

REVERSE ENGINEERING BASED APPROACH FOR TRANSFERRING LEGACY RELATIONAL DATABASES INTO XML¹

Chunyan Wang, Anthony Lo, Reda Alhadj, Ken Barker
Advanced Database Systems and Applications Lab
Department of Computer Science
University of Calgary, Alberta, CANADA
{wangch, chiul, alhadj, barker}@cpsc.ucalgary.ca

Abstract: *XML (extensible Markup Language) has emerged, and, is being gradually accepted as the standard for data interchange over the Internet. Since most data is currently stored in relational database systems, the problem of converting relational data into XML assumes special significance. Many researchers have already done some accomplishments in this direction. They mainly focus on finding XML schema (e.g., DTD, XML-Schema, and RELAX) that best describes a given relational database with a corresponding well-defined database catalog that contains all information about tables, keys and constraints. However, not all existing databases can provide the required catalog information. Therefore, these applications do not work well for legacy relational database systems that were developed following the logical relational database design methodology, without being based on any commercial DBMS, and hence do not provide well-defined metadata files describing the database structure and constraints. In this paper, we address this issue by first applying the reverse engineering approach described in [2] to extract the ER (Extended Entity Relationship) model from a legacy relational database, then convert the ER to XML Schema. The proposed approach is capable of reflecting the relational schema flexibility into XML schema by considering the mapping of binary and nary relationships. We have implemented a first prototype and the initial experimental results are very encouraging, demonstrating the applicability and effectiveness of the proposed approach.*

Keywords: *Reverse engineering, legacy relational databases, XML.*

1. INTRODUCTION

As the World Wide Web is becoming a major means of disseminating and sharing information, there is an exponential increase in the amount of data in a web-

¹ A preliminary version of this paper appeared in *Proceedings of the 6th International Conference on Enterprise Information Systems, ICEIS2004, Porto, Portugal, April 4 -17, 2004.*

compliant format such as HTML (HyperText Markup Language) and XML (extensible Markup Language). Especially, the XML model is a novel textual representation of hierarchical (tree-like) data where a meaningful piece of data is bounded by matching starting and ending tags, such as <name> and </name>, respectively. Due to the simplicity of XML compared to SGML (Standard Generalized Markup Language) and its relatively powerful expressiveness compared to HTML, XML has become ubiquitous, and XML data need to be managed in databases.

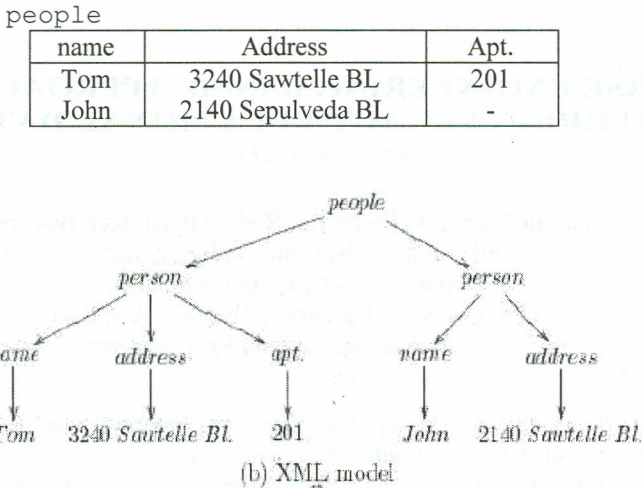


Figure 1: Representing people in both: a) relational and b) XML models.

The same data can be correctly captured in different models, including relational and XML models. However, as illustrated in Figure 1, there are subtle differences. The XML model representation can capture more detailed structural information (e.g., person) due to its hierarchical structure and does not suffer from unnecessary null information (e.g., there is no apt. branch in the second person node). In the relational model, data are represented in a flat structure where the only available constructs are tables and attributes, with foreign keys providing links between tables. A tree-like structure in XML model replaces the notion of tables, where leaf nodes are equivalent to the notion of attributes in the relational model.

XML is emerging as the standard format for data exchange. Among its great advantages are portability, extensibility and the possibility to add semantics, i.e., in particular structural constraints, to data within the document itself. Therefore, one of the important tasks that XML will solve is the exchange of data and information between different partners. Since most of the data nowadays reside in relational databases, it is important to automate the process of generating XML documents containing information from existing databases. Of course, one would like to preserve as much information as possible during the transformation process, especially constraints.

The Relational-to-XML conversion is considered complex because the two data models are significantly different. While relational data is flat, normalized into many relations, and the schema is often proprietary, XML data is nested, unnormalized, and its schema is public, mostly created by agreement between members of a community

after lengthy negotiations. The Relational-to-XML conversion involves mapping the relational tables and attributes names into XML elements and attributes names, creating XML hierarchies, and processing values in an application specific manner.

Considerable work has been done on transforming relational databases that have rich corresponding catalogs. Although a large number of the existing relational databases are

classified as legacy, the conversion of legacy relational databases has received little attention. Legacy databases are characterized by old-fashioned architecture, non-uniformity resulting from numerous extensions, and lack of the related documentation, i.e., little or no information about their design and constraints. Realizing the importance of converting legacy databases into XML documents, we have developed a method that successfully handles the process. Our approach highly benefits from our previous finding on reverse engineering of legacy databases detailed in [2]. Reverse engineering has been extensively studied as the process of discovering the characteristics of a legacy system. It leads to identifying and understanding all components of an existing system and the relationships between them. Database reverse engineering is necessary to semantically enrich and document a legacy database, and to avoid throwing away the huge amounts of data stored in existing legacy databases if the owners of existing legacy databases want to re-engineering, or maintain and adjust the corresponding database design.

Two basic steps are identified in the process of converting legacy databases into XML documents. First, reverse engineering is employed to deduce information about functional dependencies, keys and inclusion dependencies; the process involves reconstructing the ER model from an existing legacy database. Second, the obtained ER model is transformed into XML schema in a process known as forward engineering.

The rest of the paper is organized as follows. Related work is discussed in Section 2. Section 3 is an overview of the reverse engineering process to extract ER model from the existing relational database; for more details, the reader is referred to [2]. ER model to XML schema conversion is presented in Section 4. A closer look at the developed approach and the implemented prototype is given in Section 5. Section 6 is the conclusions.

2. RELATED WORK

There exist several tools that enable the composition of XML documents from relational data, such as IBM DB2 XML Extender, SilkRoute, and XPERANTO. XML Extender [7] serves as a repository for XML documents as well as their Document type declarations (DTDs), and also generate XML documents from existing data stored in relational database. DAD or XML Extender is used to define the mapping of DTD to relational tables and columns. XSLT and Xpath syntax are used to specify the transformation and the location path. SilkRoute [9] is described as a general, dynamic, and efficient tool for viewing and querying relational data in XML. In SilkRoute, data is exported into XML in two steps. The first step is to create an XML view of relational database by using a declarative query language RXL (Relational to XML Transformation Language). The resulting XML view is virtual. The second step is to formulate a query (e.g., XML-QL) over the XML view and to pass the query to the query composer. The query composer computes and generates a new executable RXL query, which is then translated by the translator into one or more SQL queries. The XML generator receives the result of the database query, and then translates them into XML documents. XPERANTO [6] (XML Publishing of Entities, Relationships, ANd Typed Objects), is a middleware solution for publishing XML; object-relational data can be published as XML documents. It can be used by developers who prefer to work in a "pure XML" environment. XPERANTO provides the XML view over the relational data of the existing database, and provides an XML query facility, which translates XML queries into corresponding SQL queries for database, and transforms the results back to XML. XPERANTO contains four major components: the Query Translation, the XML View Services, the XML Schema Generator, and the XML Tagger. In these tools, the success of the conversion is closely related to the quality of the target XML schema onto which a given input relational schema is mapped. However, the mapping from the relational schema

to the XML schema is specified by human experts. Therefore, when a large relational schema and corresponding data need to be translated into XML documents, a significant investment of human effort is required to initially design the target schema. Finally, the work described in [13] requires knowing the catalog contents in order to extract the relational database schema.

There are also some efforts on mapping non-relational models to XML model. The work presented in [14] studies the conversion from XML to ER model and vice versa. XGrammar is used for the notation of XML schema. XGrammar is an extension of the regular tree grammar definition in [15], which uses a six tuple notation to precisely describe content models of XML schema languages. Informally, XGrammar takes the structural specification feature from DTD and RELAX and the data typing feature from XML-Schema. The basic idea of this conversion is to generate XGrammar from a given XML model, then convert XGrammar to ER model or vice versa. Generation of an XML schema from the standard object-oriented design language UML (Unified Modeling Language) model is studied in [4]. They map all elements and data types in XML-Schema to classes annotated with stereotypes that reflect the semantics of the related XML-Schema concept. They use a sequence number for content model elements in order to indicate the order of document types. XML-schemas may contain anonymous groups. They introduce special stereotypes indicating that class represents an anonymous grouping of elements in UML. In addition, the conversion of Relational-to-ER-to-XML has been proposed in [10]. This reconstructs the semantic model, in the form of ER model, from the logical schema capturing user's knowledge, and then converts it to the XML document. However, many-to-many (M:N) and nary relationships are not considered properly. Finally, DB2XML [17] is a tool for transforming data from relational databases into XML documents. DTDs are generated describing the characteristics of the data making the documents self contained and usable as a data exchange format.

Our approach is different from these approaches; we focus on legacy relational databases. We adopt our reverse engineering approach proposed in [2] to extract a semantically rich ER model from the given legacy relational database, and then we convert the ER model to XML schema; we consider M:N and nary relationships.

3. EXTRACTING ER MODEL FROM LEGACY DATABASE: AN OVERVIEW

In this section, we present an overview of the reverse engineering process described in [2]. We will show the results obtained for the following running example.

Example 3.1 Consider the *COMPANY* database shown in Figure 2. This database contains six tables: *EMPLOYEE*, *DEPENDENT*, *PROJECT*, *DEPARTMENT*, *WORKS-ON*, and *DEPT.LOCATIONS*; and each table contains tuples some of which are shown in Figure 2.

The first step is to extract the basic and necessary information from a given legacy relational database. The information includes all candidates and foreign keys found within the relations. It is summarized in the following two tables. Such information is very essential for deriving the target ER model. *CandidateKeys*(*relationname*, *attributename*, *CandidateKey**if*) *ForeignKeys*(*PKrelation*, *PKattribute*, *FKrelation*, *FKattribut*, *Link**if*)

The *CandidateKeys* table contains all possible candidate keys of the relations. The *CandidateKeys*# can be used to keep track of having the same attribute participating in more than one candidate key. The *CandidateKeys* table for the example *COMPANY* database is shown in Figure 3.

Relation	Attribute	Data Type	Constraint		
DEPARTMENT	DNAME	CHAR	PK		
	DNUMBER	INTEGER	PK		
	MGRSSN	CHAR	FK		
	MGRSTARTDATE	DATE			
DEPENDENT	ESSN	CHAR	FK		
	DEPENDENT_NAME	VAR	PK		
	SEX	CHAR			
	BDATE	DATE			
DEPT_LOCATIONS	DNUMBER	INTEGER	FK		
	DLOCATION	VAR	PK		
EMPLOYEE	LNAME	VAR	PK		
	SSN	CHAR	PK		
	BDATE	DATE			
	ADDRESS	VAR			
	SEX	CHAR			
	SALARY	DEC			
	SUPERSSN	CHAR	FK		
DNO	INTEGER	FK			
PROJECT	PNAME	VAR	PK		
	PNUMBER	INTEGER	PK		
	PLLOCATION	VAR			
	DNUM	INTEGER	FK		
WORKS_ON	ESSN	CHAR	FK		
	PNO	INTEGER	FK		
			HOURS	DECIMAL	

Figure 2: An example COMPANY relational database

The *ForeignKeys* table contains all pairs of attributes such that the first attribute is part of a candidate key in a certain relation and the second attribute is part of a foreign key, a representative of the first attribute within any of the relations. *Link#* is to differentiate different foreign keys in the same relation. Foreign keys are numbered so that all attributes within the same foreign key are assigned the same sequence number. The

Relation Name	Attribute Name	Candidate Key #
DEPARTMENT	DNAME	1
DEPARTMENT	DNUMBER	2
DEPENDENT	DEPENDENT_NAME	1
DEPENDENT	BDATE	1
DEPENDENT	ESSN	2
DEPENDENT	DEPENDENT_NAME	2
DEPT_LOCATIONS	DNUMBER	1
DEPT_LOCATIONS	DLOCATION	1
EMPLOYEE	SSN	1
EMPLOYEE	FNAME	2
EMPLOYEE	LNAME	2
PROJECT	PNAME	1
PROJECT	PNUMBER	2
WORKS_ON	ESSN	1
WORKS_ON	PNO	1

Figure 3: Candidate Keys: A list of possible candidate keys of all relations in COMPANY database

ForeignKeys table for the example COMPANY database is shown in Figure 4.

CK Relation	CK Attribute	FK Relation	FK Attribute	Link number	CK Cardinality	FK Cardinality	Note
DEPARTMENT	DNUMBER	DEPT_LOCATIONS	DNUMBER	1	M >= 1	1	is not a candidate
DEPARTMENT	DNUMBER	EMPLOYEE	DNO	1	M >= 1	0 or 1	is not a candidate
DEPARTMENT	DNUMBER	PROJECT	DNUM	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	DEPENDENT	MGRSSN	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	DEPENDENT	ESSN	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	EMPLOYEE	SUPERSSN	1	M >= 0	0 or 1	is not a candidate
EMPLOYEE	SSN	WORKS_ON	ESSN	1	M >= 1	0 or 1	is not a candidate
PROJECT	PNUMBER	WORKS_ON	PNO	1	M >= 1	1	is not a candidate
WORKS_ON	ESSN	DEPARTMENT	MGRSSN	1	M >= 1	0 or 1	is not a candidate
WORKS_ON	ESSN	DEPENDENT	ESSN	1	M >= 1	0 or 1	is not a candidate
WORKS_ON	ESSN	EMPLOYEE	SUPERSSN	1	M >= 0	0 or 1	is not a candidate
WORKS_ON	PNO	PROJECT	PNUMBER	1	1	0 or 1	is a candidate

Figure 4: Foreign Keys: A list of attributes in candidate keys and their corresponding foreign keys with their cardinality

In general, a relation may have a set of candidate keys. One candidate key is chosen as the primary key by checking corresponding foreign keys. For relations that have multiple candidate keys, the primary key is selected to be the candidate key that appears in the first column of *ForeignKeys*. The *Primary Keys* table for the example *COMPANY* database is shown in Figure 5.

The employed reverse engineering process decides on the presence of candidate *key(s)* of a given relation *R* as foreign *key(s)* within *R* itself or any other relation in the relational schema. This leads to unary and binary relationships because while going from the ER model into the relational model, all relationships are mapped into these two types of relationships with some weak entities introduced and converted into relations.

Relation Name	Attribute Name	Primary Key #
DEPARTMENT	DNUMBER	2
DEPENDENT	ESSN	2
DEPENDENT	DEPENDENT_NAME	2
DEPT_LOCATIONS	DNUMBER	1
DEPT_LOCATIONS	DLOCATION	1
EMPLOYEE	SSN	1
PROJECT	PNUMBER	2
WORKS_ON	ESSN	1
WORKS_ON	PNO	1

Figure 5: A list of the primary keys for all relations in the COMPANY database

The information in *ForeignKeys* is used in constructing what is called the *Relational Intermediate Directed (RID) Graph*, which present all possible unary and binary relationships between relations in the given relational schema. In the RID graph, each node represents a relation and two nodes are connected by a link to show that a foreign key in the relation that corresponds to the first node represents the primary key of the relation that corresponds to the second node.

As described in [2], the cardinality of a relationship in the RID graph is determined as follows. A link is directed from R_2 to R_1 to reflect the presence of the primary key of R_1 as a foreign key in R_2 ; so, its cardinality is:

- 1:1 if and only if at most one tuple from R_2 holds the value of the primary key of a tuple from R_1 .
- M:1 if more than one tuple from R_2 hold the value of the primary key of a tuple from R_1 .

The employed process decides also on the minimum and maximum cardinalities at both sides of the link by investigating whether the link is optional or mandatory on each side.

Analyzing the information in Figure 4, which leads to the RID graph, it can be easily observed that it contains some extra information because a foreign key is allowed to play the role of a candidate key and this leads to two *symmetric* and *transitive* references. Such extra information is deleted as described in [2]. The *CandidateKeys* table for the example *COMPANY* database after deleting symmetric and transitive references are shown in Figures 6 and 7, respectively.

Eliminating symmetric and transitive references lead to an optimized RID graph. The optimized RID graph is analyzed further to identify relationships with attributes, M:N and nary relationships, if any. The remaining unary and binary relationships are without attributes, and are represented by direct connections between nodes in the optimized RID graph. They are all classified as 1:1, or M:1. The optimized RID graph for the example *COMPANY* database is shown in Figure 8.

Updated Foreign Key Table (removed symmetric reference)							
CK Relation	CK Attribute	FK Relation	FK Attribute	Link number	CK Cardin...	FK Cardinality	Note
DEPARTMENT	DNUMBER	DEPT_LOCATIONS	DNUMBER	1	M >= 1	1	is not a candidate
DEPARTMENT	DNUMBER	EMPLOYEE	DNO	1	M >= 1	0 or 1	is not a candidate
DEPARTMENT	DNUMBER	PROJECT	DNUM	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	DEPARTMENT	MGRSSN	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	DEPENDENT	ESSN	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	EMPLOYEE	SUPERSSN	1	M >= 0	0 or 1	is not a candidate
EMPLOYEE	SSN	WORKS_ON	ESSN	1	M >= 1	0 or 1	is not a candidate
PROJECT	PNUMBER	WORKS_ON	PNO	1	M >= 1	1	is not a candidate
WORKS_ON	ESSN	DEPARTMENT	MGRSSN	1	M >= 1	0 or 1	is not a candidate
WORKS_ON	ESSN	DEPENDENT	ESSN	1	M >= 1	0 or 1	is not a candidate
WORKS_ON	ESSN	EMPLOYEE	SUPERSSN	1	M >= 0	0 or 1	is not a candidate

Figure 6: The ForeignKeys table after eliminating symmetric references

Updated Foreign Key Table (removed transitive reference)							
CK Relation	CK Attribute	FK Relation	FK Attribute	Link number	CK Cardinality	FK Cardinality	Note
DEPARTMENT	DNUMBER	DEPT_LOCATIONS	DNUMBER	1	M >= 1	1	is not a candidate
DEPARTMENT	DNUMBER	EMPLOYEE	DNO	1	M >= 1	0 or 1	is not a candidate
DEPARTMENT	DNUMBER	PROJECT	DNUM	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	DEPARTMENT	MGRSSN	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	DEPENDENT	ESSN	1	M >= 1	0 or 1	is not a candidate
EMPLOYEE	SSN	EMPLOYEE	SUPERSSN	1	M >= 0	0 or 1	is not a candidate
EMPLOYEE	SSN	WORKS_ON	ESSN	1	M >= 1	0 or 1	is not a candidate
PROJECT	PNUMBER	WORKS_ON	PNO	1	M >= 1	1	is not a candidate

Figure 7: The ForeignKeys table after eliminating transitive references

4. CONVERTING ER MODEL TO XML SCHEMA

In this section, we present the proposed process for translating a conceptual schema (presented as RID graph) into XML schema. The process in pseudo-code is depicted in Algorithm 4.1.

Algorithm 4.1 (ER Model to XML Schema Conversion)

Input: The RID graph

Output: The corresponding XML schema

Step:

Step 1: Translate each entity in the ER model into a complex-type in XML schema.

Step 2: Map each attribute in every entity into a subelement within the corresponding complex-type.

Step 3: Create a root element and insert each entity in the ER model as a subelement with the corresponding complex-type.

Step 4: Use "key" and "keyref" to map each relationship between any two entities.

EndAlgorithm

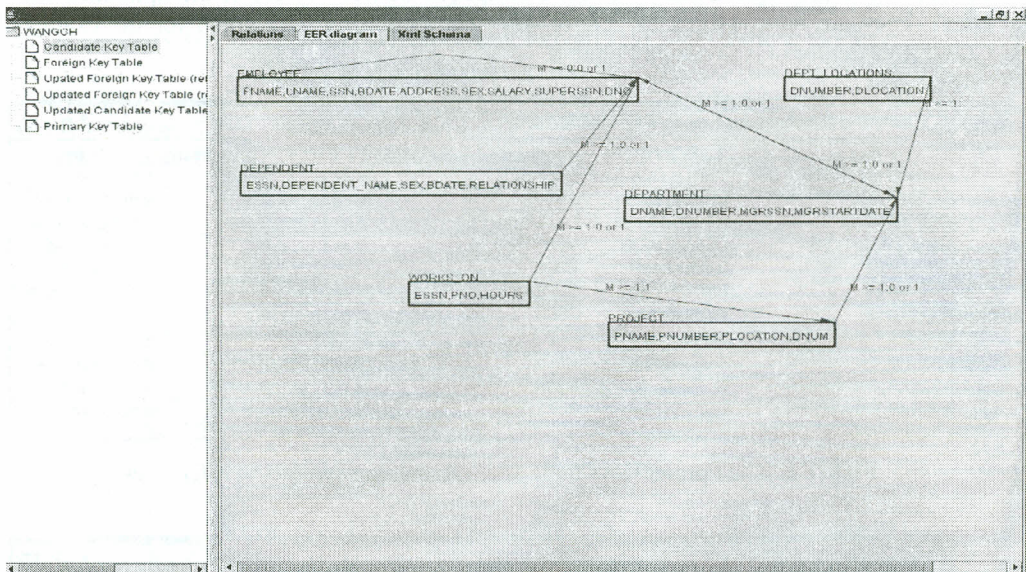


Figure 8: The RID graph of the COMPANY database

In order to convert an ER model to XML schema by Algorithm 4.1, we need to go through the four steps as detailed next:

- Each entity E of the ER model is translated into an XML complex-type of the same name E in the XML schema. In each complex-type E, there is only one empty element. There will be several subelements inside the empty element. For example, the *PROJECT* entity is translated into a complex-type named "PROJECTJRelation". The empty element is called "PROJECT".


```

<complexType name = "PROJECT_Relation" >
  <sequence>
    <element name = "PROJECT" type = "r:PROJECT_Type"
      maxOccurs = "unbounded" />
  </sequence>
</complexType> <complexType name = "PROJECT.Tuple" >
  <sequence>
    .....
  <sequence>
</complexType>

```

The cardinality constraint in the ER model can be explicated by associating two XML built-in *attributes*, also called indicator, namely "minOccurs" and "maxOccurs", with subelements under the XML complex-Type. The "maxOccurs" indicator specifies the maximum number of times a subelement can occur. "maxOccurs" = "unbounded" indicates the element may appear more than once. The "minOccurs" indicator specifies the minimum number of times a subelement can occur. The default value for both the "maxOccurs" and the "minOccurs" attributes is 1. If we want to specify a value only for "minOccurs", it must be either 0 or 1. Similarly, if we want to specify a value only for the "maxOccurs", it should be greater than or equal to 1. If both "minOccurs" and "maxOccurs" are omitted, then the subelement must appear exactly once.

- In step 2 of Algorithm 4.1, each attribute A_i of the entity E is mapped into a subelement of the corresponding complex-type E. For example, the *PROJECT* entity is mapped into a complex-type named "PROJECT_Tuple". Inside the "PROJECT.Tuple" complex-type, there are several subelements such as *PNAME*, *PNUMBER*, *PLOCATION*, and *DNUM*. They are the attributes of the "PROJECT" entity. The XML schema of the *PROJECT* entity is:

```

<complexType name = "PROJECT_Tuple" >
  <sequence>
    <element ref = "nPNAME" />
    <element ref = "nPNUMBER" />
    <element ref = "r:PLOCATION"/>
    <element ref = "r:DNUM"/>
  </sequence>
</complexType>

```

The <sequence> specification in the XML schema captures the sequential semantics of a set of subelements. For instance, in the <sequence> given above, the subelement *PNAME* comes first, followed by *PNUMBER*, and then *PLOCATION*, with *DNUM* at the end. These subelements must appear in instance documents in the same sequential order as they are declared here. XML schema also provides another constructor called <all>, which allows elements to appear in any order, and all the elements must appear once or not at all.

- In step 3 of Algorithm 4.1, each entity is mapped into the XML schema. We first need to create a root element that represents the entire given legacy relational database. We create the root element as a complex-type in the XML schema, and then insert each entity as a subelement of the root element. Next is an example which

contains the six entity objects *DEPARTMENT*, *EPENDENT*, *DEPT.LOCATIONS*, *EMPLOYEE*, *PROJECT*, *WORKS-ON* We call the root element *COMPANY*:

```
<element name = "COMPANY" >
  <complexType>
    <sequence>
      <element name = "r:DEPARTMENT_Ttype" />
      <element name = "r:DEPENDENT_Type" />
      <element name = "r:DEPT_LOCATIONS_Type" />
      <element name = "r:EMPLOYEE_Type" />
      <element name = "r:PROJECT_Type" />
      <element name = "r:WORKS_ON_Type" />
    </sequence>
  </complexType>
</element >
```

Compared to DTD, the XML schema provides a more flexible and powerful mechanism through "key" and "keyref, which share the same syntax as "unique" and also make referential constraints possible in XML documents.

- In step 4 of Algorithm 4.1, we use the elements "key" and "keyref to enforce the uniqueness and referential constraints among the data. According to [20], the "key" element specifies an attribute or element value as a key (unique, non-nullable, and always present) within the containing element in an instance document; and the "keyref element specifies foreign keys, i.e., an attribute or element value correspond to that of an already specified key or unique element. The "key" and "keyref elements replace and extent the capability of "ID", "IDREF" and "IDREFs" in DTD. They are among the great features introduced in XML schema. Also, we can use "key" and "keyref to specify the uniqueness scope and multiple attributes to create the composite keys. Here is an example:

```
<key name = "PROJECTPrimaryKey" >
  <selector xpath = "r:PROJECT/r:PROJECT"/>
  <field xpath = "PNUMBER"/>
</key>
<key name = "WORKS-ON" >
  <selector xpath = "r:WORKS_ON/r:WORKS_ON"/>
  <field xpath = "ESSN"/>
  <field xpath = "PNO"/>
</key>
<keyref name = "PROJECTPNUMBER_WORKS_ONPNOReference" refer =
  "r:PROJECTPrimaryKey" >
  <selector xpath = r:WORKS_ON/r:WORKS_ON"/>
  <field xpath = "PNUMBER"/>
</keyref>
```

In this example, we first specify the primary key for each entity in the ER model. From the *ForeignKeys* table, we know that, *PNUMBER* is the primary key of *PROJECT* entity; *ESSN*, and *PNO* together form a composite primary key of *WORKS-ON* entity. *PNUMBER* is a foreign key of *WORKS-ON*, so we use *Keyref* to specify the foreign key relationship between *PROJECT* and *WORKS-ON* entities.

5. A CLOSER LOOK AT THE DEVELOPED APPROACH AND THE IMPLEMENTED PROTOTYPE

In this section we describe the overall structure of our implementation. The purpose of this section is not to describe details of the code, but to grant the readers an overview of the system. We have two main components: extracting ER model from the given legacy relational database system and converting ER model to XML schema.

5.1 AN OVERVIEW OF THE DESIGN STRUCTURE

The block diagram of our framework is shown in Figure 9. It consists of four major components: a given relational legacy database, REGNM (short for reverse engineering module), COVM (short for conversion module), and XML schema. The purpose of the REGNM is to extract the ER model from the legacy database by running the algorithms proposed in [2]. Then COVM runs the conversion algorithm to convert the extracted ER model to a XML Schema. The result XML schema is displayed on the graphical output interface. Figure 2 shows the input legacy relational database, and Figure 8 shows the RID graph (ER diagram) obtained as output.

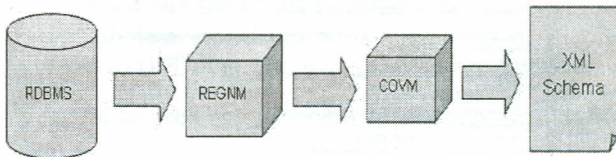


Figure 9: The structure of the framework

The prototype has been implemented using Java. In addition to the fact that we are familiar with Java, reasons for choosing Java include:

- It is an Object-Oriented language, and hence it is easy to program in Java.
- We can use JDBC driver to connect to the database, also there are some useful functions we can use for doing operations in the database.
- We can use JDOM to obtain the XML schema.

5.2 EXAMPLES AND RESULTS

We have tested our algorithms on the contents the *COMPANY* database in Example 3.1 and the *SAMPLE* database in DB2. Here, we only show the output related to the *COMPANY* database because the *SAMPLE* database is quite large; the corresponding generated output is long.

Now let us take a look at the results. We divide the testing process into the following steps:

1. Find all candidate keys of the *COMPANY* database; the result is shown a in Figure 3.
2. Find all candidate foreign keys of the *COMPANY* database; the result is shown in Figure 4, and it is used to construct the initial RID graph.
3. Find the cardinality for each relationship in the *COMPANY* database; the result is shown in Figure 4.
4. Eliminate the symmetric and transitive references; the results are shown in Figure 6 and Figure 7, respectively.
5. Find *M:M* and *n-ary* relationships; the obtained optimized RID graph is shown in Fig. 8.

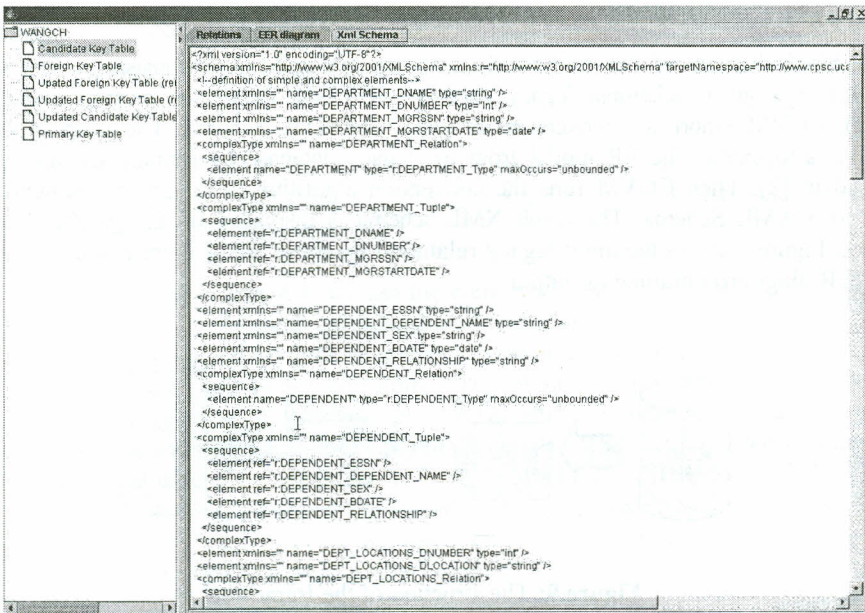


Figure 10: The output XML schema

6. Convert the ER model to the XML Schema; the output XML schema screen is shown in Fig. 10. The complete XML schema is given the appendix.

5.3. DISCUSSION OF THE RESULTS

In the last section, we run these algorithms to extract an ER model from the example legacy relational database - *COMPANY*, then the ER model is converted to XML schema by using Algorithm 4.1. The optimized RID graph for the *COMPANY* database supports the correctness, effectiveness and applicability of our approach.

We also test our approach on the contents of the *SAMPLE* database in DB2 and the *NorthWind* database in MS Access 2000, we neglected the catalog contents for both databases in order to test the reverse engineering process. It takes around 5 minutes for the *SAMPLE* database, and more than an hour for the *NorthWind* database. This is normal because we expect the time to increase when the size of the tested database increases. Compared to the

SAMPLE database, it takes much longer time to test the *NorthWind*. The main reason is that *NorthWind* contains 8 tables, many attributes in some of tables, and a lot of records in each table. Most of the time is spent on analyzing the contents of the tables and deriving the ER model. Even if a human is asked to do the same job, the process becomes unmanageable manually as the size and complexity of the database increases.

To summarize, the proposed framework could automatically extract the ER model from the given legacy relational database, and then transfer it to an XML schema. The framework consists of these major components:

Data layer, which is a legacy relational database that stores all the data to be analyzed and converted into XML.

Reverse engineering layer, which extracts an ER model from the input database by applying some reverse engineering techniques.

Transformation layer, which transfers the ER model to XML schema.

Graphical output layer, which shows the result for each step (i.e., foreign key table, candidate key table, primary key table, all tables in the given database, a **RID** graph (equivalent to the ER diagram), and an XML schema, etc).

Undoubtedly, reconstructing an ER model from a legacy database, and writing an XML schema file both are heavy and tedious jobs, especially for a large real application. The users could be relieved of this heavy load by using our framework. On the other hand, the users' knowledge could also be involved in this system. However, compared to reconstructing an ER model and writing a long XML schema file from scratch, the human's mental workload is greatly reduced with our framework.

Our framework presented in this paper has the following advantages compared to the work described in [12], where the authors show how to obtain a DTD for data whose structure is described by a conceptual data model. In brief, they present the translation of all constructs of the ER model to DTDs and integrate them into an algorithm.

- Our framework could be used not only for a normal relational database system, but also for a legacy relational database system.
- We choose the XML schema instead of the **DTD**. The XML schema provides a more flexible and powerful mechanism than the DTD. We can easily present each entity in the ER model by using XML complex-Type. And also we can use "key" and "keyref" to declare the attributes uniqueness, composite key, and referential constraint.

In our framework, we provide a user-friendly graphical user interface and also the output can be visualized in a user-friendly manner, i.e., our prototype gives users a direct visualization of the output obtained from each phase of the process.

The expected human workload is considerably reduced compared to the approach described in [12].

CONCLUSION

In this paper, we presented a novel approach to extract an ER model from a legacy relational database, and then convert the ER model to a corresponding XML schema; i.e., by applying reverse engineering followed by forward engineering. We preserve as much information as we can from the given relational schema to the XML schema. Our approach not only works for the commercial relational databases but also for the legacy relational databases. We use the XML schema instead of the DTD schema; the advantages of this is that we can use a complex-type to represent each relational table; "key" and "keyref" are great features introduced in XML schema. They replace and extend the capability of "ID", and "IDREF" and "IDREFS" in DTD. We use "key" and "keyref" to specify the relationship between tables, the uniqueness scope and multiple attributes to create the composite keys. We can also determine $M:N$ and n -ary relationships, so we produce a XML schema and XML documents for the data stored in databases without knowing anything about the catalog information. Currently, we are working on improving the prototype to provide flexible visual querying facility by allowing the user to choose from the displayed RID graph the tables and even the attributes to be displayed in XML format.

REFERENCES

- [1] S. Abiteboul, D. Quass, J. McHugh, J. Widom, and J. L. Wiener. "The Lorel Query Language for Semistructured Data," Department of Computer Science, Stanford University, Stanford, <http://www-db.Stanford.edu/lore/pubs/lore196.pdf>.
- [2] Reda Alhajj, "Extracting the Extended Entity-Relationship Model from a legacy Relational Database, *information Systems*, Vol.28, No.6, pp.597-618, 2003.
- [3] A. Bonifati, D. Lee. "Technical Survey of XML Schema and Query Languages," *Technical report, UCLA Computer Science Department*, June 2001.
- [4] G. Booch, M. Christerson, M. Fuchs, and J. Koistinen. "UML for XML Schema Mapping Specification," <http://www.rational.com/media/uinl/resources/media/iinljxinlscliema33.pdf>.
- [5] R. Bourret. "XML and Databases," Web page, 2003, <http://www.rpbourret.com/xml/XMLAndDatabases.htm>.
- [6] M. Carey, D. Florescu, Z. Ives, Y. Lu, J. Schanmugasundaram, E. Shekita, and S. Subramanian. "XPER-ATO: Publishing Object-Relational Data as XML," *Proceedings of the International Workshop on Web and Databases*, May 2000.
- [7] J. Cheng and J. Xu. *IBM DB2 XML Extender*, IBM Silcom Valley, February, 2000.
- [8] R. Elmasri, S. B. Navathe. *Fundamentals of Database Systems*, Reading, Addison-Wesley, 4th Edition, 2004.
- [9] M. F. Fernandez, W. C. Tan, and D. Suciu. "SilkRoute: Trading between Relational and XML," *Proceedings of the International Conference on World Wide Web*, May 2000.
- [10] J. Fong, F. Pang, and C. Bloor. "Converting Relational Database into XML Document," *Proceedings of International Workshop on Electronic Business Hubs*, pp61-65, Sep. 2001.
- [11] G. Kappel, E. Kapsammer, S. Rausch-Schott, and W. Retschitzegger. "X-Ray - Towards Integrating XML and Relational Database Systems," *Proceedings of the International Conference on Conceptual Modeling*, pp. 339-353, Salt Lake City, UT, Oct. 2000.

- [12] C. Kleiner, U. W. Lipeck. "Automatic Generation of XML DTDs from Conceptual Database Schemas," *University of Hannover, Germany*, Sept 2001.
- [13] D. Lee, M. Mani, F. Chiu, and W. W. Chu, "Nesting based Relational-to-XML Schema Translation," *Proceedings of the International Workshop on Web and Databases*, May 2001.
- [14] M. Mani, D. Lee, and Richard R. Muntz. "Semantic Data Modeling using XML Schemas," *Department of Computer Science, University of California, Los Angeles*, 2001.
- [15] M. Mani, D. Lee, and Richard R. Muntz. "Taxonomy of XML Schema Language using Formal Language Theory," *In Extreme Markup Languages*, Montreal, Canada, August, 2001. <http://www.cs.ucla.edu/~dongwon/paper>
- [16] M. Murata. "RELAX (REGular LAnguage description for XML)," Web page, 2000, <http://www.xml.gr.jp/relax/>.
- [17] V. Turau, "Making Legacy Data Accessible for XML applications," 1999, <http://www.informatik.fh-wiesbaden.de/turau/ps/legacy.pdf>.
- [18] *Extensible Markup Language (XML)* <http://www.w3sch.ools.com/XML/>
- [19] *Document Type Definition (DTD)* <http://www.w3schools.com/DTD/>
- [20] *XML Schema* <http://www.w3schools.com/schema/>
- [21] *Xpath* <http://www.w3schools.com/xpath/>
- [22] *DOM* <http://www.w3schools.com/DOM/>
- [23] *SOAP* <http://www.w3schools.com/SOAP/>
- [24] *XML-QL: A Query Language for XML* <http://www.w3.org/TR/N0TE-xml-ql/>
- [25] *XQuery 1.0: An XML Query Language* <http://www.w3.org/TR/xquery/>

Received: 23 May 2004

Accepted: 08 December 2004