

# NeT & CoT: Translating Relational Schemas to XML Schemas using Semantic Constraints

Dongwon Lee

Murali Mani

Frank Chiu

Wesley W. Chu

UCLA / CSD

{dongwon,mani,frankchiu,wwc}@cs.ucla.edu

## Abstract

Two algorithms, called NeT and CoT, to translate relational schemas to XML schemas using various semantic constraints are presented. The XML schema representation we use is a language-independent formalism named *XSchema*, that is both precise and concise. A given *XSchema* can be mapped to a schema in any of the existing XML schema language proposals. Our proposed algorithms have the following characteristics: (1) NeT derives a nested structure from a flat relational model by repeatedly applying the *nest* operator on each table so that the resulting XML schema becomes hierarchical, and (2) CoT considers not only the structure of relational schemas, but also semantic constraints such as inclusion dependencies during the translation - it takes as input a relational schema where multiple tables are interconnected through inclusion dependencies and converts it into a *good XSchema*. To validate our proposals, we present experimental results using both real schemas from the UCI repository and synthetic schemas from TPC-H.

## 1 Introduction

XML [3] is rapidly becoming one of the most widely adopted technologies for information exchange and representation on the World Wide Web. With XML emerging as *the* data format of the Internet era, there is a substantial increase in the amount of data encoded in XML. However, the majority of everyday data is still

---

*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 28th VLDB Conference,  
Hong Kong, China, 2002

stored and maintained in relational databases. Therefore, we expect the needs to convert such relational data into XML documents will grow substantially as well. In this paper, we study the problems in this conversion. Especially, we are interested in finding XML schema<sup>1</sup> (e.g., DTD [3], RELAX-NG [5], XML-Schema [22]) that *best* describes the existing relational schema. Having an XML schema that precisely describes the semantics and structures of the original relational data is important to further maintain the converted XML documents in future.

At present, there exist several tools that enable the composition of XML documents from relational data, such as XML Extender from IBM<sup>2</sup>, XML-DBMS<sup>3</sup>, DB2XML [23], SilkRoute [8], and XPERANTO [4]. In these tools, the success of the conversion is closely related with 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 large amount of relational schemas and data need to be translated into XML documents, a significant investment of human effort is required to initially design target schemas. To make matters worse, in the context of merging legacy relational data to existing XML documents, devising a good XML schema that does not violate existing structures and constraints is a non-trivial task. Being able to *automatically* infer a precise XML schema out of relational schema would be very useful in such settings.

In this paper, therefore, we are interested in finding a method that can infer the *best* XML schema from the given relational schema automatically. We particularly focus on two aspects of the translation: (1) **Structural aspect:** We want to find the most intuitive and precise XML schema structure from the given relational schema. We especially try to use the hidden

---

<sup>1</sup>We differentiate two terms - XML schema(s) and XML-Schema. The former refers to a general term for schema in XML model while the latter refers to one particular kind of XML schema language proposed by W3C [22].

<sup>2</sup><http://www-4.ibm.com/software/data/db2/extenders/xmllex/>

<sup>3</sup><http://www.rpbouret.com/xmldbms/index.htm>

characteristics of data using *nest* operator, and (2) **Semantic aspect**: During the translation, we want to use semantic constraints that could be either acquired from database directly or provided by human experts explicitly.

We first present a straightforward relational to XML translation algorithm, called *Flat Translation* (FT). Since FT maps the flat relational model to the flat XML model in a one-to-one manner, it does not utilize the regular expression operators (e.g., “\*”, “+”) supported in the content models of XML. Then, we present our first proposal called *Nesting-based Translation* (NeT), to remedy the problems found in FT. NeT derives nested structures from a flat relational model by the use of the *nest* operator so that the resulting XML schema is more intuitive and precise than otherwise. Although NeT infers hidden characteristics of data by nesting, it is only applicable to a single table at a time. Therefore, it is unable to capture a correct “big picture” of relational schema where many tables are interconnected. To remedy this problem, we present the second proposal called *Constraints-based Translation* (CoT); CoT considers inclusion dependencies during the translation. Such constraints can be acquired from database through ODBC/JDBC interface or provided by human experts who are familiar with the semantics of the relational schema being translated. CoT is capable of generating a more intuitive XML schema than what NeT. Figure 1 illustrates the overview of our approach.

**Related Work**: There have been different approaches for the conversion from relational model to XML model, such as XML Extender from IBM, XML-DBMS, SilkRoute [8], XPERANTO [4], DB2XML [23] and NeT [14]. All the above tools (except NeT) require the user to specify the mapping from the given relational model to the XML model. In XML Extender, the user specifies the mapping through a language such as DAD or XML Extender Transform Language. In XML-DBMS, a template-driven mapping language is provided to specify the mappings. SilkRoute provides a declarative query language (RXL) for viewing relational data in XML. XPERANTO uses XML query language for viewing relational data in XML. Note that in SilkRoute and XPERANTO, the user has to specify the query in the appropriate query language. DB2XML uses an algorithm similar to FT (and hence suffers from similar problems). NeT does not require user-input for mapping the relational model to XML model, however it does not use semantic constraints specified in the relational model.

There also have been work in mapping from non-relational models to XML model, and XML to relational and other models. In [20, 21], the authors study the conversion from XML to relational models. [15] studies the conversion from XML to ER model and vice versa. Generation of an XML schema from a

UML model is studied in [18]. Given a set of XML documents, generating an XML schema for them is studied in [10].

**Roadmap**: In Section 2, we first present formalisms to represent relational as well as XML schemas in a language independent notation. In Section 3, we propose NeT algorithm that uses the *nest* operator developed in the nested relational model community. In Section 4, we propose an improved CoT algorithm that considers various semantic constraints during the translation to generate a better XML schema, in addition to applying *nest* operations for each table. In Section 5, we discuss some issues related to correctness and goodness of the schema that NeT and CoT generate. In Section 6, we report the results from our experimentations. Finally, concluding remarks and future directions are discussed in Section 7.

## 2 Input & Output Models

We first briefly define the input and output models for the translation. In relational databases, schema is typically created by SQL DDL (e.g., CREATE) statements. Therefore, by examining such DDL statements, one can find out the original schema information. Even if such DDL statements are not available, one can still infer the schema information - table and column names, key and foreign key information, etc - by querying the database through an ODBC/JDBC interface or by examining the database directly. In this paper, regardless of how one acquired the schema information, we assume that the schema information is encoded in a vector  $\mathbb{R}$  defined below.

Let us assume the existence of a set  $\hat{T}$  of table names, a set  $\hat{C}$  of column names and a set  $\hat{b}$  of atomic base types defined in the standard SQL (e.g., integer, char, string). When name collision occurs, a column name  $c \in \hat{C}$  is qualified by a table name  $t \in \hat{T}$  using the “[ ]” notation (e.g.,  $t[c]$ ).

**Definition 1 (Relational Schema)** A relational schema is denoted by 4-tuple  $\mathbb{R} = (T, C, P, \Delta)$ , where:

- $T$  is a finite set of table names in  $\hat{T}$ ;  $C$  is a function from a table name  $t \in T$  to a set of column names  $c \in \hat{C}$ ,
- $P$  is a function from a column name  $c$  to its column type definition: i.e.,  $P(c) = \alpha$ , where  $\alpha$  is a 5-tuple  $(\tau, u, n, d, f)$ , where  $\tau \in \hat{b}$ ,  $u$  is either “ $v$ ” (unique) or “ $\neg v$ ” (not unique),  $n$  is either “?” (nullable) or “ $\neg?$ ” (not nullable),  $d$  is a finite set of valid domain values of  $c$  or  $\epsilon$  if not known, and  $f$  is a default value of  $c$  or  $\epsilon$  if not known, and
- $\Delta$  is a finite set of relational integrity constraints that can be either retrieved from databases directly or provided by human experts.  $\square$

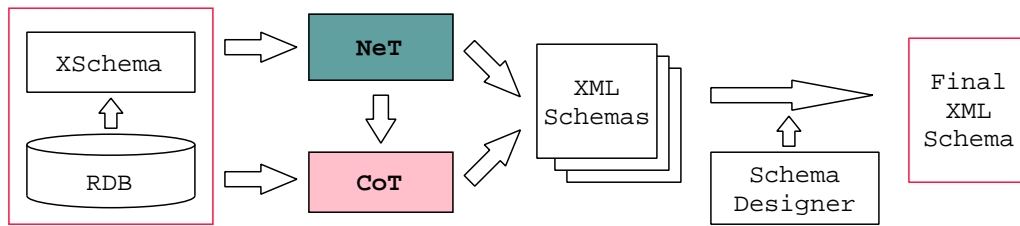


Figure 1: Overview of our approach. Two algorithms, NeT and CoT, can be used independently or jointly; (1)  $XSchema \rightarrow NeT \rightarrow XML\ Schemas$ , (2)  $XSchema \rightarrow CoT \rightarrow XML\ Schemas$ , or (3)  $XSchema \rightarrow NeT \rightarrow CoT \rightarrow XML\ Schemas$ .

$ \begin{aligned} T &= \{student, professor\} \\ C(student) &= \{Sname, Advisor, Course\} \\ C(professor) &= \{Pname, Age\} \\ P(Sname) &= (string, \neg v, \neg?, \epsilon, \epsilon) \\ P(Advisor) &= (string, \neg v, \neg?, \epsilon, "J.Smith") \\ P(Course) &= (string, \neg v, \neg?, \epsilon, \epsilon) \\ P(Pname) &= (string, v, \neg?, \epsilon, \epsilon) \\ P(Age) &= (integer, \neg v, \neg?, \epsilon, \epsilon) \\ \Delta &= \{\{Sname, Advisor, Course\} \xrightarrow{key} student, \\ &\quad Pname \xrightarrow{key} professor, Advisor \subseteq Pname\} \end{aligned} $	<table border="1" style="border-collapse: collapse; width: 100%; text-align: center;"> <thead> <tr> <th colspan="3">student:</th> </tr> <tr> <th>Sname</th> <th>Advisor</th> <th>Course</th> </tr> </thead> <tbody> <tr><td>John</td><td>Prof. Muntz</td><td>Database Systems</td></tr> <tr><td>John</td><td>Prof. Muntz</td><td>Data Mining</td></tr> <tr><td>John</td><td>Prof. Chu</td><td>Database Systems</td></tr> <tr><td>John</td><td>Prof. Chu</td><td>Data Mining</td></tr> <tr><td>John</td><td>Prof. Muntz</td><td>Queueing Theory</td></tr> <tr><td>John</td><td>Prof. Zaniolo</td><td>Logic</td></tr> </tbody> </table> <table border="1" style="border-collapse: collapse; width: 100%; text-align: center;"> <thead> <tr> <th colspan="2">professor:</th> </tr> <tr> <th>Pname</th> <th>Age</th> </tr> </thead> <tbody> <tr><td>Prof. Muntz</td><td>60</td></tr> <tr><td>Prof. Chu</td><td>55</td></tr> <tr><td>Prof. Zaniolo</td><td>-</td></tr> <tr><td>Prof. Parker</td><td>49</td></tr> </tbody> </table>	student:			Sname	Advisor	Course	John	Prof. Muntz	Database Systems	John	Prof. Muntz	Data Mining	John	Prof. Chu	Database Systems	John	Prof. Chu	Data Mining	John	Prof. Muntz	Queueing Theory	John	Prof. Zaniolo	Logic	professor:		Pname	Age	Prof. Muntz	60	Prof. Chu	55	Prof. Zaniolo	-	Prof. Parker	49
student:																																					
Sname	Advisor	Course																																			
John	Prof. Muntz	Database Systems																																			
John	Prof. Muntz	Data Mining																																			
John	Prof. Chu	Database Systems																																			
John	Prof. Chu	Data Mining																																			
John	Prof. Muntz	Queueing Theory																																			
John	Prof. Zaniolo	Logic																																			
professor:																																					
Pname	Age																																				
Prof. Muntz	60																																				
Prof. Chu	55																																				
Prof. Zaniolo	-																																				
Prof. Parker	49																																				

Table 1: Example relational schema and data.

**Example 1.** Consider two tables `student(Sname, Advisor, Course)` and `professor(Pname, Age)` where keys are underlined, and `Advisor` is a foreign key referencing `Pname` column. The column `Age` is an integer type, while the rest of the columns are string types. Also `Age` may be null. When student’s advisor has not yet been decided, professor “J. Smith” will be the initial advisor. Student can have many advisors and take zero or more courses. The corresponding relational schema and data fragment is given in Table 1.  $\square$

Next, let us define the output model. Lately, there have been about a dozen competing XML schema language proposals. Although XML-Schema is being shaped by W3C and will replace DTD soon, it is likely that different applications will choose different XML schema languages that best suit their particular purposes. Therefore, instead of choosing one language proposal, we formalize a core set of important features into a new notion of *XSchema* and use it as our output modeling language. The benefits of such formalization is that it is both concise and precise. More importantly, it breaks the tie between the translation algorithm that we are developing and the final schema language notations. Informally, *XSchema* borrows structural features from DTD and RELAX-NG, and data types and constraint specification features from XML-Schema. From a formal language and database perspective [17], *XSchema* is a local tree grammar ex-

tended with attribute, datatype and constraint specifications.

Starting from the notations in [7], we define *XSchema* below. We first assume the existence of a set  $\hat{E}$  of element names, a set  $\hat{A}$  of attribute names and a set  $\hat{\tau}$  of atomic data types defined in [1] (e.g., ID, IDREF, string, integer, date, etc). When needed, an attribute name  $a \in \hat{A}$  is qualified by the element names using the *path expression* notation  $e_1.e_2 \cdots e_n.a$ , where  $e_i \in \hat{E}, 1 \leq i \leq n$ .

**Definition 2 (XSchema)** An *XSchema* is denoted by 6-tuple  $\mathbb{X} = (E, A, M, P, r, \Sigma)$ , where:

- $E$  is a finite set of element names in  $\hat{E}$ ;  $A$  is a function from an element name  $e \in E$  to a set of attribute names  $a \in \hat{A}$ ,
- $M$  is a function from an element name  $e \in E$  to its element type definition: i.e.,  $M(e) = \alpha$ , where  $\alpha$  is a regular expression:  $\alpha ::= \epsilon \mid \tau \mid \alpha + \alpha \mid \alpha, \alpha \mid \alpha^? \mid \alpha^* \mid \alpha^+$ , where  $\epsilon$  denotes the empty element,  $\tau \in \hat{\tau}$ , “+” for the union, “,” for the concatenation, “ $\alpha^?$ ” for zero or one occurrence, “ $\alpha^*$ ” for the Kleene star, and “ $\alpha^+$ ” for “ $\alpha, \alpha^*$ ”,
- $P$  is a function from an attribute name  $a$  to its attribute type definition: i.e.,  $P(a) = \beta$ , where  $\beta$  is a 4-tuple  $(\tau, n, d, f)$ , where  $\tau \in \hat{\tau}$ ,  $n$  is either “?” (nullable) or “ $\neg?$ ” (not nullable),  $d$  is a finite

set of valid domain values of  $a$  or  $\epsilon$  if not known, and  $f$  is a default value of  $a$  or  $\epsilon$  if not known, and

- $r \subseteq E$  is a finite set of root elements;  $\Sigma$  is a finite set of integrity constraints for XML model  $\square$

Translation from *XSchema* to the actual XML schema language notations is relatively straightforward and not discussed further in this paper. It is worthwhile to note, however, that depending on the chosen XML schema language, some of the features specifiable in *XSchema* might not be translatable at the end. For instance, any “non-trivial type” or composite key information would be lost if one decides to use DTD as the final XML schema language.

### 3 Flat Translation and Nesting-based Translation

XML model uses two basic building blocks to construct XML documents – attribute and element. A few basic characteristics inherited from XML model include: (1) the attributes of a node are not ordered, while the child elements of a node are ordered, (2) both support data types as specified in [1], and (3) elements can express multiple occurrences better than attributes. The detailed capabilities of those, however, vary depending on the chosen XML schema language. In translating  $\mathbb{R}$  to  $\mathbb{X}$ , therefore, one can either use attribute or element in  $\mathbb{X}$  to represent the same entity in  $\mathbb{R}$  (e.g, a column with string type in  $\mathbb{R}$  can be translated to either attribute or element with string type in  $\mathbb{X}$ ).

To increase the flexibility of the algorithms, we assume that there are two modes – **attribute-oriented** and **element-oriented**. Depending on the mode, an algorithm can selectively translate an entity in  $\mathbb{R}$  to either attribute or element if both can capture the entity correctly. However, if the chosen XML schema language requires attribute or element for an entity (e.g., a key column in  $\mathbb{R}$  needs to be translated to an attribute with type ID in  $\mathbb{X}$ ), we assume that the algorithm follows the limitations.

#### 3.1 Flat Translation

The simplest translation method is to translate (1) tables in  $\mathbb{R}$  to elements in  $\mathbb{X}$  and (2) columns in  $\mathbb{R}$  to attributes (in attribute-oriented mode) or elements (in element-oriented mode) in  $\mathbb{X}$ . These two modes are analogous except that element-oriented mode adds additional order semantics to the resulting schema. Since  $\mathbb{X}$  represents the “flat” relational tuples faithfully, this method is called *Flat Translation* (FT). The general procedure of the **Flat Translation** is omitted in the interest of space and can be found in [14]. One example is shown in Appendix (Example 8).

FT is a simple and effective translation algorithm, but it has some problems. As the name implies, FT

translates the “flat” relational model to a “flat” XML model in a one-to-one manner. The drawback of FT is that it does not utilize several basic “non-flat” features provided by XML for data modeling such as representing *repeating sub-elements* through regular expression operators (e.g., “\*”, “+”). We remedy this problem in the NeT algorithm below.

#### 3.2 Nesting-based Translation

To remedy the problems of FT, one needs to utilize various *element content models* of XML. Towards this goal, we propose to use the *nest* operator [12]. Our idea is to find a “best” element content model that uses  $\alpha^*$  or  $\alpha^+$  using the *nest* operator. First, let us define the *nest* operator. Informally, for a table  $t$  with a set of columns  $C$ , *nesting* on a non-empty column  $X \in C$  collects all tuples that agree on the remaining columns  $C - X$  into a set<sup>4</sup>. Formally,

**Definition 3 (Nest)** [12]. Let  $t$  be a  $n$ -ary table with column set  $C$ , and  $X \in C$  and  $\overline{X} = C - X$ . For each  $(n - 1)$ -tuple  $\gamma \in \Pi_{\overline{X}}(t)$ , we define an  $n$ -tuple  $\gamma^*$  as follows:  $\gamma^*[\overline{X}] = \gamma$ , and  $\gamma^*[X] = \{\kappa[X] \mid \kappa \in t \wedge \kappa[\overline{X}] = \gamma\}$ . Then,  $nest_X(t) = \{\gamma^* \mid \gamma \in \Pi_{\overline{X}}(t)\}$ .  $\square$

After  $nest_X(t)$ , if column  $X$  has only a set with “single” value  $\{v\}$  for all the tuples, then we say that **nesting failed** and we treat  $\{v\}$  and  $v$  interchangeably (i.e.,  $\{v\} = v$ ). Thus when nesting failed, the following is true:  $nest_X(t) = t$ . Otherwise, if column  $X$  has a set with “multiple” values  $\{v_1, \dots, v_k\}$  with  $k \geq 2$  for at least one tuple, then we say that **nesting succeeded**. The general procedure for nesting is given in Table 9 of Appendix.

**Example 2.** Consider a table  $R$  in Table 2. Here we assume that the columns  $A, B, C$  are non-nullable. In computing  $nest_A(R)$  at (b), the first, third, and fourth tuples of  $R$  agree on their values in columns  $(B, C)$  as (a, 10), while their values of the column  $A$  are all different. Therefore, these different values are grouped (i.e., nested) into a set  $\{1,2,3\}$ . The result is the first tuple of the table  $nest_A(R) - (\{1,2,3\}, a, 10)$ . Similarly, since the sixth and seventh tuples of  $R$  agree on their values as (b, 20), they are grouped to a set  $\{4,5\}$ . In computing  $nest_B(R)$  at (c), there are no tuples in  $R$  that agree on the values of the columns  $(A, C)$ . Therefore,  $nest_B(R) = R$ . In computing  $nest_C(R)$  at (d), since the first two tuples of  $R - (1, a, 10)$  and  $(1, a, 20)$  – agree on the values of the columns  $(A, B)$ , they are grouped to  $(1, a, \{10,20\})$ . Nested tables (e) through (j) are constructed similarly.  $\square$

Since the *nest* operator requires scanning of the entire set of tuples in a given table, it can be quite expensive. In addition, as shown in Example 2, there are various ways to nest the given table. Therefore, it is

<sup>4</sup>Here, we only consider single attribute nesting.

	A	B	C
#1	1	a	10
#2	1	a	20
#3	2	a	10
#4	3	a	10
#5	4	b	10
#6	4	b	20
#7	5	b	20

(a)  $R$

A <sup>+</sup>	B	C
{1,2,3}	a	10
1	a	20
4	b	10
{4,5}	b	20

(b)  $nest_A(R)$

A	B	C
1	a	10
1	a	20
2	a	10
3	a	10
4	b	10
4	b	20
5	b	20

(c)  $nest_B(R) = R$

A	B	C <sup>+</sup>
1	a	{10,20}
2	a	10
3	a	10
4	b	{10,20}
5	b	20

(d)  $nest_C(R)$

A <sup>+</sup>	B	C
{1,2,3}	a	10
1	a	20
4	b	10
{4,5}	b	20

(e)  $nest_B(nest_A(R)) = nest_C(nest_A(R))$

A <sup>+</sup>	B	C <sup>+</sup>
1	a	{10,20}
{2,3}	a	10
4	b	{10,20}
5	b	20

(f)  $nest_A(nest_C(R))$

A	B	C <sup>+</sup>
1	a	{10,20}
2	a	10
3	a	10
4	b	{10,20}
5	b	20

(g)  $nest_B(nest_C(R))$

A <sup>+</sup>	B	C
{1,2,3}	a	10
1	a	20
4	b	10
{4,5}	b	20

(h)  $nest_C(nest_B(nest_A(R))) = nest_B(nest_C(nest_A(R)))$

A <sup>+</sup>	B	C <sup>+</sup>
1	a	{10,20}
{2,3}	a	10
4	b	{10,20}
5	b	20

(i)  $nest_B(nest_A(nest_C(R))) = nest_A(nest_B(nest_C(R)))$

Table 2: A relational table  $R$  and its various nested forms. Column names containing a set after nesting (i.e., nesting succeeded) are appended by “+” symbol.

important to find an efficient way (that uses the  $nest$  operator minimum number of times) of obtaining an acceptable element content model.

First, to find out the total number of ways to nest, let us use the following two properties [12]:

$$\begin{aligned}
 P1: \quad nest_A(nest_B(t)) &\neq nest_B(nest_A(t)) \\
 P2: \quad nest_X(nest_{All_L}(t)) &= nest_{All_L}(t), \\
 &\text{if } X \in L.
 \end{aligned}$$

Here,  $nest_{All_L}(t)$  represents performing nesting on the columns on  $L$  in the order as shown below:

$$nest_{All_L} = \langle c_1, c_2, \dots, c_n \rangle (t) = nest_{c_1}(nest_{c_2}(\dots(nest_{c_n}(t))))$$

P1 states that “commutativity” of nesting does not hold in general and P2 states that nesting along the same column repeatedly has the property of “idempotency”. Using the two properties, the number of permutations to nest tables can be described as follows:

**Remark 1** Using the falling factorial power notation “ $x$  to the  $m$  falling” as  $x^{\underline{m}}$  in [11], the total number of different nestings  $T$  for a table with  $n$  columns is given by:  $T = \sum_{k=1}^n n^{\underline{k}}$   $\square$

According to Remark 1, there are 15 meaningful ways of nesting along the columns  $A, B, C$  in Table 2. Then, the next questions are (1) how to decrease  $T$  by avoiding unnecessary nesting, and (2) which nesting should be chosen as *the* translation. To answer these questions, let us first describe a few useful properties of the  $nest$  operator as follows:

**Lemma 1.** Suppose we are nesting a table  $t$ , with  $C(t) = \{C_t\}$ , and with candidate keys  $K_1, K_2, \dots, K_n$ , where  $K_1, K_2, \dots, K_n \subseteq C_t$ . Applying the  $nest$  operator on a column  $X \notin K_1 \cap K_2 \cap \dots \cap K_n$  yields no changes.  $\blacksquare$

**Lemma 2.** For any nested table  $nest_X(t)$ ,  $\overline{X} \rightarrow X$  holds.  $\blacksquare$

Lemma 2 states that after applying the  $nest$  operator of column  $X$ , the remaining columns  $\overline{X}$  become a super key. Fischer et al. [9] have proved that functional dependencies are preserved against nesting as follows:

**Lemma 3.** [9] If  $X, Y, Z$  are columns of  $t$ , then:  $t : X \rightarrow Y \implies nest_Z(t) : X \rightarrow Y$   $\blacksquare$

Now, we arrive at the following useful property:

**Lemma 4.** For a table  $t$  with  $n$  columns and  $m$  columns that participate in all candidate keys ( $m \leq n$ ), the maximum number of nestings is  $\sum_{k=1}^m m^{\underline{k}}$   $\blacksquare$

Using Lemma 4, one can avoid unnecessary nestings illustrated as follows:

**Example 3.** Consider a table  $R$  in Table 2 again. Suppose attributes  $A$  and  $C$  constitute a key for  $R$ . Since nesting on the same column repeatedly is not useful by property P2 there is no need to construct, for instance,  $nest_A(nest_A(R))$ . Since nesting on a non-key column is not useful by Lemma 1, nesting along column  $B$  (e.g.,  $nest_B(R)$  at (c)) can be avoided. Furthermore, the functional dependency (i.e.,  $AC \xrightarrow{key} R = AC \rightarrow \overline{AC} = AC \rightarrow B$ ) persists after nesting on either column  $A$  or  $C$  by Lemma 3. Consequently, one needs to construct only the following nested tables:  $nest_A(R)$  at (b),  $nest_C(R)$  at (d),  $nest_C(nest_A(R))$  at (e),  $nest_A(nest_C(R))$  at (f).  $\square$

As we have shown, when candidate key information is available, the number of nestings to be performed can be reduced. However, when such information is not known, the  $nest$  operator must be applied for all possible combinations in Remark 1. After applying the

*nest* operator to the given table repeatedly, there can be still several nested tables where nesting succeeded. In general, the choice of the final schema should take into consideration the semantics and usages of the underlying data or application and this is where user intervention is beneficial. By default, without further input from users, NeT chooses as the final schema the nested table where the most number of nestings succeeded - this is a schema which provides low “data redundancy” - as given in Table 9 of Appendix.

**Example 4.** Using NeT with the element-oriented mode,  $\mathbb{R}_1$  in Example 1 would be translated to  $\mathbb{X}_4 = (E, A, M, P, r, \Sigma)$ , where

$$\begin{aligned} E &= \{student, professor\} \\ A(professor) &= \{Pname\} \\ M(student) &= (Sname, Advisor^+, Course^+) \\ M(professor) &= (Age^?) \\ P(Pname) &= (ID, \neg?, \epsilon, \epsilon) \\ r &= \{student, professor\} \\ \Sigma &= \{\{Sname, Advisor, Course\} \xrightarrow{key} student, \\ &Pname \xrightarrow{key} professor, Advisor \subseteq Pname\} \end{aligned}$$

□

We expect that the NeT algorithm will be especially useful in two scenarios, outlined below.

- The given relation is in 3NF (or BCNF) but not in 4NF. Non-fully normalized relations occur quite commonly in legacy databases, and they exhibit data redundancy. The NeT algorithm helps to decrease the data redundancy in such cases.

As an example, consider the relation  $ctx(\text{Course}, \text{Teacher}, \text{Text})$ , which gives the set of teachers and the set of text books for each course. Assume that the following multivalued dependencies hold,  $\text{Course} \twoheadrightarrow \text{Teacher}$ , and  $\text{Course} \twoheadrightarrow \text{Text}$ . Suppose the relation  $ctx$  is represented as such (i.e.,  $ctx$  is not in 4NF). The key for this relation is given by  $\{\text{Course}, \text{Teacher}, \text{Text}\} \xrightarrow{key} ctx$ . When we do nesting on  $ctx$ , we will get the following table  $ctx'(\text{Course}, \text{Teacher}^+, \text{Text}^+)$ . Thus NeT helps in removing data redundancies arising from multivalued dependencies.

- It is sometimes possible to represent the given relation “more intuitively” as a nested table by performing grouping on one or more of the attributes. As an example, consider the relation  $emp(empNum, branch)$  where the key is given by  $empNum \xrightarrow{key} emp$ . This relation gives the employees and the branch where they work. When NeT is applied on the above relation, we might get the new nested relation as  $emp'(empNum^+, branch)$ . This relation has grouped the list of employees by their branch.

Thus we observe that NeT is useful for decreasing data redundancy and obtaining a “more intuitive” schema by (1) removing redundancies caused by multivalued dependencies and (2) performing grouping on attributes. However NeT considers tables one by one, and *cannot* obtain a *big picture* of the relational schema where many tables are interconnected with each other through various other dependencies such as inclusion dependencies. To remedy this problem, we propose to use other semantic constraints of relational schema.

## 4 Translation using Inclusion Dependencies

In this section, we consider one kind of semantic constraints called *Inclusion Dependency (IND)* in database theory. Other kinds of semantic constraints such as *Functional Dependency (FD)* or *Multi-Valued Dependency (MVD)* are not considered in this section for the following reasons: In general, most CASE tools for relational database design generate schemas at least in 3NF. If there existed any MVDs, for instance, then the table would have been split up (i.e., normalized to 4NF) to avoid excessive data redundancy. Even if it was not normalized to 4NF, if we apply the NeT algorithm, most of the MVDs would have been removed.

General forms of INDs are difficult to acquire from the database automatically. However, we shall consider the most pervasive form of INDs - foreign key constraints - which can be queried through ODBC/JDBC interface. We study the translation of inclusion dependencies incrementally in three steps. In the first step, we consider the simplest case - one foreign key constraint defined between two tables. In the second step, we consider the case when there exist two foreign key constraints among three tables. In the third step, we consider the general case of mapping any given relational schema.

### 4.1 One Foreign Key between two Tables

Foreign key constraints are a special kind of INDs where the attributes being referenced form the *primary key* of the referenced relation. For two distinct tables  $s$  and  $t$  with lists of columns  $X$  and  $Y$ , respectively, suppose we have a foreign key constraint  $s[\alpha] \subseteq t[\beta]$ , where  $\alpha \subseteq X$  and  $\beta \subseteq Y$ . Also suppose that  $K_s \subseteq X$  is the key for  $s$ . Then, rewriting this in  $\mathbb{R}$  notation, we have:  $T = \{s, t\}, C(s) = \{X\}, C(t) = \{Y\}, \Delta = \{s[\alpha] \subseteq t[\beta], \beta \xrightarrow{key} t, K_s \xrightarrow{key} s\}$ .

Different cardinality binary relationships between  $s$  and  $t$  can be expressed in the relational model by a combination of the following: (1)  $\alpha$  is unique/not-unique (2)  $\alpha$  is nullable/non-nullable.

The translation of two tables  $s, t$  with a foreign key constraint into *XSchema*, summarized in Table 3, works as follows:

- If  $\alpha$  is non-nullable (i.e., none of the columns of  $\alpha$  can take null values), then:

- If  $\alpha$  is unique, then there is a 1 : 1 relationship between  $s$  and  $t$ . This can be captured as a sub-element  $M(t) = (Y, s?)$ .
- If  $\alpha$  is not-unique, then there is a 1 :  $n$  relationship between  $s$  and  $t$ , and this is captured as a sub-element  $M(t) = (Y, s^*)$ .

If  $s$  is represented as a sub-element of  $t$ , then the key for  $s$  will change from  $K_s$  to  $(K_s - \alpha)$ . The key for  $t$  will remain the same.

- If  $\alpha$  is nullable, then the IND is represented as such in *XSchema*. Here we do flat translation on  $s$ , and copy the IND  $s[\alpha] \subseteq t[\beta]$  to  $\Sigma$ .

Let us study the case when  $\alpha$  is nullable more closely with the following example. Consider the relation  $t(w_1, w_2, w_3)$  with key  $(w_1, w_2)$ . Let  $t$  have the following tuples:  $\{(1, 1, 1)\}$ . Now consider  $s(v_1, v_2, v_3)$  with key  $(v_2, v_3)$ , and IND  $s[v_1, v_2] \subseteq t[w_1, w_2]$ . Let  $s$  have the following tuples:  $\{(null, 1, 1), (null, 1, 2), (null, 2, 1), (1, 1, 3)\}$ . We can observe that we *cannot* represent  $s$  as  $s(v_3)$ , and obtain the values of  $(v_1, v_2)$  for an  $s$  tuple by representing this  $s$  tuple as a child of a  $t$  tuple, or by having an IDREF attribute for the  $s$  tuple that refers to a  $t$  tuple. This is because  $v_1$  is nullable.

In such a case, we represent the IND as such in *XSchema*. In this paper, we are concerned mostly with the usage of sub-elements and IDREF attribute for translation, and therefore, we will focus on the case when  $\alpha$  is non-nullable, unless stated otherwise.

$\alpha$	$s : t$	<i>XSchema</i>
$v, \neg?$	$(1, 1) : (0, 1)$	$M(t) = (Y, s?),$ $M(s) = (X - \alpha),$ $\Sigma = \{(K_s - \alpha) \xrightarrow{key} s, \beta \xrightarrow{key} t\}$
$v, ?$	$(0, 1) : (0, 1)$	$M(t) = (Y), M(s) = (X),$ $\Sigma = \{s[\alpha] \subseteq t[\beta],$ $K_s \xrightarrow{key} s, \beta \xrightarrow{key} t\}$
$\neg v, \neg?$	$(1, 1) : (0, n)$	$M(t) = (Y, s^*),$ $M(s) = (X - \alpha),$ $\Sigma = \{(K_s - \alpha) \xrightarrow{key} s, \beta \xrightarrow{key} t\}$
$\neg v, ?$	$(0, 1) : (0, n)$	$M(t) = (Y), M(s) = (X),$ $\Sigma = \{s[\alpha] \subseteq t[\beta],$ $K_s \xrightarrow{key} s, \beta \xrightarrow{key} t\}$

Table 3: Different values taken by  $\alpha$ , the corresponding cardinality of the binary relationship between  $s$  and  $t$ , and the corresponding translation to *XSchema*.  $v$  and  $?$  denote “unique” and “nullable”, respectively.

**Example 5.** Consider two tables `student` and `professor` of Example 1 again. There is a foreign key  $Advisor \subseteq Pname$  and  $Advisor$  is not unique. Using the above rules, the schema will be mapped to the following XML schema in DTD notation:

```
<!ELEMENT professor (Pname, Age, student*)>
<!ELEMENT student (Sname, Course)>
```

Note the usage of  $*$  attached to the sub-element `student`. Note further that to identify a unique `student` element for a given professor, one needs now only `Sname` and `Course` pair (`Advisor` attribute was removed from the original key attribute list).  $\square$

## 4.2 Two Foreign Key among three Tables

Now consider the case where two foreign key constraints exist among three tables  $s, t_1, t_2$  with a list of columns  $X, Y_1, Y_2$ , respectively, such that  $s[\alpha] \subseteq t_1[\beta_1]$  and  $s[\gamma] \subseteq t_2[\beta_2]$ , where  $\alpha, \gamma \subseteq X$  and are non-nullable,  $\beta_1 \subseteq Y_1$  and  $\beta_2 \subseteq Y_2$ . If one applies the mapping rules for the case of a foreign key between two tables in Section 4.2 one at a time, then one will have a combination of the following depending on whether  $\alpha$  and  $\gamma$  are unique or not: (1)  $M(t_1) = (Y_1, s?)$  or  $M(t_1) = (Y_1, s^*)$ , (2)  $M(t_2) = (Y_2, s?)$  or  $M(t_2) = (Y_2, s^*)$ .

The above translation has redundancy, and it exhibits the phenomenon known in database theory as “update anomaly” for  $s$ . That is, when one wants to update data for  $s$ , he/she needs to update  $s$  in two different places – fragment of  $s$  data under both  $t_1$  and  $t_2$ . On the contrary, the original relational schema is “better” because one needs to update tuples of  $s$  in a single place. The same problem occurs for the case of “delete” as well. To avoid these anomalies, one of the two foreign key constraints should be captured either using INDs or using IDREF attributes. For example, let us assume that the first foreign key constraint  $s[\alpha] \subseteq t_1[\beta_1]$  is represented as  $M(t_1) = (Y_1, s^*)$ ,  $M(s) = (X - \alpha)$ . Then the second foreign key constraint  $s[\gamma] \subseteq t_2[\beta_2]$  can be represented using IDREF attribute as follows:

$$A(t_2) = \{ID\_t_2\}, \quad P(ID\_t_2) = (ID, \neg?, \epsilon, \epsilon)$$

$$A(s) = \{Ref\_t_2\}, \quad P(Ref\_t_2) = (IDREF, \neg?, \epsilon, \epsilon)$$

$$M(t_2) = (Y_2), \quad M(s) = (X - \alpha - \gamma)$$

Let us denote the old and new keys for  $s$  as  $K_s$  and  $K'_s$ , respectively. Then,  $K'_s$  is determined as follows: (1) if  $\alpha \cap K_s = \phi$ , then  $K'_s = K_s$ , and (2) if  $\alpha \cap K_s \neq \phi$ , then  $K'_s = (K_s - \alpha) \cup Ref\_t_2$

**Example 6.** In addition to two tables `student` and `professor` of Example 5, consider a third table `class(Cname, Room)` with a second foreign key  $student[Course] \subseteq class[Cname]$ . Then, using the above rules, the schema will be mapped to the following XML schema in DTD notation:

```
<!ELEMENT professor (Pname, Age, student*)>
<!ELEMENT student (Sname)>
<!ELEMENT class (Cname, Room)>
<!ATTLIST student Ref_class IDREF>
<!ATTLIST class ID_class ID>
```

Note the addition of two new attributes - `Ref_class` of type `IDREF` and `ID_class` of type `ID`. The new key for `student` is given by  $\{Sname, Ref\_class \xrightarrow{key} student\}$ , which cannot be represented in DTD.  $\square$

Note that between two foreign keys, deciding which one is represented as sub-element and which one is represented as IDREF attribute can best be done based on further semantics.

### 4.3 A General Relational Schema

Now let us consider the most general case with set of tables  $\{t_1, \dots, t_n\}$  and INDs  $t_i[\alpha_i] \subseteq t_j[\beta_j]$ , where  $i, j \leq n$ . We consider only those INDs that are foreign key constraints (i.e.,  $\beta_j \xrightarrow{key} t_j$ ), and where  $\alpha_i$  is non-nullable. The relationships among tables can be captured by a graph representation, termed as IND-Graph.

**Definition 4 (IND-Graph)** An IND-Graph  $G = (V, E)$  consists of a node set  $V$  and a directed edge set  $E$ , such that for each table  $t_i$ , there exists a node in  $V$ , and for each distinct IND  $t_i \subseteq t_j$ ,  $t_j \rightarrow t_i$  exists in  $G$ .  $\square$

Note the edge direction is reversed from the IND direction for convenience. Given a set of INDs, such IND-Graph can be easily constructed. Once IND-Graph is constructed, one needs to decide the starting point to apply translation rules. For that purpose, we use the notion of **top nodes** similar to the one in [20, 13], where an element is a top node if it *cannot* be represented as a sub-element of any other element. Such top nodes can be identified as follows:

1. An element  $s$  is a top node, if there exists no other element  $t$ ,  $t \neq s$ , where there is a IND of the form  $s[\alpha] \subseteq t[\beta]$ , and  $\alpha$  is non-nullable.
2. Consider a set of elements  $S = s_1, s_2, \dots, s_k$  that form a cyclic set of INDs and none of the elements in  $S$  is a top node by 1. Suppose there exists no element  $t \notin S$ , such that there is a IND of the form  $s_j[\alpha] \subseteq t[\beta]$ , and  $\alpha$  is non-nullable. In this case, choose any one of the elements in  $S$  as a top node.

Let  $T$  denote the set of top nodes. After identifying the top nodes, we traverse  $G$ , using say Breadth-First Search (BFS), until we traverse all the nodes and edges, and represent the INDs as sub-elements or IDREF attributes. The algorithm for **Constraint-based Translation** (CoT) is given in Table 4.

**Example 7.** Consider a schema and its associated INDs in Table 5. The IND-Graph is shown in Figure 2. Two top nodes are identified (1) `course`: There is no node  $t$ , where there is an IND of the form `course` $[\alpha]$

<code>student(Sid, Name, Advisor)</code> <code>emp(Eid, Name, ProjName)</code> <code>prof(Eid, Name, Teach)</code> <code>course(Cid, Title, Room)</code> <code>dept(Dno, Mgr)</code> <code>proj(Pname, Pmgr)</code>
<code>student(Advisor) <math>\subseteq</math> prof(Eid)</code> <code>emp(ProjName) <math>\subseteq</math> proj(Pname)</code> <code>prof(Teach) <math>\subseteq</math> course(Cid)</code> <code>prof(Eid, Name) <math>\subseteq</math> emp(Eid, Name)</code> <code>dept(Mgr) <math>\subseteq</math> emp(Eid)</code> <code>proj(Pmgr) <math>\subseteq</math> emp(Eid)</code>

Table 5: An example schema with associated INDs.

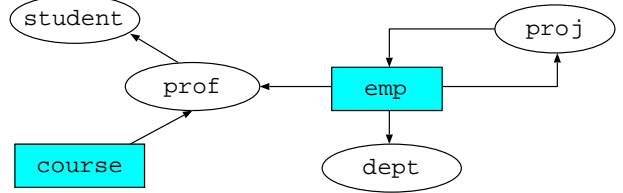


Figure 2: The IND-Graph representation of the schema of Table 5 (*top nodes* denoted by rectangular nodes).

$\subseteq t[\beta]$ , and (2) `emp`: There is a cyclic set of INDs between `emp` and `proj`, and there exists no node  $t$  such that there is an IND of the form `emp` $[\alpha] \subseteq t[\beta]$  or `proj` $[\alpha] \subseteq t[\beta]$ . Therefore of `emp` and `proj` we decided to choose `emp` arbitrarily. Table 6 shows one of the possible orders in which the different INDs are visited, the choice made to represent the IND (either sub-element or IDREF attribute), and the resulting changes in `XSchema`.  $\square$

It is worthwhile to point out that there are several places in CoT where human experts can determine better mapping based on the semantics and usages of the underlying data or application.

- The CoT algorithm identifies a minimal set of top-nodes, breaking any ties (when there are cyclic INDs) arbitrarily. A better mapping might have more top-nodes than this minimal set, or might choose to break a tie in a particular manner.
- Given a set of foreign-key constraints on one table, CoT chooses one foreign-key constraint to be represented as a sub-element, and represents the remaining using IDREF attributes. Human experts might be able to provide better input as to which constraint should be represented as sub-element, and which as IDREF attributes.

Examples so far have shown the conversion flow of  $\mathbb{X} \rightarrow \text{CoT} \rightarrow \text{DTD}$ . We can also have the conversion flow  $\mathbb{X} \rightarrow \text{NeT} \rightarrow \text{CoT} \rightarrow \text{DTD}$  as shown in Example 10 in Appendix. However this imposes a restriction; when



CoT:  $R = (T, C, P, \Delta) \implies X = (E, A, M, P, r, \Sigma)$

1. Construct IND-Graph  $G = (V, E)$  from the given INDs; Identify  $T$ , the set of top nodes. Define  $S = T$  to keep track of top-nodes and nodes that are represented as sub-elements.
2. For each top-node  $t \in T$ , do BFS. Suppose we reach a node  $w$  from  $v$  (i.e., IND:  $w[\alpha] \subseteq v[\beta]$ ); Let  $C(w) = C_w$ , and  $C(v) = C_v$ .
  - (a) If  $w \notin S$  (i.e.,  $w$  is *not* yet a sub-element of some other node), translate the IND as in Section 4.1.
    - i. If  $\alpha$  is unique, then  $M(v) = (C_v, w?)$ .
    - ii. If  $\alpha$  is not-unique, then  $M(v) = (C_v, w^*)$ .
    - iii.  $M(w) = (C_w - \alpha)$ .
    - iv.  $S = S \cup w$ .
  - (b) If  $w \in S$  (i.e.,  $w$  is already a sub-element of some other node), translate the IND as IDREF attribute as in Section 4.2.
    - i.  $A(v) = \{ID_v\}$ ,  $A(w) = \{Ref_v\}$ ,  $M(v) = (C_v)$ ,  $M(w) = (C_w - \alpha)$ ,  $\Sigma = K_w' \xrightarrow{key} w$ .
3. Copy the remaining integrity constraints in  $\Delta$  to  $\Sigma$ . Also set  $r = T$ .

Table 4: CoT algorithm.

No.	IND	sub-element vs. IDREF	Representation in $XSchema$
1	$\text{prof}(\text{Teach}) \subseteq \text{course}(\text{Cid})$	sub-element	$M(\text{course}) = (\text{Cid}, \text{Title}, \text{Room}, \text{prof}^*)$ , $M(\text{prof}) = (\text{Eid}, \text{Name})$
2	$\text{student}(\text{Advisor}) \subseteq \text{prof}(\text{Eid})$	sub-element	$M(\text{prof}) = (\text{Eid}, \text{Name}, \text{student}^*)$ , $M(\text{student}) = (\text{Sid}, \text{Name})$
3	$\text{dept}(\text{Mgr}) \subseteq \text{emp}(\text{Eid})$	sub-element	$M(\text{emp}) = (\text{Eid}, \text{Name}, \text{ProjName}, \text{dept}^*)$ , $M(\text{dept}) = (\text{Dno})$
4	$\text{proj}(\text{Pmgr}) \subseteq \text{emp}(\text{Eid})$	sub-element	$M(\text{emp}) = (\text{Eid}, \text{Name}, \text{ProjName}, \text{dept}^*, \text{proj}^*)$ $M(\text{proj}) = (\text{Pname})$
5	$\text{emp}(\text{ProjName}) \subseteq \text{proj}(\text{Pname})$	IDREF	$M(\text{emp}) = (\text{Eid}, \text{Name}, \text{dept}^*, \text{proj}^*)$ $A(\text{proj}) = \{\text{ID\_proj}\}$ , $A(\text{emp}) = \{\text{Ref\_proj}\}$
6	$\text{prof}(\text{Eid}, \text{Name}) \subseteq \text{emp}(\text{Eid}, \text{Name})$	IDREF	$M(\text{prof}) = (\text{student}^*)$ , $A(\text{emp}) = \{\text{ID\_emp}\}$ , $A(\text{prof}) = \{\text{Ref\_emp}\}$ $\Sigma = \{\text{Ref\_emp} \xrightarrow{key} \text{prof}\}$

Table 6: The order where INDs are chosen by the CoT algorithm, and how the translation is done.

NeT followed by CoT are applied, nesting can be done only on attributes that do not participate in any IND.

## 5 Discussion

All three algorithms – FT, NeT, and CoT – are “correct” in the sense that they all have preserved the original information of relational schema. For instance, using the notion of information capacity [16], a theoretical analysis for the correctness of our translation procedures is possible; we can actually show that NeT and CoT algorithms are *equivalence preserving transformations*. However, we defer this detailed analysis to a later version.

With respect to the “goodness” of XML schema that the proposed algorithms generate, it is not obvious to bluntly state whether or not they are good, since there has not been any unanimous normalization theory for XML model yet. Some early work for nested relational model (e.g., [19]) is related, but more recently a few proposals have been made for normal forms of XML model (e.g., [6, 24]). To a greater or lesser ex-

tent, the crux of such normal forms is an attempt to reduce data redundancy so that various anomalies can be avoided. Although the output schema from NeT or CoT does not exactly fit into normal forms defined by [6, 24], they share similar properties. For instance, identifying multivalued attributes and making them repeating sub-elements in NeT is essentially a necessary step towards “object class normal form” in [24]. The use of reference attributes in CoT for handling multiple foreign key constraints defined on one table (Section 4.2) can be explained similarly. Therefore, we would like to point out that although it is early to formally prove the goodness of our proposals, it is evident that our proposals lead to less *redundant* yet *correct* XML schema.

Test Set	# of attributes / # of tuples	# of successful nesting / # of attempted nesting	# of values in nested table / # original data values	Nested Table Size / Original Size (KB)	# of nested attributes	time (sec.)
Balloons1	5 / 16	42 / 64	22 / 80	0.152 / 0.455	3	1.08
Balloons2	5 / 16	42 / 64	22 / 80	0.150 / 0.455	3	1.07
Balloons3	5 / 16	40 / 64	42 / 80	0.260 / 0.455	3	1.14
Balloons4	5 / 16	42 / 64	22 / 80	0.149 / 0.455	3	1.07
Hayes	6 / 132	1 / 6	522 / 792	1.219 / 1.758	1	1.01
Bupa	7 / 345	0 / 7	2387 / 2387	7.234 / 7.234	0	4.40
Balance	5 / 625	56 / 65	1120 / 3125	2.259 / 6.265	4	21.48
TA_Eval	6 / 110	253 / 326	534 / 660	1.281 / 1.559	5	24.83
Car	7 / 1728	1870 / 1957	779 / 12096	3.157 / 51.867	6	469.47
Flare	13 / 365	11651 / 13345	2834 / 4745	5.715 / 9.533	4	6693.41

Table 7: Summary of NeT experimentations.

## 6 Experimental Results

### 6.1 NeT Results

We implemented the NeT algorithm described in Table 9<sup>5</sup>. We used two additional optimization rules in our implementation: (1) if  $nest_X(t) = t$ , then  $nest_X(nest_{All_L}(t)) = nest_{All_L}(t)$  for any list of columns,  $L$ , and (2) if  $nest_X(nest_{All_L}(t)) = nest_{All_L}(t)$  for any column  $X$  and all possible list of columns  $L$  of length  $l$ , then  $nest_X(nest_{All_M}(t)) = nest_{All_M}(t)$  for any column  $X$  and all possible list of columns  $M$  of length  $m$ , where  $m \geq l$ .

Our preliminary results comparing the goodness of the  $XSchema$  obtained from NeT, and FT with that obtained from DB2XML v 1.3 [23] is given in Appendix A.3. We further applied our NeT algorithm on several test sets drawn from UCI KDD<sup>6</sup> / ML<sup>7</sup> repositories, which contain a multitude of single-table relational schemas and data. Sample results are shown in Table 7. Two metrics are used as follows:

$$\text{NestRatio} = \frac{\# \text{ of successful nesting}}{\# \text{ of total nesting}}$$

$$\text{ValueRatio} = \frac{\# \text{ of data values } D \text{ in the nested table}}{\# \text{ of original data values}}$$

where  $D$  is the number of individual data values present in the table. For example,  $D$  for the row  $(\{1, 2, 3\}, a, 10)$  of a nested table is 5.

Note that  $NestRatio$  denotes the efficiency of our optimization rules, while  $ValueRatio$  implies whether it was useful to perform nesting.

In our experimentation, we observed that most of the attempted nestings are successful, and hence our optimization rules are quite efficient. In Table 7, we see that nesting was useful for all the data sets except for the Bupa data set. Also nesting was *especially* useful for the Car data set, where the size of the nested table

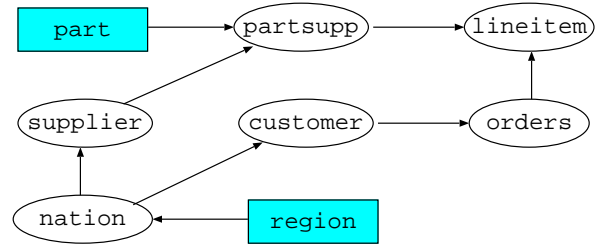


Figure 3: The IND-Graph representation of the TPC-H schema.

is only 6% of the original data set. Time required for nesting is an important parameter, and it depends on the number of attempted nestings, and the number of tuples. The number of attempted nestings depend on the number of attributes, and increase drastically as the number of attributes increases. This is observed for the Flare data set, where we have to do nesting on 13 attributes. The NeT algorithm could nest only up to 4 attributes for feasibility reasons, it is actually possible to nest more for this data set.

### 6.2 CoT Results

For testing CoT, we need some well-designed relational schema where tables are interconnected via inclusion dependencies. For this purpose, we use the TPC-H schema v 1.3.0<sup>8</sup>, which is an ad-hoc, decision support benchmark and has 8 tables and 8 inclusion dependencies. The IND-Graph for the TPC-H schema is shown in Figure 3.

CoT identifies two top-nodes - part and region. Suppose we start the scan of the top-nodes from region, we get the  $XSchema$  shown in Appendix A.4. Six of the eight inclusion dependencies are mapped using sub-element, and the remaining two are mapped using IDREF attributes. We believe that the  $XSchema$  produced by CoT is “more intuitive” than the relational schema we started with.

<sup>5</sup> Available at <http://www.cs.ucla.edu/~mani/xml>

<sup>6</sup> <http://kdd.ics.uci.edu/>

<sup>7</sup> <http://www.ics.uci.edu/~mlern/MLRepository.html>

<sup>8</sup> <http://www.tpc.org/tpch/spec/h130.pdf>

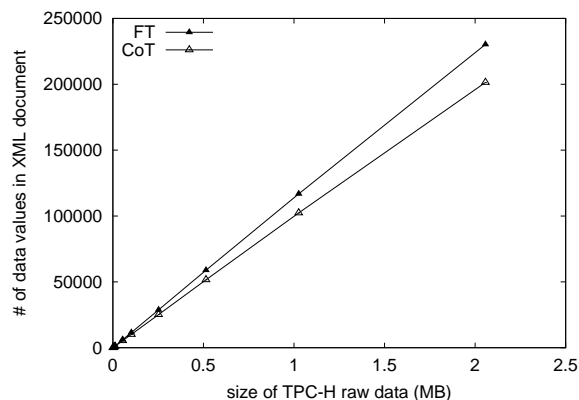


Figure 4: Comparison of XML documents generated by FT and CoT algorithms for TPC-H data.

Figure 4 plots the number of data values in the XML document generated by FT and CoT depending on the size of the original data. Because FT is a flat translation, the number of data values in the XML document generated by FT is the same as the number of data values in the original data. However, CoT is able to decrease the number of data values in the generated XML document by more than 12%.

## 7 Conclusion

We have presented two relational-to-XML conversion algorithms - NeT and CoT. The naive translation algorithm FT translates the “flat” relational model to “flat” XML model in a one-to-one manner. Thus FT does *not* use the non-flat features of the XML model, possible through regular expression operators such as “\*” and “+”. To remedy this problem, we first presented NeT, which uses the *nest* operator to generate a more precise and intuitive XML Schema from relational inputs. When poorly designed or legacy relational schema needs to be converted to XML format, NeT can suggest a fairly intuitive XML schema. However NeT is only applicable to a single table at a time, and *cannot* obtain a big picture of a relational schema where many tables are interconnected with each other. Our next algorithm CoT addresses this problem - CoT uses *semantic constraints* (especially inclusion dependencies) specified in the relational model to come up with a more intuitive XML Schema for the entire relational schema.

Thus our approaches have the following properties: (1) automatically infer a “good” XML Schema from a given relational schema, (2) remove redundancies that might be present in poorly designed or legacy relational schema (3) maintain semantic constraints during translation. With a rapid adoption of XML standards among industries and majority of data still stored in relational databases, the need to correctly and effectively convert relational data into XML for-

mat is imminent. We believe our proposed methods are good additions to such a practical problem.

Our work in this paper concentrates on obtaining a “good” and “correct” XML schema. However there are still several other issues to be studied. Implementation issues (e.g., I/O cost, tagging strategy, nesting strategy) are very important. Early investigation on these issues is done in [2]. Since our work in this paper proposes algorithms which can result in a fairly complex target XML schema as an output, studying an efficient implementation of our NeT and CoT algorithms is an important direction. Another direction of future research is studying the normalization theory of XML schema. By formally defining what is a “good” XML schema, one can devise better relational-to-XML conversion algorithms that result in normalized XML schema. Also our NeT algorithm performed only single attribute nesting. Multiple attribute nesting is another interesting research direction.

## References

- [1] P. V. Biron and A. Malhotra (Eds). “XML Schema Part 2: Datatypes”. W3C Recommendation, May 2001. <http://www.w3.org/TR/xmlschema-2/>.
- [2] R. Bourret. “Data Transfer Strategies: Transferring Data between XML Documents and Relational Databases”. Web page, 2000. <http://www.rpbouret.com/xml/DataTransfer.htm>.
- [3] T. Bray, J. Paoli, and C. M. Sperberg-McQueen (Eds). “Extensible Markup Language (XML) 1.0 (2nd Edition)”. W3C Recommendation, Oct. 2000. <http://www.w3.org/TR/2000/REC-xml-20001006>.
- [4] M. Carey, D. Florescu, Z. Ives, Y. Lu, J. Shanmugasundaram, E. Shekita, and S. Subramanian. “XPERANTO: Publishing Object-Relational Data as XML”. In *Int’l Workshop on the Web and Databases (WebDB)*, Dallas, TX, May 2000.
- [5] J. Clark and M. Murata (Eds). “RELAX NG Specification”. OASIS Committee Specification, Dec. 2001. <http://www.oasis-open.org/committees/relax-ng/spec-20011203.html>.
- [6] D. W. Embley and W. Y. Mok. “Developing XML Documents with Guaranteed “Good” Properties”. In *Int’l Conf. on Conceptual Modeling (ER)*, Yokohama, Japan, Nov. 2001.
- [7] W. Fan and J. Siméon. “Integrity Constraints for XML”. In *ACM PODS*, Dallas, TX, May 2000.

- [8] M. F. Fernandez, W.-C. Tan, and D. Suciu. “SilkRoute: Trading between Relations and XML”. In *Int’l World Wide Web Conf. (WWW)*, Amsterdam, Netherlands, May 2000.
- [9] P. C. Fischer, L. V. Saxton, S. J. Thomas, and D. V. Gucht. “Interactions between Dependencies and Nested Relational Structures”. *J. Computer and System Sciences (JCSS)*, 31(3):343–354, Dec. 1985.
- [10] M. N. Garofalakis, A. Gionis, R. Rastogi, S. Shadri, and K. Shim. “XTRACT: A System for Extracting Document Type Descriptors from XML Documents”. In *ACM SIGMOD*, Dallas, TX, May 2000.
- [11] R. L. Graham, D. E. Knuth, and O. Patashnik. “Concrete Mathematics: A Foundation for Computer Science”. Addison-Wesley Pub., 1994.
- [12] G. Jaeschke and H.-J. Schek. “Remarks on the Algebra of Non First Normal Form Relations”. In *ACM PODS*, Los Angeles, CA, Mar. 1982.
- [13] D. Lee and W. W. Chu. “Constraints-preserving Transformation from XML Document Type Definition to Relational Schema”. In *Int’l Conf. on Conceptual Modeling (ER)*, pages 323–338, Salt Lake City, UT, Oct. 2000.
- [14] D. Lee, M. Mani, F. Chiu, and W. W. Chu. “Nesting-based Relational-to-XML Schema Translation”. In *Int’l Workshop on the Web and Databases (WebDB)*, Santa Barbara, CA, May 2001.
- [15] M. Mani, D. Lee, and R. R. Muntz. “Semantic Data Modeling using XML Schemas”. In *Int’l Conf. on Conceptual Modeling (ER)*, Yokohama, Japan, Nov. 2001.
- [16] R. J. Miller, Y. E. Ioannidis, and R. Ramakrishnan. “Schema Equivalence in Heterogeneous Systems: Bridging Theory and Practice (Extended Abstract)”. In *EDBT*, Cambridge, UK, Mar. 1994.
- [17] M. Murata, D. Lee, and M. Mani. “Taxonomy of XML Schema Languages using Formal Language Theory”. In *Extreme Markup Languages*, Montreal, Canada, Aug. 2001.
- [18] C. Nentwich, W. Emmerich, A. Finkelshtein, and A. Zisman. “BOX: Browsing Objects in XML”. *Software Practice and Experience*, 30(15):1661–1676, 2000. <http://www.cs.ucl.ac.uk/staff/c.nentwich/Box/>.
- [19] Z. M. Özsoyoglu and L. Y. Yuan. “A New Normal Form for Nested Relations”. *ACM Trans. on Database Systems (TODS)*, 12(1):111–136, Mar. 1987.
- [20] J. Shanmugasundaram, K. Tufte, G. He, C. Zhang, D. DeWitt, and J. Naughton. “Relational Databases for Querying XML Documents: Limitations and Opportunities”. In *VLDB*, Edinburgh, Scotland, Sep. 1999.
- [21] T. Shimura, M. Yoshikawa, and S. Uemura. “Storage and Retrieval of XML Documents using Object-Relational Databases”. In *Int’l Conf. on Database and Expert Systems Applications (DEXA)*, pages 206–217, Florence, Italy, Aug. 1999.
- [22] H. S. Thompson, D. Beech, M. Maloney, and N. Mendelsohn (Eds). “XML Schema Part 1: Structures”. W3C Recommendation, May 2001. <http://www.w3.org/TR/xmlschema-1/>.
- [23] V. Turau. “Making Legacy Data Accessible for XML Applications”. Web page, 1999. <http://www.informatik.fh-wiesbaden.de/~turau/veroeff.html>.
- [24] X. Wu, T. W. Ling, M. L. Lee, and G. Dobbie. “Designing Semistructured Database Using ORASS Model”. Unpublished Manuscript, 2001.

## A Appendix

### A.1 Proofs

PROOF. (**Remark 1**) The number of the *first* nesting along  $n$  columns is the same as the number of 1-element sequences:  $n$ . The number of the *second* nesting along  $n$  columns is again the same as the number of 2-elements sequences by P2:  $n(n-1)$ . Continuing this, the number of the *last* nesting along  $n$  columns is again the same as the number of  $n$ -elements sequences:  $n + n(n-1) + \dots + n(n-1) \dots (2)(1) = n^1 + n^2 + \dots + n^n = \sum_{k=1}^n n^k$  (q.e.d)

PROOF. (**Lemma 1**) Consider a table  $t$  with column set  $C$ , and candidate keys,  $K_1, K_2, \dots, K_n \subseteq C$ . Consider a column  $X \in C$ , such that  $X$  is not an attribute of at least one of the candidate keys, say  $X \notin K_i$ . Now  $\overline{X} \supseteq K_i$ , and hence  $\overline{X}$  is unique. Thus, no two tuples can agree on  $\overline{X}$ . Therefore, by the definition of the *nest* operator, nesting on  $X$  will fail. (q.e.d)

PROOF. (**Lemma 4**) Let  $C$  denote the set of columns of  $t$ , let the candidate keys be  $K_1, K_2, \dots, K_n \subseteq C$ . Let  $K = K_1 \cap K_2 \cap \dots \cap K_n$ , where  $|K| = m$ . The first column to be nested, say  $X$  is chosen such that  $X \in K$  by Lemma 1, in one of the  $m$  ways. Now after the first nesting, by Lemma 2, we have a new candidate key  $\overline{X}$ . The next column to be nested is chosen from  $K \cap \overline{X}$ , where  $|K \cap \overline{X}| = m - 1$ . Thus we have  $m - 1$  ways of choosing the second column

for nesting. Continuing this, we have total number of nesting is  $m + m(m-1) + \dots + m(m-1) \dots (2)(1) = \sum_{k=1}^m m^k$ . (q.e.d)

## A.2 More Conversion Examples

**Example 8.**  $\mathbb{R}_1$  in Example 1 would be translated to  $\mathbb{X}_8 = (E, A, M, P, r, \Sigma)$  via FT, where

$$\begin{aligned}
E &= \{student, professor\} \\
A(student) &= \{Sname, Advisor, Course\} \\
A(professor) &= \{Pname, Age\} \\
M(student) &= \epsilon \\
M(professor) &= \epsilon \\
P(Sname) &= (string, ?, \epsilon, \epsilon) \\
P(Advisor) &= (IDREF, ?, \epsilon, "J.Smith") \\
P(Course) &= (string, ?, \epsilon, \epsilon) \\
P(Pname) &= (ID, \neg?, \epsilon, \epsilon) \\
P(Age) &= (integer, ?, \epsilon, \epsilon) \\
r &= \{student, professor\} \\
\Sigma &= \{\{Sname, Advisor, Course\} \xrightarrow{key} student, \\
&\quad Pname \xrightarrow{key} professor, Advisor \subseteq Pname\} \quad \square
\end{aligned}$$

**Example 9.**  $\mathbb{X}_7$  of Example 7 can be further rewritten in element-oriented mode to DTD notations as follows:

```

<!ELEMENT course (Cid, Title, Room, prof*)>
<!ELEMENT prof (Name, student*)>
<!ATTLIST prof Eid IDREF>
<!ELEMENT student (Sid, Name)>
<!ELEMENT emp (Name, ProjName, dept*,
proj*)>
<!ATTLIST emp Eid ID, ProjName IDREF>
<!ELEMENT dept (Dno)>
<!ELEMENT proj EMPTY>
<!ATTLIST proj Pname ID>

```

**Example 10.** Consider  $\mathbb{R}_1$  in Example 1. Let us first apply NeT and then CoT on this. When we apply NeT, we perform nesting only on the column *Course* of *student*, and obtain a content model for *student* as  $M(student) = (Sname, Advisor, Course^+)$ . Now applying CoT on the above schema, we get the output  $\mathbb{X}_{10}$  as  $(E, A, M, P, r, \Sigma)$ , where

$$\begin{aligned}
E &= \{student, professor\} \\
A(professor) &= \{Pname\} \\
M(student) &= (Sname, Course^+) \\
M(professor) &= (Age^?, student^*) \\
P(Pname) &= (ID, \neg?, \epsilon, \epsilon) \\
r &= \{student, professor\} \\
\Sigma &= \{\{Sname, Course\} \xrightarrow{key} student, \\
&\quad Pname \xrightarrow{key} professor\} \quad \square
\end{aligned}$$

## A.3 NeT Experimentation

We compare the results of FT and NeT with that of DB2XML v 1.3 [23]. Consider the *Orders* table (containing 830 tuples) found in MS Access NorthWind sample database.

*Orders* (CustomerID, EmployeeID, ShipVia, ShipAddress, ShipCity, ShipCountry, ShipPostalCode)

Table 8 shows the DTDs generated by DB2XML, FT in attribute-oriented mode, NeT in both element-oriented and attribute-oriented modes, respectively. In (a), DB2XML always uses element to represent columns of a table. To represent whether the column is nullable or not, DB2XML adds a special attribute ISNULL to every element: i.e., "ISNULL = true" means the column is nullable. In (b), FT in attribute-oriented mode uses #IMPLIED or REQUIRED to represent whether the column is nullable or not. Observe that both DB2XML and FT share the same problem of translating "flat" relational schema to "flat" XML schema. In (c) and (d), NeT finds two columns *EmployeeID* and *ShipVia* can be nested. Intuitively, the new schema infers that for each *CustomerID*, multiple non-zero *EmployeeID* and multiple *ShipVia* can exist. Also NeT finds that *CustomerID* is a mandatory column. To ensure this property, (c) adds no suffix such as ? or \* to *CustomerID* sub-element and (d) uses #REQUIRED construct explicitly. We observe that the DTDs found by NeT are more succinct and more intuitive than the ones found by DB2XML.

## A.4 CoT Experimentation

CoT converted the TPC-H schema into the following  $\mathbb{X}_{Schema} \mathbb{X}_{tpc} = (E, A, M, P, r, \Sigma)$ . We show the definitions only for attributes of types ID and IDREF.

$$\begin{aligned}
E &= \{part, partsupp, lineitem, orders, \\
&\quad customer, supplier, nation, region\} \\
A(part) &= \{PartKey, Name, Mfgr, \\
&\quad Brand, Type, Size, Container, \\
&\quad RetailPrice, Comment\} \\
M(part) &= \epsilon \\
P(PartKey) &= (ID, \neg?, \epsilon, \epsilon) \\
A(region) &= \{RegionKey, Name, Comment\} \\
M(region) &= \{nation^*\} \\
P(RegionKey) &= (ID, \neg?, \epsilon, \epsilon) \\
A(nation) &= \{NationKey, Name, Comment\} \\
M(nation) &= \{supplier^*, customer^*\} \\
P(NationKey) &= (ID, \neg?, \epsilon, \epsilon) \\
A(supplier) &= \{SuppKey, Name, Address, Phone, \\
&\quad AcctBal, Comment\} \\
M(supplier) &= \{partsupp^*\} \\
P(SuppKey) &= (ID, \neg?, \epsilon, \epsilon)
\end{aligned}$$

```

<!ELEMENT Orders (CustomerID,EmployeeID,ShipVia,ShipAddress,
ShipCity,ShipCountry,ShipPostalCode)>
<!ELEMENT CustomerID (#PCDATA)>
<!ATTLIST CustomerID ISNULL (true|false) #IMPLIED>
<!ELEMENT EmployeeID (#PCDATA)>
<!ATTLIST EmployeeID ISNULL (true|false) #IMPLIED>
<!ELEMENT ShipVia (#PCDATA)>
<!ATTLIST ShipVia ISNULL (true|false) #IMPLIED>
...
<!ELEMENT ShipCountry (#PCDATA)>
<!ATTLIST ShipCountry ISNULL (true|false) #IMPLIED>
<!ELEMENT ShipPostalCode (#PCDATA)>
<!ATTLIST ShipPostalCode ISNULL (true|false) #IMPLIED>

```

(a) DB2XML

```

<!ELEMENT Orders (EMPTY)>
<!ATTLIST Orders
CustomerID CDATA #IMPLIED
EmployeeID CDATA #IMPLIED
ShipVia CDATA #IMPLIED
ShipAddress CDATA #IMPLIED
ShipCity CDATA #IMPLIED
ShipCountry CDATA #IMPLIED
ShipPostalCode CDATA #IMPLIED>

```

(b) FT in attribute-oriented mode

```

<!ELEMENT Orders (CustomerID+,EmployeeID+,ShipVia*,ShipAddress?,
ShipCity?,ShipCountry?,ShipPostalCode?)>
<!ELEMENT CustomerID (#PCDATA)>
<!ELEMENT EmployeeID (#PCDATA)>
<!ELEMENT ShipVia (#PCDATA)>
<!ELEMENT ShipAddress (#PCDATA)>
<!ELEMENT ShipCity (#PCDATA)>
<!ELEMENT ShipCountry (#PCDATA)>
<!ELEMENT ShipPostalCode (#PCDATA)>

```

(c) NeT in element-oriented mode

```

<!ELEMENT Orders (EmployeeID+,ShipVia*)>
<!ATTLIST Orders
CustomerID CDATA #REQUIRED
ShipAddress CDATA #IMPLIED
ShipCity CDATA #IMPLIED
ShipCountry CDATA #IMPLIED
ShipPostalCode CDATA #IMPLIED>
<!ELEMENT EmployeeID (#PCDATA)>
<!ELEMENT ShipVia (#PCDATA)>

```

(d) NeT in attribute-oriented mode

Table 8: DTDs generated by different algorithms.

NeT:  $R = (T, C, P, \Delta) \implies X = (E, A, M, P, r, \Sigma)$

- Each table  $t_i$  in  $R$  is translated to an element  $e_i$  in  $X$ :  $E = \bigcup_{v_i} \{e_i\}$ .
- For each table  $t_i$  in  $R$ , apply the *nest* operator repeatedly until no nesting succeeds. Choose the best nested table based on the selected criteria. Denote this table as  $t'_i(c_1, \dots, c_{k-1}, c_k, \dots, c_n)$ , where nesting succeeded on the columns  $\{c_1, \dots, c_{k-1}\}$ . If  $k = 1$  (i.e., no nesting succeeded), follow the flat translation. Otherwise, do the following:
  - For each column  $c_i$  ( $1 \leq i \leq k - 1$ ), where  $P(c_i) = (\tau, u, n, d, f)$ , if  $n = ?$ , then the content model is  $M(e_i) = (\dots, c_i^*, \dots)$ , otherwise  $M(e_i) = (\dots, c_i^+, \dots)$ .
  - For each column  $c_j$  ( $k \leq j \leq n$ ), where  $P(c_j) = (\tau, u, n, d, f)$ , do flat translation
    - (element-oriented mode) if  $n = ?$ , the content model is  $M(e_i) = (\dots, c_j^?, \dots)$ , otherwise  $M(e_i) = (\dots, c_j, \dots)$ .
    - (attribute-oriented mode) if  $c_j$  is translated to  $a_j$ , then  $A(e_i) = \bigcup_{v_j} \{a_j\}$  and  $P(a_j) = (\tau, n, d, f)$ .
- All elements  $e_i$  in  $X$  become roots:  $r = \bigcup_{v_i} \{e_i\}$ .
- Copy  $\Delta$  in  $R$  into  $\Sigma$  in  $X$ .

Table 9: NeT algorithm.

$A(partsupp)$	$\{Ref\_part, AvailQty, SupplyCost, Comment\}$	$P(CustKey)$	$(ID, \neg?, \epsilon, \epsilon)$
$M(partsupp)$	$\{lineitem^*\}$	$A(orders)$	$\{OrderKey, OrderStatus, TotalPrice, OrderDate, OrderPriority, Clerk, ShipPriority, Comment\}$
$P(Ref\_part)$	$(IDREF, \neg?, \epsilon, \epsilon)$	$M(orders)$	$\epsilon$
$A(lineitem)$	$\{Ref\_orders, LineNumber, Quantity, ExtendedPrice, Discount, Tax, ReturnFlag, LineStatus, ShipDate, CommitDate, ReceiptDate, ShipInstruct\}$	$P(OrderKey)$	$(ID, \neg?, \epsilon, \epsilon)$
$M(lineitem)$	$\epsilon$	$r$	$\{part, region\}$
$P(Ref\_orders)$	$(IDREF, \neg?, \epsilon, \epsilon)$	$\Sigma$	$\{PartKey \xrightarrow{key} part, SuppKey \xrightarrow{key} supplier, Ref\_part, SuppKey\} \xrightarrow{key} partsupp,$
$A(customer)$	$\{CustKey, Name, Address, Phone, AcctBal, MktSegment, Comment\}$		$\{Ref\_orders, LineNumber\} \xrightarrow{key} lineitem,$
$M(customer)$	$\{orders^*\}$		$OrderKey \xrightarrow{key} orders, CustKey \xrightarrow{key} customer,$
			$NationKey \xrightarrow{key} nation, RegionKey \xrightarrow{key} region\}$