element types, and (regular) XPath queries on S can be rewritten into equivalent (regular) XPath queries on S' in PTIME [7].

Schema Graphs. We represent a DTD S as a labeled graph G_S , referred to as the *graph of S*. For each element type A in S , there is a unique node labeled A in G_S , referred to as the *A node*. From the A -node there are edges to nodes representing child types in $P(A)$, determined by the production $A \rightarrow P(A)$ of A . There are three different types of edges indicating different DTD constructs. Specifically, if $P(A)$ is B_1, \dots, B_n then there is a *solid edge* from the A node to each B_i node; it is labeled with a position k if B_i is the k -th occurrence of a type B in $P(A)$ (the label can be omitted if B_i 's are distinct). If $P(A)$ is $B_1 + \dots + B_n$ then there is a *dashed edge* from the A node to each B_i node (w.l.o.g. assume that B_i 's are distinct in disjunction). If $P(A)$ is B^* , then there is a *solid edge* with a ‘*’ label from the A node to the B node. Note that a DTD is *recursive* if its graph is *cyclic*. When it is clear from the context, we shall use the DTD and its graph interchangeably, both referred to as S ; similarly for A element type and A node.

For example, Fig. 1 shows graphs representing three DTDs, where Figs. 1(a) and 1(c) depict recursive DTDs.

An XML *instance* of a DTD S is a node-labeled tree that conforms to S . We denote by $\mathcal{I}(S)$ the set of all instances of S . A DTD S is *consistent* if it has no useless element types, i.e., each type of S has an instance. In the sequel we only consider consistent DTDs, w.l.o.g. since any DTD S can be converted to a consistent S' in $O(|S|^2)$ time such that $\mathcal{I}(S') = \mathcal{I}(S)$, by dropping all useless types from S .

2.3 Invertibility and Query Preservation

For XML DTDs S_1 and S_2 , a (data) *instance mapping* $\sigma_d : \mathcal{I}(S_1) \rightarrow \mathcal{I}(S_2)$ is *invertible* if there exists an inverse σ_d^{-1} of σ_d such that for any XML instance $T \in \mathcal{I}(S_1)$, $\sigma_d^{-1}(\sigma_d(T)) = T$, where $f(T)$ denotes the result of applying a function (or mapping, query) f to T . In other words, the composition $\sigma_d^{-1} \circ \sigma_d$ is equivalent to the identity mapping id , which maps an XML document to itself.

For an XML query language \mathcal{L} , a mapping σ_d is *query preserving w.r.t. \mathcal{L}* if there exists a computable function $F : \mathcal{L} \rightarrow \mathcal{L}$ such that for any XML query $Q \in \mathcal{L}$ and any $T \in \mathcal{I}(S_1)$, $Q(T) = F(Q)(\sigma_d(T))$, i.e., $Q = F(Q) \circ \sigma_d$.

In a nutshell, invertibility is the ability that the original source XML document can be recovered from the target document; query preservation w.r.t. \mathcal{L} indicates whether or not *all* queries of \mathcal{L} on any source T of S_1 can be effectively answered over $\sigma_d(T)$, i.e., the mapping σ_d does not lose information of T when \mathcal{L} queries are concerned.

The notions of invertibility and query preservation are inspired by (calculus) *dominance* and *query dominance* that were proposed in [16] for relational mappings and later studied in [3, 27, 28]. In contrast to query dominance, query preservation is defined w.r.t. a given XML query language that does not necessarily support query composition. Invertibility is defined for XML mappings and it only requires σ_d^{-1} to be a partial function defined on $\sigma_d(\mathcal{I}(S_1))$.

We say that a mapping $\sigma_d : \mathcal{I}(S_1) \rightarrow \mathcal{I}(S_2)$ is *information preserving w.r.t. \mathcal{L}* if it is both invertible and query preserving w.r.t. \mathcal{L} .

3 Information Preservation

In this section we establish basic results for separation and equivalence of the invertibility and query preservation of XML mappings, as well as complexity of determining whether a given XML mapping is information preserving.

Invertibility and Query Preservation. It was shown [16] that calculus dominance and query dominance are equivalent for relational mappings. In contrast, invertibility and query preservation do not necessarily coincide for XML mappings and query languages. Recall the class \mathcal{X} of *XPath queries* defined in Section 2, which supports neither query composition, nor identify mapping, nor the ability to navigate a recursive DTD based on certain patterns that are expressible in terms of the Kleen closure p^* .

Theorem 3.1: *There exists an invertible XML mapping that is not query preserving w.r.t. \mathcal{X} ; and there exists an XML mapping that is not invertible but is query-preserving w.r.t. the class of \mathcal{X} queries without position() qualifier. \square*

We identify sufficient conditions for the two to coincide: the definability of *the identity mapping*, and *query compossibility* (i.e., for any Q_1, Q_2 in \mathcal{L} , $Q_2 \circ Q_1$ is in \mathcal{L}).

Theorem 3.2: *Let \mathcal{L} be any XML query language and σ_d be a mapping from $\mathcal{I}(S_1) \rightarrow \mathcal{I}(S_2)$.*

- *If the identity mapping id is definable in \mathcal{L} and σ_d is query preserving w.r.t. \mathcal{L} , then σ_d is invertible.*
- *If \mathcal{L} is composable, σ_d is invertible and σ_d^{-1} is expressible in \mathcal{L} , then σ_d is query preserving w.r.t. \mathcal{L} . \square*

Recall the class \mathcal{X}_R of *regular XPath queries* defined in Section 2. Although the identity mapping id is not definable in \mathcal{X}_R , we show below that query preservation w.r.t. \mathcal{X}_R is a stronger property than invertibility: every node in a source document can uniquely identified by an \mathcal{X}_R query on the target document, and thus can be retracted.

Theorem 3.3: *If an XML mapping σ_d is query preserving w.r.t. \mathcal{X}_R , then σ_d is invertible. Conversely, there exists σ_d that is invertible but is not query preserving w.r.t. \mathcal{X}_R . \square*

Complexity. It is common to find XML mappings defined in XQuery or XSLT. A natural and important question is to decide whether or not an XML mapping is invertible or query preserving w.r.t. a query language \mathcal{L} . Unfortunately, this is impossible for XML mappings defined in any \mathcal{L} that subsumes first-order logic (FO , or relational algebra- RA), e.g., XQuery, XSLT, even when \mathcal{L} consists of projection queries only. Thus it is beyond reach to answer the question for XQuery or XSLT mappings.

Theorem 3.4: *It is undecidable to determine, given an XML mapping σ_d defined in any language subsuming FO , whether or not (a) σ_d is invertible; and (b) σ_d is query preserving w.r.t. projection queries. \square*

This can be verified by reduction from the equivalence problem for RA queries. The undecidability suggests that we start with languages simpler than XQuery and XSLT when studying information preserving XML mappings. Indeed, understanding (regular) XPath query preservation is a necessary step toward a full treatment of XML mappings defined in XQuery or XSLT, in which XPath is embedded.

4 Schema Embeddings for XML

The negative results in Section 3 tell us that it is already hard to determine whether or not an XML mapping is information preserving, not to mention finding one. This motivates us to look for a class of XML mappings that are *guaranteed* to be information preserving.

We approach this problem by specifying XML mappings at the schema level embeddings, and providing an automated derivation of instance-level mappings from these embeddings. Our notion of *schema embeddings* is novel, and extends the conventional notion of graph similarity by allowing edges in a source DTD schema to be mapped to a path in a target DTD with a “larger information capacity”. For example, a STAR edge can only be mapped to a path with at least one STAR edge.

In this section we define XML schema embeddings, present an algorithm for deriving an instance-level mapping from a schema embedding, and verify that the resulting mappings ensure information preservation.

4.1 Schema Level Embeddings

Consider a source XML DTD schema $S_1 = (E_1, P_1, r_1)$ and a target DTD $S_2 = (E_2, P_2, r_2)$. In a nutshell, a schema embedding σ is a pair of functions (λ, path) that maps each A type in E_1 to a $\lambda(A)$ type in E_2 , and each edge (A, B) in S_1 to a unique path $\text{path}(A, B)$ from $\lambda(A)$ to $\lambda(B)$ in S_2 , such that the S_2 paths mapped from sibling edges in S_1 are sufficiently distinct to allow information to be preserved. To define λ and path we first introduce a few notations.

\mathcal{X}_R Paths. An \mathcal{X}_R path over a DTD $S = (E, P, r)$ is an \mathcal{X}_R query of the form $\rho = \eta_1 / \dots / \eta_k$, where $k \geq 1$, η_i is

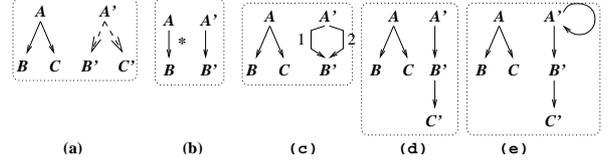


Figure 2: Path mappings for DTDs

of the form $A[q]$, and q is either *true* or a *position()* qualifier, such that ρ is a path in S and it carries all the position labels on the path. An \mathcal{X}_R path is called an *AND path* (resp. *OR path*, and *STAR path*) if it is nonempty and consists of only solid or star edges (resp. of solid edges and at least one dashed edge, and of solid edges and at least one edge labeled $*$). Referring to Fig. 1(c), for example, *basic/class/semester/title* is an AND path as well as a STAR path, and *mandatory/regular* is an OR path.

Name Similarity. A *similarity matrix* for S_1 and S_2 is an $|E_1| \times |E_2|$ matrix att of numbers in the range $[0, 1]$. For any $A \in E_1$ and $B \in E_2$, $\text{att}(A, B)$ indicates the suitability of mapping A to B , as determined by human domain experts or computed by an existing algorithm, e.g., [5, 13, 21].

Type Mapping. A *type mapping* λ from S_1 to S_2 is a (total) function from E_1 to E_2 ; it maps the root of S_1 to the root of S_2 , i.e., $\lambda(r_1) = r_2$. A type mapping λ is *valid* w.r.t. a similarity matrix att if for any $A \in E_1$, $\text{att}(A, \lambda(A)) > 0$.

Path Mapping. A *path mapping* from S_1 to S_2 , denoted by $\sigma : S_1 \rightarrow S_2$, is a pair (λ, path) , where λ is a type mapping and path is a function that maps each edge (A, B) in S_1 to an \mathcal{X}_R path $\text{path}(A, B)$ that is from $\lambda(A)$ to $\lambda(B)$ in S_2 .

For a particular element type A in E_1 , we say that σ is *valid* for A if the following conditions hold, based on the production $A \rightarrow P_1(A)$ in S_1 :

- if $P_1(A) = B_1, \dots, B_l$, then for each i , $\text{path}(A, B_i)$ is an AND path from $\lambda(A)$ to $\lambda(B_i)$ that is not a prefix of $\text{path}(A, B_j)$ for any $j \neq i$;
- if $P_1(A) = B_1 + \dots + B_l$, then for each i , $\text{path}(A, B_i)$ is an OR path from $\lambda(A)$ to $\lambda(B_i)$ that is not a prefix of $\text{path}(A, B_j)$ for any $j \neq i$ ¹;
- if $P_1(A) = B^*$, then $\text{path}(A, B_i)$ is a STAR path;
- if $P_1(A) = \text{str}$, then $\text{path}(A, \text{str})$ is an AND path ending with *text()*.

The validity requires a *path type* condition and a *prefix-free* condition, which, as will be seen shortly, are important for deriving the instance-level mapping from σ .

Example 4.1: Consider pairs of source (on the left) and target (on the right) DTDs depicted in Fig. 2, for which a type mapping λ is defined as $\lambda(X) = X'$ for X in $\{A, B, C\}$, except in Fig. 2(c) where both $\lambda(C) = B'$ and $\lambda(B) = B'$. Observe the following. For Fig. 2(a), there is no valid path embedding from the source DTD to the target; intuitively, B and C must coexist in a source document while only one of B' and C' exists in the target.

¹Abusing our normal form of DTDs, an optional type B can be specified as, e.g., $A \rightarrow B + \epsilon$; here $\text{path}(A, B_i)$ simply needs to be an OR path since ϵ is not an element type and thus $\text{path}(A, \epsilon)$ is undefined.

For Fig. 2(b), the source cannot be mapped to the target since there are possibly multiple B elements in the source, which cannot be accommodated by the target. For Fig. 2(c), a valid embedding is $\text{path}(A, B) = B'[\text{position}() = 1]$ and $\text{path}(A, C) = B'[\text{position}() = 2]$. For Fig. 2(d), there is no valid embedding since $\text{path}(A, B)$ is a prefix of $\text{path}(A, C)$, violating the prefix-free condition. For Fig. 2(e), a valid embedding is $\text{path}(A, B) = A'/B'$ (by unfolding the cycle once) and $\text{path}(A, C) = B'/C'$. \square

Finally, we define XML schema embeddings as follows.

Schema Embedding. A *schema embedding* from S_1 to S_2 *valid* w.r.t. a similarity matrix att is a path mapping $\sigma = (\lambda, \text{path})$ from S_1 to S_2 such that λ is valid w.r.t. att , and σ is valid for every element A in E_1 .

Example 4.2: Assume a similarity matrix att such that $\text{att}(A, A') = 1$ for all A in the DTD S_0 of Fig. 1(a) and A' in S of Fig. 1(c). The source DTD S_0 can be embedded in the target S via $\sigma_1 = (\lambda_1, \text{path}_1)$ defined as follows:

$\lambda_1(\text{db}) = \text{school},$	$\lambda_1(\text{class}) = \text{course},$	$\lambda_1(\text{type}) = \text{category},$
$\lambda_1(A) = A$	/* A : cno, title, regular, project, prereq, str */	
$\text{path}_1(\text{db}, \text{class})$	= courses/current/course	
$\text{path}_1(\text{class}, \text{cno})$	= basic/cno	
$\text{path}_1(\text{class}, \text{title})$	= basic/class/semester/title	
$\text{path}_1(\text{class}, \text{type})$	= category	
$\text{path}_1(\text{type}, \text{regular})$	= mandatory/regular	
$\text{path}_1(\text{type}, \text{project})$	= advanced/project	
$\text{path}_1(\text{regular}, \text{prereq})$	= required/prereq	
$\text{path}_1(\text{prereq}, \text{class})$	= course	
$\text{path}_1(A, \text{str})$	= text() /* A for cno, title */	

Note that $\text{path}_1(A, B)$ is a path in S denoting how to reach $\lambda_1(B)$ from $\lambda_1(A)$, *i.e.*, the path is *relative to* $\lambda_1(A)$. For example, $\text{path}_1(\text{type}, \text{project})$ indicates how to reach *project* from a *category* context node in S , where *category* is mapped from *type* in S_0 by λ_1 . Here the similarity matrix att imposes no restrictions: any name in the source can be mapped to any name in the target; thus the embedding here is decided solely on the DTD structures.

In contrast, one *cannot* map S_0 to S by graph similarity, which requires that node A in the source is mapped (similar) to B in the target only if all *children* of A are mapped (similar) to *children* of B . In other words, graph similarity maps an edge in the source to an edge in the target. \square

The definition of schema embedding can be extended to support further restructuring “across hierarchies” such that a child B of a source type A is not necessarily mapped to a descendant of $\lambda(A)$ in the target; this can be achieved via, e.g., upward modality in $\text{path}(A, B)$. It is also possible that an AND edge does not have to be mapped to an AND path. We focus on the main idea of schema embeddings in this paper and defer the full treatment to the full version.

Embedding Quality. There are many possible metrics. In this paper we consider only a simple one: the quality of a schema embedding $\sigma = (\lambda, \text{path})$ w.r.t. att is the sum of $\text{att}(A, \lambda(A))$ for $A \in E_1$, and we say that σ is *invalid* if λ is invalid w.r.t. att . We refer to this metric as $\text{qual}(\sigma, \text{att})$.

4.2 Instance Level Mapping

For a valid schema embedding $\sigma = (\lambda, \text{path})$ from S_1 to S_2 , we give its semantics by defining a (data) instance-level mapping $\sigma_d : \mathcal{I}(S_1) \rightarrow \mathcal{I}(S_2)$, referred to as the *XML mapping* of σ .

We define σ_d by presenting an algorithm that, given an instance T_1 of S_1 , computes an instance $T_2 = \sigma_d(T_1)$ of S_2 . In a nutshell, σ_d constructs T_2 top down starting from the root r_2 of T_2 , mapped from the root r_1 of T_1 (recall $\lambda(r_1) = r_2$). Inductively, for each $\lambda(A)$ element u in T_2 that is mapped from an A element v in T , σ_d generates a distinct $\lambda(B)$ node u' in T_2 for each distinct B child v' of v in T_1 , such that u' is reached from u via $\text{path}(A, B)$ in T_2 , *i.e.*, u' is uniquely identified by the \mathcal{X}_R path from u . More specifically, the construction is based on the production $A \rightarrow P_1(A)$ in S_1 as follows.

- (1) $P_A(A)$ is B_1, \dots, B_n . For each child v_i of v , σ_d creates a node u_i bearing the same id as v_i . These nodes are added to T_2 as follows. For each $i \in [1, n]$, u_i is added to T_2 by creating $\text{path}(A, B_i)$ emanating from u to u_i , such that the path shares any prefix already in T_2 which were created for, e.g., $\text{path}(A, B_j)$ for $j < i$. The definition of $\text{path}()$ ensures that u_i and u_j are not the same node in T_2 , since $\text{path}(A, B_i)$ is not a prefix of $\text{path}(A, B_j)$ and vice versa.
- (2) $P_1(A)$ is $B_1 + \dots + B_n$. Here v in T_1 must have a unique child v_i . For v_i , σ_d creates a node u_i bearing the same id as v_i , and adds u_i to T_2 via $\text{path}(A, B_i)$ as above.
- (3) $P_1(A)$ is B^* . By the definition of valid path function, $\text{path}(A, B)$ is of the form $\text{path}(A, A_1)/B_1/\text{path}(B_1, B)$, where A_1 is the first type defined in terms of Kleene star in P_2 , *i.e.*, $P_2(A_1) = B_1^*$. Let $[v_1, \dots, v_k]$ be the list of all the children of v . Then σ_d creates u_1, \dots, u_k bearing the same id's as v_1, \dots, v_k , and adds these nodes to T_2 as follows. It first generates a single $\text{path}(A, A_1)$ from u to an A' node u' if it does not already exist in T_2 , and for each $i \in [1, k]$, it creates a distinct i -th B_1 child if it is not already in T_2 . From the i -th B_i node it generates $\text{path}(B_1, B)$ leading to u_i , in the same way as in (1) above. Note that the order of the children of v is preserved by σ_d .
- (4) $P_1(A)$ is str . The treatment is the same as (1) except the last node of $\text{path}(A, \text{str})$ in T_2 is a text node holding the same value as the text node in T_1 .

We repeat the process until all nodes in T_1 are mapped to nodes in T_2 . We finally complete $\sigma_d(T)$ by adding *necessary* default elements such that $\sigma_d(T)$ conforms to S_2 . Recall from Section 2 that we can assume w.l.o.g. consistent DTDs. Thus for each element type A in S_2 , we can pick a fixed instance I_A of A and use it as A 's *default* element. The choice of default elements is arbitrary since as will be seen shortly, the inverse σ_d^{-1} of σ_d exists and it can distinguish T_2 nodes mapped from T_1 from default elements.

Example 4.3: Consider the XML mapping σ_d of the embedding defined in Example 4.2. Given an instance T_1 of S_0 of Fig. 1(a), σ_d generates a tree T_2 of S of Fig. 1(c) as follows: σ_d first creates the root *school* of T_2 , bearing the node id of the root *db* of T_1 . Then, σ_d creates a sin-

gle courses child x of school, a single current child y of x , and for each class child c of db , σ_d creates a distinct course child z of y bearing the id of c , such that the course children of y are in the same order as the class children of db . It then maps the *cno*, *title*, *type* children of c to *cno*, *title*, *category* descendants of z in T_2 , based on path_1 . In particular, to map *title* in S_0 , it creates a single class child x_c of the basic element, a single semester child x_s under x_c (although class is defined with a Kleene star), and then a title child under x_s . For the category element w mapped from the type child t of c , σ_d creates a distinct path *advanced/project* under w if t has a *project* child, or a *mandatory/regular* path otherwise, but not both. The process proceeds until all nodes in T_1 are mapped to T_2 . Finally, default elements of *history*, *credit*, *year*, *term*, *instructor* and *gpa* are added to T_2 such that T_2 conforms to S . At the last stage, no children of disjunctive types *category*, *mandatory* or *advanced* are added, and no children are created under *history*. That is, default elements are added only *when necessary*. \square

We next show that σ_d is *well defined*. That is, given any T_1 in $\mathcal{I}(S_1)$, $\sigma_d(T_1)$ is an XML tree that conforms to S_2 . This is nontrivial due to the interaction between different paths defined for disjunction types in the schema mapping σ , among other things. Consider, for example, $\text{path}(\text{type}, \text{regular})$ in Example 4.2. The path requires the existence of a *regular* child under a *mandatory* element m , which is in turn a child under a *category* element c in an instance of S . Thus it rules out the possibility of adding an *advanced* child under c or a *lab* child under m , perhaps requested by a *conflicting path* in σ . However, Theorem 4.1 below shows that the *prefix-free* condition in the definition of valid path functions ensures that conflicting paths do not exist.

Theorem 4.1 also shows that σ_d is *injective*: it maps distinct nodes in T_1 to distinct nodes in $\sigma_d(T_1)$, a property necessary for information preservation. Indeed, σ determines an injective *path-mapping* function δ such that, for each \mathcal{X}_R path $\rho = A_1[q_1]/\dots/A_k[q_k]$ in S_1 from r_1 , $\delta(\rho)$ is $\text{path}(r_1, A_1)[q_1]/\dots/\text{path}(A_{k-1}, A_k)[q_k]$, an \mathcal{X}_R path in S_2 from r_2 , by substituting $\text{path}(A_i, A_{i+1})$ for each A_{i+1} in ρ . Since each node in T_1 is uniquely determined by an \mathcal{X}_R path from the root, it follows that σ_d is injective.

Theorem 4.1: *The XML mapping σ_d of a valid schema embedding $\sigma : S_1 \rightarrow S_2$ is well defined and injective.* \square

4.3 Properties of Schema Embeddings

We have shown that the XML mapping σ_d of a valid schema embedding σ is guaranteed to type check. We next show that σ_d and σ also have all the other desired properties.

Information Preservation. In contrast to Theorem 3.4, information preservation is guaranteed by schema embeddings. Recall regular XPath \mathcal{X}_R from Section 2.

Theorem 4.2: *The XML mapping σ_d of a valid schema embedding $\sigma : S_1 \rightarrow S_2$ is invertible and is query preserving w.r.t. \mathcal{X}_R . More precisely, (a) there exists an inverse σ_d^{-1} of σ_d that, given any $\sigma_d(T)$, recovers T in $O(|\sigma_d(T)|^2)$ time;*

and (b) there is a query translation function F that given any \mathcal{X}_R query Q over S_1 , computes an \mathcal{X}_R query $F(Q)$ equivalent w.r.t. σ_d over S_2 in $O(|Q| |\sigma| |S_1|)$ time. \square

Example 4.4: The \mathcal{X}_R query Q below, over S_0 of Fig. 1(a), is to find all the classes that are (direct or indirect) prerequisites of CS331. It is translated to an \mathcal{X}_R query Q' over S of Fig. 1(c), which is equivalent w.r.t. the mapping σ_d given in Example 4.3, i.e. $Q(T) = Q'(\sigma_d(T))$ for any $T \in \mathcal{I}(S_0)$, when evaluated on T with the root as the context node.

Q : `class[cno/text()='CS331']/(type/regular/prereq/class)*`.

Q' : `courses/current/course[basic/cno/text()='CS331']/(category/mandatory/regular/required/prereq/course)*`. \square

In contrast, the notion of graph similarity ensures neither invertibility nor query preservation w.r.t. \mathcal{X}_R . As a simple example, the source and target schemas in Fig. 2(a) are bisimilar by the conventional definition of graph similarity, which does not consider cardinality constraints of different DTD constructs. However, there exists no instance-level mapping from the source to the target, not to mention inverse mappings and query translation.

Multiple sources. In contrast to graph similarity, it is possible to embed multiple source DTD schemas to a single target DTD, as illustrated by the example below. This property is particularly useful in data integration.

Example 4.5: The embedding $\sigma_2 = (\lambda_2, \text{path}_2)$ below maps S_1 of Fig. 1(b) to the target DTD S of Fig. 1(c).

$\lambda_2(\text{db}) = \text{school}$
 $\lambda_2(A) = A$ /* A: student, ssn, name, taking, cno */
 $\text{path}_2(\text{db}, \text{student}) = \text{students/student}$
 $\text{path}_2(\text{student}, B) = B$ /* B: ssn, name, taking */
 $\text{path}_2(\text{taking}, \text{cno}) = \text{cno}$
 $\text{path}_2(C, \text{str}) = \text{text}()$ /* C: ssn, name, cno */

Taken together with σ_1 of Example 4.2, this allows us to integrate a *course* document of S_0 and a *student* document of S_1 into a single *school* instance of the target DTD S . \square

In general, given multiple source DTDs S_1, \dots, S_n and a single target DTD S , one can define schema embeddings $\sigma_i : S_i \rightarrow S$ to simultaneously map S_i to S . Their XML mappings $\sigma_d^1, \dots, \sigma_d^n$ are invertible and query preserving w.r.t. \mathcal{X}_R as long as δ_i, δ_j are *pairwise disjoint*, where δ_i is the path mapping function derived from σ_i to map \mathcal{X}_R paths from root in S_i to \mathcal{X}_R paths from root in S . The instance-level XML mapping σ_d is a composition of individual $\sigma_d^1, \dots, \sigma_d^n$. Here σ_d^i increments the document constructed by σ_d^j 's for $j < i$ by modifying default elements or introducing new elements, instead of constructing a new document of S constructed starting from scratch.

Small model property. The result below gives us an upper bound on the length $|\text{path}(A, B)|$, and allows us to reduce the search space when defining or finding an embedding.

Theorem 4.3: *If there exists a valid schema embedding $\sigma : S_1 \rightarrow S_2$, then there exists one such that for any edge (A, B) in S_1 , $|\text{path}(A, B)| \leq (k + 1) |E_2|$, where $S_2 = (E_2, P_2, r_2)$, and k is the size of the production $P_2(A)$.* \square

5 Computing Schema Embeddings

In this section we address the computation of XML schema embeddings as defined by the following problem, stated in terms of two XML DTD schemas $S_1 = (E_1, P_1, r_1)$ and $S_2 = (E_2, P_2, r_2)$, and a similarity matrix att :

PROBLEM: Schema-Embedding
INPUT: Two DTDs S_1 and S_2 and matrix att .
OUTPUT: A schema embedding $\sigma : S_1 \rightarrow S_2$ valid w.r.t. att if one exists.

In practice, a reasonable goal is to find an embedding $\sigma : S_1 \rightarrow S_2$ with as high a value for $\text{qual}(\sigma, \text{att})$ as possible. The ability to efficiently find good solutions to this problem will lead to an automated tool that, given two DTD schemas, compute candidate embeddings to recommend to users.

However desirable, this problem is intractable. Worse, it remains NP-hard for nonrecursive DTDs even when they are defined in terms of concatenation types only.

Theorem 5.1: *The Schema-Embedding problem is NP-complete. It remains NP-hard for nonrecursive DTDs.* \square

In light of the intractable results we develop two efficient yet accurate heuristic algorithms for computing schema embedding candidates in the rest of the section.

Notations. Recall that a schema embedding is a path mapping σ that is valid for each element type A in S_1 . Since the validity conditions for A involve only A 's immediate children, it is useful to talk about mappings local to A . A *local mapping* for A is simply a *partial* path mapping $(\lambda_0, \text{path}_0)$ such that (a) λ_0 and path_0 are defined exactly on all the element types appearing in A 's production $A \rightarrow P_1(A)$, including A itself; and (b) it is *valid*, *i.e.*, it satisfies the path type and prefix-free conditions given in the last section.

Consider two partial mappings, $\sigma_0 = (\lambda_0, \text{path}_0)$ and $\sigma_1 = (\lambda_1, \text{path}_1)$. We say that λ_0 and λ_1 *conflict on* A if both $\lambda_0(A)$ and $\lambda_1(A)$ are defined, but $\lambda_0(A) \neq \lambda_1(A)$, and similarly for path_0 and path_1 . We say σ_0 and σ_1 are *consistent* if they do not conflict, either on λ or path. The *union of consistent partial mappings*, denoted by $\sigma_0 \oplus \sigma_1$, is a partial embedding $(\lambda_1 \oplus \lambda_2, \text{path}_1 \oplus \text{path}_2)$, where

$$\lambda_1(A) \oplus \lambda_2(A) = \begin{cases} \lambda_1(A) & \text{if } \lambda_2(A) \text{ is } \perp \text{ (undefined)} \\ \lambda_2(A) & \text{if } \lambda_1(A) \text{ is } \perp \\ \lambda_1(A) & \text{otherwise} \end{cases}$$

similarly for $\text{path}_1(A, B) \oplus \text{path}_2(A, B)$.

Outline. In the rest of the section we first present a technique for finding local embeddings, already a nontrivial yet interesting problem. Making use of this algorithm, we then provide three heuristics for finding embedding candidates. The first two are based on randomized programming and the last is by reduction from our problem to the Max-Weight-Independent-Set problem for which a well-developed heuristic tool [10] is available.

5.1 Finding Valid Local Mappings

We start by giving an algorithm to find a local embedding $\sigma_0 = (\lambda_0, \text{path}_0)$ when the partial type mapping λ_0

Algorithm findPathsDAG (G, s, L_{tar})

Input: Directed Acyclic Graph G , source node s ,
a bag of target nodes $L_{\text{tar}} = \{t_1, \dots, t_k\}$.

Output: Paths ρ_1, \dots, ρ_k satisfying the prefix-free condition.

1. path $\rho := \langle \text{empty} \rangle$;
2. $\mathcal{P} = \emptyset$;
3. marked (n) := false for all n ;
4. traverse ($G, s, \rho, L_{\text{tar}}, \mathcal{P}$);
5. if L_{tar} is nonempty
6. return \emptyset ;
7. else return \mathcal{P} ;

Figure 3: Algorithm findPathsDAG

is fixed, as this is a key building block of our schema-embedding algorithms. We then extend the algorithm to handle the general case when λ_0 is not given. To simplify the presentation we focus on nonrecursive DTDs, *i.e.*, DTDs with a *directed acyclic graph* (DAG) structure, but we show that our technique also works on recursive (cyclic) DTDs.

Finding Valid Paths. Let $A \in E_1$ be a source element type with production $A \rightarrow P_1(A)$, in which the element types appearing in $P_1(A)$ are B_1, \dots, B_k . Assume that the type mapping λ_0 is already given as a partial function from E_1 to E_2 that is defined on B_1, \dots, B_k and A . The Valid-Paths problem is to find paths $\text{path}_0(A, B_1), \dots, \text{path}_0(A, B_k)$ such that $(\lambda_0, \text{path}_0)$ is a valid local mapping for A .

The validity conditions stated for embeddings in Section 4.1 require that (a) target paths for each edge are of the appropriate *type* (AND, OR, or STAR path), and (b) that the target path for an edge is *not a prefix* of a *sibling's* target path. We abstract the second condition as a directed-graph problem: Given a directed graph $G = (V, E)$, a source vertex s and a *bag* of target vertices $L_{\text{tar}} = \{t_1 \dots t_k\}$, find paths ρ_1, \dots, ρ_k such that no path is the prefix of another. That is, for all $i \neq j$, $\rho_j \neq \rho_i / \rho_{ij}$ for any ρ_{ij} including the empty path. In contrast to most sub-problems of Schema-Embedding, this can be solved in *PTIME*. We introduce our solution by giving an algorithm that works only on a DAG and discuss extending it to handle cycles below.

We present our algorithm, findPathsDAG, in Fig. 3, for finding prefix-free paths in a DAG. The algorithm depends on the recursive procedure *traverse*, shown in Fig. 4. The intuition of this algorithm is to modify a simple (but exponential) algorithm to recursively enumerate all paths in a DAG in such a way that prefix-free paths are found, but excessive running time is avoided. In a nutshell, *traverse* conducts a depth-first-search on the input graph G , enumerating paths from the source node s to target nodes in L_{tar} , and identifies prefix-free ones. It uses a (global) boolean array *marked* (n) to keep track of whether the subgraph rooted at a node n has been searched and yielded no matches for nodes in L_{tar} , and if so, it does not re-enter the subgraph. A (local) variable *ret* is used to indicate whether the search of a subgraph finds any matches to nodes in L_{tar} .

To see that *traverse* is correct, consider removing line 5 in which the algorithm returns early, and line 11 in which nodes are marked to avoid revisiting them. It is clear that the resulting algorithm considers every possible path lead-

Algorithm `traverse` ($G, n, \rho, L_{\text{tar}}, \mathcal{P}$)

Input: Directed Acyclic Graph G , node n ,
a bag of target nodes $L_{\text{tar}} = \{t_1, \dots, t_k\}$,
 ρ , the current path to the root,
and \mathcal{P} , the output set of prefix-free paths.

Global variables: `marked`: maps nodes to $\{\text{true}, \text{false}\}$

Output: a list of paths.

1. if (marked (n)) return false;
2. if ($n \in L_{\text{tar}}$)
3. remove n from L_{tar} ;
4. add ρ to \mathcal{P}
5. return true;
6. else ret = false;
7. for each edge $e = (n, m)$ outgoing from n
8. append e to ρ ;
9. ret := ret or `traverse` ($G, m, \rho, L_{\text{tar}}, \mathcal{P}$);
10. remove e from ρ ;
11. if (not ret) `marked` (n):=true;
12. return ret;

Figure 4: Algorithm `traverse`

ing to nodes in L_{tar} , and assigns one path to each $n \in L_{\text{tar}}$, but it does not avoid assigning one node the prefix of another path. However, the prefix-free condition is assured by the return at line 5 *without affecting correctness*, since a suffix of the path assigned to n could only be generated by continuing the recursion from this node. Thus it remains to argue that the algorithm is still correct if line 11 is in place. The intuition of line 11 is simple: if no new target nodes were found in the subtree of a node when it was explored by the recursive calls of lines 7-10, then the current node will not be on any path to any n' remaining in L_{tar} .

Example 5.1: Consider the schema embedding problem shown in Fig. 1. Assume that `att` (*regular, seminar*), and `att` (*project, advanced*) in S_0 are 0.75. This means that the bag of possible target matchings for source tags $\{\text{regular}, \text{project}\}$ in S_0 can be $\{\text{seminar}, \text{advanced}\}$ from S . We then invoke `traverse` with S , *category*, ρ (which is empty), and L_{tar} as $\{\text{seminar}, \text{advanced}\}$. The first call to `traverse` would result in all edges from *category* to be recursed. Say, our algorithm first picks the edge to *advanced*. Line 2 of `traverse` would check *advanced* to be in L_{tar} and add the path to *advanced* into \mathcal{P} . It would then return back from the recursion and try the other edges from *category* in lines 7 though 10. This would result in a prefix-free path *mandatory/seminar* which would also be added to \mathcal{P} . \square

To analyze the performance of `findPathsDAG`, consider `traverse` as a sequence of forward and backward traversals of edges in the graph. A forward traversal occurs at line 9 and a backward traversal at lines 1, 5 and 12. Clearly, the number of forward traversals and backward traversals in a run are the same. Further, observe that one returns from an un-marked node at line 5 only on the path *back* from some node newly removed from L_{tar} . Thus, there can be at most $|L_{\text{tar}}| |V|$ such backward steps, and at most $|E|$ other backward steps (which mark the child of the edge traversed). Since G is a DAG, the algorithm is in $O(|L_{\text{tar}}| |V|)$ time.

To use `findPathsDAG` in our algorithms for schema embedding, we must further ensure that the paths returned

match the types needed for $n \in L_{\text{tar}}$. That is easy to accomplish, as the type of a path can be maintained incrementally as it is lengthened and shortened (by storing counts of nodes of each type), and be checked at line 2.

Schema Embeddings with a Given λ . This algorithm can be used to directly find a schema embedding $\sigma = (\lambda, \text{path})$ from S_1 to S_2 when the type mapping λ is a given total function from E_1 to E_2 . As remarked earlier, the validity conditions for any A in E_1 involve only A 's children; thus to find path we only need to find valid paths for each A in E_1 and take the union of these valid local embeddings. This yields an $O(|S_1| |S_2|)$ algorithm to find embeddings in this special setting, which is not so uncommon since one may know in advance which target type a source type should map to, based on, e.g., machine-learning techniques [13].

Handling Multiple Targets. However, to find valid local mappings when λ is not given, we must consider that there are *multiple* possible target nodes for each source node. The general Local-Embedding problem is to find a local embedding $(\lambda_0, \text{path}_0)$ when λ_0 may not be fixed. This problem is no longer tractable as shown below.

Theorem 5.2: *The Local-Embedding problem is NP-complete for nonrecursive DTDs.* \square

One heuristic approach to finding local embeddings is to extend `findPathsDAG` as follows. We compute the set of all pairings of source nodes A and possible matches for A from `att` and pass it as L_{tar} . We also modify line 3 of `traverse` to (a) pick an arbitrary pair with the current node as the target from L_{tar} at line 2 and (b) remove all pairs associated with source node A from L_{tar} at line 3. While this may work, it is essentially a greedy algorithm and may not find a solution if one exists. To compensate for this, we actually use a randomized variant `findPathsRand` (not shown) which (a) picks a random source node associated with n at line 2 of `traverse`, and (b) tries outgoing edges from n at line 7 in random order. The ability of `findPathsRand` to find embeddings varies with the size of L_{tar} , and will be investigated in Section 6.

Handling Cycles. Of course, schemas are frequently cyclic (recursive), and the algorithms as presented so far only handle DAGs. In fact, handling cycles generally is somewhat more complicated, but not hard – it is easy to see that an arbitrary number of paths can be generated by repeated loops around some cycle on the path to a target, and careful use of these paths can guarantee the prefix-free property (Figure 2(e) gives such an example, in which the cycle is unfolded once to get a prefix-free path, in contrast to Fig. 2(d)). While we present this full algorithm in [8], the complication is not warranted here since long cyclic paths are almost certainly semantically uninteresting. In practice, we have extended `findPathsDAG` once again to allow limited exploration of cycles limited by (a) no more than k trips through visited nodes and (b) no more than l total path length. A bound on k and l is given in Theorem 4.3 and usually k and l are set to small numbers.

Algorithm Ordered (S_1, S_2, O, C)

Input: Schemas S_1 and S_2 , an ordered set of source tags O , and C , a set of local embeddings for each source tag.
Output: a schema embedding from S_1 to S_2 if one is found.

```

1.  $\sigma :=$  empty solution ( $\emptyset, \emptyset$ );
2. for  $A$  in  $O$ 
3.   for  $\sigma_A$  in  $C(A)$ 
4.      $c :=$  conflict between  $\sigma$  and  $\sigma_A$ ;
5.     if  $c$  is null
6.        $\sigma = \sigma \oplus \sigma_A$ ; break;
7.     if  $c$  is not null
8.       findPathsRand ( $G, A, L_{\text{tar}}(A) - c$ );
9.     if  $c$  is not null return  $\emptyset$ ;
10. return  $\sigma$ ;
```

Figure 5: Algorithm Ordered

5.2 Three Methods for Finding Schema Embeddings

We next give three heuristic embedding-search algorithms: QualityOrdered, RandomOrdered and RandomMaxInd.

Finding Solutions with Ordered Algorithms. Our first two heuristics are based on a common subroutine Ordered, shown in Fig. 5. A key data structure is a table, C , where $C(A)$ is a set of known local embeddings for a source node A . The initialization of this table is discussed later. Given C and an ordered set O of source types, Ordered tries to assemble a consistent mapping σ by considering each A in O order (line 2), and trying to find a local embedding σ_A in $C(A)$ which can be merged with the existing σ without a conflict (lines 3-8). If a conflict occurs it finds new local embeddings for A by invoking findPathsRand (lines 7-8).

Our first Ordered-based algorithm, QualityOrdered, is shown in Fig. 6. Here $C(A)$ is initialized with a single randomly chosen local embedding for each source node A , and O is sorted by the *quality* of the local embedding.

In our second algorithm RandomOrdered (not shown), C is the complete set of local embeddings discovered so far for each source node (lines 4 and 5 in Fig. 6), while O is a random ordering of source nodes (line 6 in Fig. 6).

A Reduction Approach. We now discuss our third heuristic, RandomMaxInd. To understand this heuristic, consider the following problem defined on the table C of local mappings defined above:

PROBLEM:	Assemble-Embedding
INPUT:	Two DTDs S_1 and S_2 , a similarity matrix att , and a table C .
OUTPUT:	A schema embedding $\sigma : S_1 \rightarrow S_2$, valid w.r.t. att , formed as the union of a subset of embeddings in C if one exists.

Composing σ from local embeddings in C is nontrivial:

Theorem 5.3: *The Assemble-Embedding problem is NP-complete for nonrecursive DTDs.* \square

To cope with this, the RandomMaxInd heuristic takes the approach of reducing the Assemble-Embeddings problem to the problem of finding high-weight independent sets in a graph. It uses an existing heuristic solution [10] to produce

Algorithm QualityOrdered (S_1, S_2)

Input: Schemas S_1 and S_2 .

Output: a schema embedding from S_1 to S_2 if one is found.

```

1. count := 0;
2. while (count < MAX_TRIES) do
3.   count++;
4.   for each source node  $A$ 
5.      $C(A) :=$  {a local embedding,  $\sigma_A$  for  $A$ 
6.               as found by findPathsRand };
7.    $O :=$  All source nodes, ordered by  $\text{qual}(\sigma_A, \text{att})$ ;
8.    $\sigma :=$  Ordered ( $S_1, S_2, O, C$ );
9.   if  $\sigma \neq \emptyset$ 
10.    return  $\sigma$ ;
11. return  $\emptyset$ ;
```

Figure 6: Algorithm QualityOrdered

partial or complete solutions to this problem, which can be used to create partial or complete embeddings.

Before describing our reduction, we review the definition of Max-Weight-Independent-Set. That problem is defined on an undirected graph $G = (V, E)$ (not to be confused with a schema graph) with node weights $w[v], v \in V$. The goal is to find a subset V' of V such that for v_i and v_j in V' , there is no edge from v_i to v_j ; i.e., $(v_i, v_j) \notin E$ and the weight of V' , defined as $\sum_{v \in V'} w[v]$, is maximized.

Given an instance of the Assemble-Embedding problem, it is straightforward to construct an instance of Max-Weight-Independent-Set. First, for each local mapping $\sigma_a \in C(A)$ for any $A \in E_1$, we construct a vertex v_{σ_a} in V . Second, for each pair σ_a, σ_b of such mappings, we construct an edge between v_{σ_a} and v_{σ_b} if σ_a and σ_b conflict. The weight of v_{σ_a} is given as $\text{qual}(\sigma_a, \text{att})$.

To complete the algorithm on the resulting graph, we use an existing heuristic tool for Max-Weight-Independent-Set, which returns a subset V' of V . Finally, we construct an embedding σ by adding local embedding σ_a to σ for each $v_{\sigma_a} \in V'$. The quality of σ is warranted by the heuristic tool used, and its correctness is verified below.

Theorem 5.4: *If $|V'| = |E_1|$, σ constructed as above is a schema embedding from S_1 to S_2 .* \square

If σ is not a full embedding, we use findPathsRand to generate new local mappings, if any are available, for tags A not mapped by σ , and repeating the process until either it finds a valid embedding or it reaches a threshold of tries.

6 Experimental Study

In this section, we present an experimental evaluation of our schema embedding algorithms. Our approach is to vary the difficulty of the matching task by introducing artificial noise into a target schema, and measuring the ability of our algorithms to find an embedding.

Our experiments are based on real-world DTDs taken from a publicly available repository [30], plus the DTD of the XMark benchmark [33]. Each DTD was normalized into our graph representation. The XMark schema is the largest, with 57 productions after normalization. The XMark schema is apparently the most involved schema

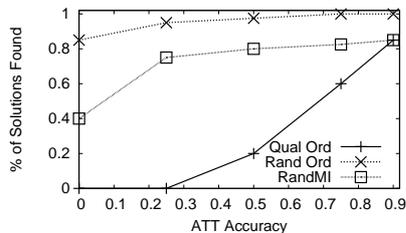


Figure 7: Varying accuracy

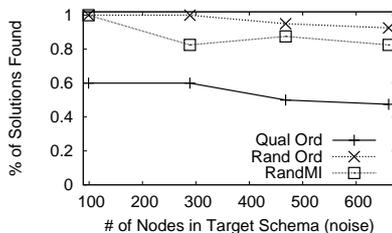


Figure 8: Noise vs. success rate

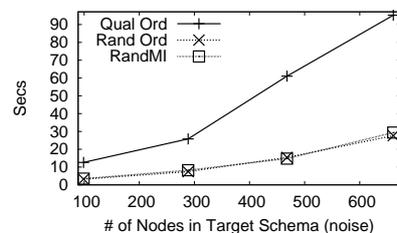


Figure 9: Noise vs. time

as the others scale better (see Fig. 10), and accordingly, we evaluate our algorithms for all the schemas but use the XMark schema for more detailed experiments.

Generating Target Schemas. Target schemas are generated from source schemas with added complexity and noise. As we introduce noise, we take care to preserve this matching, but make it harder to find in a number of ways, so as to attribute any failure to find a matching to the algorithm rather than the data. Particular target schemas are generated according to a probability noise in two steps: First, for each edge in the schema, with probability noise, the edge in the target is replaced with a path of between 1 and 5 nodes. When new nodes are added, with probability .5, the name of the node is formed as a small mutation of an existing name. Also, the type of the deleted edge (AND, OR, STAR) is used as the type of the first introduced edge to ensure that the original mapping is still possible.

In the second step, each node in the target (including newly-added nodes) are visited again, and with probability noise, a new subtree is added under it. The new subtree adds between 1 and 10 nodes. After each subtree addition, each leaf in the new subtree is visited, and with probability .5, an edge is added to an existing leaf outside the newly-added subtree. (This leaf may later have a subtree added under it.) The intuition for this last step is that confusion between different parts of the tree is more likely to arise if the same “attributes” (leaf nodes) appear in multiple places.

Generating the att. The similarity array, att, is initialized by computing pairwise string-edit distances between source and target tags (string edit distance with unit cost is also known as Damerau-Levenshtein distance). Furthermore, if a minimum threshold, sel, of similarity is not met by a pair, the similarity of that pair is set to 0, and as a result the tags cannot be matched. Note that the “similar names” introduced above range in similarity from .5 for short strings to over .8 for longer strings, and will be counted as potential matches in many experiments. There are also similar names in the schemas themselves, caused by the conversion of the schema to our graph format.

Clearly, sel, referred to as the *selectivity* of att, is an important parameter, as it directly determines the size of the candidate pool of target tags matching each source tag. Larger selectivities make the problem easier, and for our experimental data if sel is 1.0 (exact matches only), finding a schema embedding reduces to finding valid prefix-free paths for each local embedding in the source schema.

A second important parameter is the *accuracy* of att. This matters greatly for heuristic algorithms, since the valid

embedding in our generated data always has the highest average quality. Accuracy is implemented with a parameter c , which varies between 0 and 1. Each entry m in att is replaced by $cm + (1-c)rnd$, where rnd is a random number from 0 to 1. A low accuracy tends to mislead heuristics that rely heavily on att. Combining a low accuracy with a very low selectivity makes the problem very difficult to solve.

Experimental Setting. Experiments are conducted by copying the source schema, adding some amount of noise based on the parameter noise, and adjusting the att according to sel and c . Then the three algorithms given in Section 5 (RandomOrdered, QualityOrdered and RandomMaxInd) are used to try to find embeddings. For the ordered algorithms, the set C is initialized by finding 3 random mappings for each A , and discarding the two with the lowest qual ratings. When not otherwise stated, experiments are run with sel = 0.6, $c = 0.75$ (accuracy) and noise = 0.25. Since all algorithms (and the noise introduction) have a random component, they are repeated with 40 different random seeds, and an average is used.

The software is written in Java, except for the external heuristic for maximum independent sets [9], which is an optimized C program. Experiments are run on a variety of machines with Pentium III processors running at either 933MHZ or 1.0GHZ, with 256MB of RAM.

Accuracy Results. Figure 7 shows how the three algorithms perform while varying accuracy, with noise = 0.25. The y axis shows the percentages of runs for which a successful embedding is found. For this noise amount, the target schema is approximately three times as large as the source schema. This graph shows that QualityOrdered is extremely sensitive to the quality of the att values. It uses att extensively in its search pattern, and thus cannot find solutions unless att is accurate. Figure 7 also shows that RandomOrdered finds correct solutions more frequently than RandomMaxInd. While RandomOrdered takes into account att when it is seeking its solution set, it tries to find alternative solutions based on the conflicts it detects, independent of the att values. RandomMaxInd seeks alternative solutions for nodes based solely on their weights, as defined by att. It does not use conflicts to guide its search.

Varying Target Schema Size. We also consider target schemas with different numbers of erroneous nodes and edges introduced. These results are shown in Fig. 8. Because this graph shows results when accuracy is 0.75, QualityOrdered does not do well, as expected. RandomOrdered and RandomMaxInd both find the correct

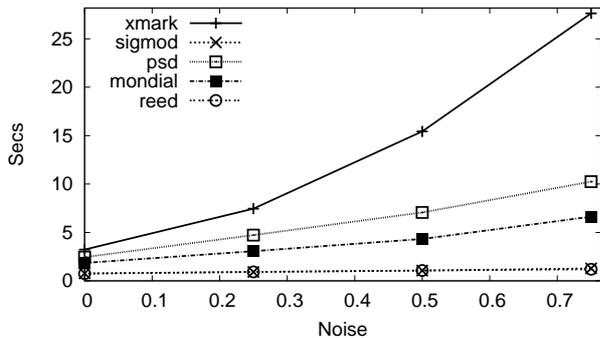


Figure 10: Time required for different source schemas

solution the majority of the time, decreasing somewhat as noise increases. The running times are shown in Fig. 9.

Different Source Schemas. We also run tests with different source schemas. We vary noise over five different source schemas, using RandomOrdered and accuracy=0.75. Figure 10 shows the running times for the various source schemas. For all runs across the different schemas, a solution was found more than 90% of the time (not shown).

Varying Selectivity. We also run experiments with different values of selectivity. Both RandomOrdered and RandomMaxInd find solutions less frequently as selectivity decreases (not shown). QualityOrdered is relatively indifferent to the selectivity level, finding approximately the same number of solutions at sel = 0.3 as at sel = 0.7. The running time increases dramatically, however, once sel falls below 0.4. The results are shown in Fig. 11.

Discussion. Our experimental results show that, when a feasible matching exists, it is likely to be almost completely found for schema sizes of up to a few hundred nodes. While this does not demonstrate that similar results can be obtained with differing target schemas and the use of real-world tools to produce att, it is certainly promising. Further, we found that the randomized algorithm RandomOrdered performs better than RandomMaxInd, and that QualityOrdered only does well with a highly accurate att. Based on these results, we plan to integrate RandomOrdered and RandomMaxInd, since the external independent set heuristic is very fast in practice. Finally, we note that QualityOrdered may be important in practice, where the att values may in fact be reliable.

7 Related Work

A wide variety of techniques have been developed to solve different forms of schema matching for relational, ER and object-oriented models (e.g., [5, 12, 18, 21, 31]; see [32] for a recent survey). While these are not focused on XML DTD schema matching, some techniques, such as linguistic analyses and machine learning, are useful for finding name/label similarity, which our algorithms take as input.

Closer to XML schema matching are [6, 13, 22, 24, 25, 26, 29]. LSD [13] proposes machine-learning techniques that make use of instance-level information to determine XML DTD tag mapping. Systems of [22, 24, 25] target

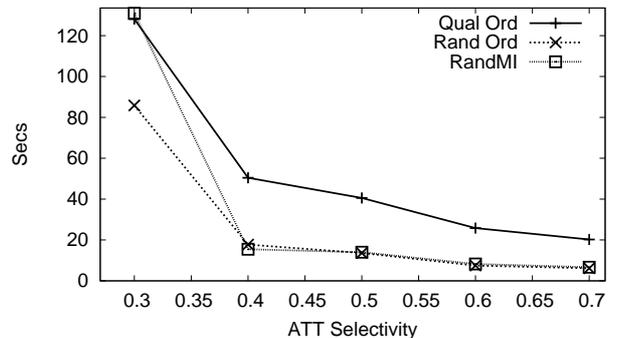


Figure 11: Varying Selectivity

a wide class of schemas and can be tailored to a variety of data models. The similarity flooding algorithm of [24] provides a novel schema matching tool based on graph-similarity. Cupid [22] is a generic system that encompasses a variety of techniques such as linguistic analyses and context dependencies. Rondo [25] proposes a powerful set of model mapping operators. For structure-level schema matching, these systems adopt graph similarity to map a single source schema to a target. TransScm [29] considers instance-level mappings based on schema matching, and uses a semi-automatic mechanism to match highly similar schemas. Clio [26] also focuses on deriving instance translation from schema mappings. The recent work [6] studies invertible XML-to-relation mappings that guarantee the source XML document remains valid in the presence of updates to the mapped relations. To our knowledge, no previous work has considered information preservation for XML DTD schema mappings. Our notion of schema embedding extends graph similarity and allows multiple source DTD schemas to be mapped to a single structurally different target DTD. Furthermore, from a schema embedding an instance mapping can be *automatically* derived and it *guarantees* both invertibility and query preserving w.r.t. regular XPath queries. The ability of finding information-preserving XML mappings is important for data integration (see, e.g., [19]) and P2P systems (e.g., [14, 17, 34]).

Information preservation has been studied for nested relational and complex data models (e.g., [3, 16, 27, 28]). [16] proposed several notions of dominance and studied their relationships, which were revisited in [27]. The focus of [3, 28] has mainly been on the information capacity of type constructs and structural transformation rules. Our study of information preservation is inspired by the prior work: our notions of invertibility and query preservation are mild extensions of calculus dominance and query dominance [16]. We revise these notions and study their basic properties for XML DTD schemas and XML queries, and our focus is to develop the notion of DTD schema embedding that preserves information by ensuring both effective invertible mapping and efficient XML query translation.

Query preservation is related to query rewriting using views, which has been extensively studied for conjunctive and datalog queries for relational databases and regular path queries on semistructured data (e.g., [2, 11, 20]; see [15, 19] for surveys). View-based query rewriting

mainly studies whether a given query on the source can be answered using materialized data from a set of views (lossless), by translating the query to an equivalent query in a particular language on the views. In contrast, query preservation deals with the issue whether *all* queries in an (infinite) query language on an XML source can be rewritten to equivalent queries over XML target (view). Moreover, the focus of this work is to generate XML “views” that automatically preserves all the queries in an XML query language, rather than to determine the losslessness of views. Note that Theorem 3.2 establishes a connection between invertibility and query rewriting; *e.g.*, if the query language \mathcal{L} includes the identity query id , then a view σ_d is invertible and σ_d^{-1} is in \mathcal{L} iff id has a rewriting in \mathcal{L} using σ_d .

8 Conclusions

We have revised information-preservation criteria for XML mappings and established separation, equivalence and complexity results. We have introduced a novel notion of schema embedding for XML DTD schemas, from which an instance-level XML mapping is automatically derived and is guaranteed to be information preserving, type checking, and able to accommodate multiple source schemas. While we show that finding a schema embedding is NP-complete, we have provided heuristic algorithms to compute embeddings, which are efficient and accurate as shown by our experimental results. These yield a practical approach to computing lossless XML data migration and integration.

We plan to extend the notion of schema embedding to (a) accommodate more general XML schemas with constraints and inheritance, (b) allow one source type to map to different target types in *different contexts*, (c) allow certain queries in XQuery in the path function, and (d) preserve XQuery fragments as query languages.

References

- [1] S. Abiteboul, P. Buneman, and D. Suciu. *Data on the Web. From Relations to Semistructured Data and XML*. Morgan Kaufman, 2000.
- [2] S. Abiteboul and O. M. Duschka. Complexity of answering queries using materialized views. In *PODS*, 1998.
- [3] S. Abiteboul and R. Hull. Restructuring hierarchical database objects. *TCS*, 62(1-2), 1988.
- [4] N. Alon, T. Milo, F. Neven, D. Suciu, and V. Vianu. XML with data values: Typechecking revisited. In *PODS*, 2001.
- [5] V. Athitsos, M. Hadjieleftheriou, G. Kollios, and S. Sclaroff. Query-sensitive embeddings. In *SIGMOD*, 2005.
- [6] D. Barbosa, J. Freire, and A. Mendelzon. Designing information-preserving mapping schemes for XML. In *VLDB*, 2005.
- [7] M. Benedikt, W. Fan, and F. Geerts. XPath satisfiability in the presence of DTDs. In *PODS*, 2005.
- [8] P. Bohannon, W. Fan, M. Flaster, and P. Narayan. Information preserving XML schema embedding. Full version, 2005.
- [9] S. Busygin. QUALEX: QUick ALmost EXact maximum weight clique/independent set solver. <http://www.busygin.dp.ua/npc.html>.
- [10] S. Busygin, S. Butenko, and P. M. Pardalos. A heuristic for the maximum independent set problem based on optimization of a quadratic over a sphere. *J. Comb. Optim.*, 6(3):287–297, 2002.
- [11] D. Calvanese, G. D. Giacomo, M. Lenzerini, and M. Y. Vardi. Lossless regular views. In *PODS*, 2002.
- [12] S. Castano, V. D. Antonellis, and S. D. C. di Vimercati. Global viewing of heterogeneous data sources. *TKDE*, 13(2):277–297, 2001.
- [13] A. Doan, P. Domingos, and A. Y. Halevy. Reconciling schemas of disparate data sources: A machine-learning approach. In *SIGMOD*, 2001.
- [14] A. Fuxman, P. Kolaitis, R. Miller, and W. Tan. Peer data exchange. In *PODS*, 2005.
- [15] A. Y. Halevy. Theory of answering queries using views. *SIGMOD Record*, 29(4), 2001.
- [16] R. Hull. Relative information capacity of simple relational database schemata. *SIAM J. Comput.*, 15(3):239–265, 1986.
- [17] A. Kementsietsidis, M. Arenas, and R. Miller. Mapping data in peer-to-peer systems: Semantics and algorithmic issues. In *SIGMOD*, 2003.
- [18] L. Lakshmanan, F. Sadri, and I. N. Subramanian. SchemaSQL – a language for interoperability in relational multi-database systems. In *VLDB*, 1996.
- [19] M. Lenzerini. Data integration: A theoretical perspective. In *PODS*, 2002.
- [20] A. Y. Levy, A. O. Mendelzon, Y. Sagiv, and D. Srivastava. Answering queries using views. In *PODS*, 1995.
- [21] W.-S. Li and C. Clifton. SemInt: A tool for identifying attribute correspondences in heterogeneous databases using neural networks. *Data Knowl Eng*, 33(1):49–84, 2000.
- [22] J. Madhavan, P. A. Bernstein, and E. Rahm. Generic schema matching with Cupid. In *VLDB*, 2001.
- [23] M. Marx. XPath with conditional axis relations. In *EDBT*, 2004.
- [24] S. Melnik, H. Garcia-Molina, and E. Rahm. Similarity flooding: A versatile graph matching algorithm. In *ICDE*, 2002.
- [25] S. Melnik, E. Rahm, and P. A. Bernstein. Rondo: A programming platform for generic model management. In *SIGMOD*, 2003.
- [26] R. J. Miller, M. A. Hernández, L. M. Haas, L.-L. Yan, C. T. H. Ho, R. Fagin, and L. Popa. The Clio project: Managing heterogeneity. *SIGMOD Record*, 30(1):78–83, 2001.
- [27] R. J. Miller, Y. E. Ioannidis, and R. Ramakrishnan. The use of information capacity in schema integration and translation. In *VLDB*, 1993.
- [28] R. J. Miller, Y. E. Ioannidis, and R. Ramakrishnan. Schema equivalence in heterogeneous systems: bridging theory and practice. *IS*, 19(1):3–31, 1994.
- [29] T. Milo and S. Zohar. Using schema matching to simplify heterogeneous data translation. In *VLDB*, 1998.
- [30] U. of Washington. XML repository. <http://www.cs.washington.edu/research/xmldatasets>.
- [31] L. Palopoli, D. Sacca, and D. Ursino. Semi-automatic semantic discovery of properties from database schemas. In *IDEAS*, 1998.
- [32] E. Rahm and P. A. Bernstein. A survey of approaches to automatic schema matching. *VLDB Journal*, 2001.
- [33] A. Schmidt, F. Waas, M. Kersten, M. J. Carey, I. Manolescu, and R. Busse. XMark: A Benchmark for XML Data Management. In *VLDB*, 2002.
- [34] I. Tatarinov et al. The Piazza peer data management project. *SIGMOD Record*, 32(3), 2003.