

WPI-CS-TR-01-04

March 2001

**Sangam<sup>1</sup> : Modeling Transformations For Integrating Now and  
Tomorrow**

by

Kajal T. Claypool and Elke A. Rundensteiner

**Computer Science  
Technical Report  
Series**



---

**WORCESTER POLYTECHNIC INSTITUTE**

---

Computer Science Department  
100 Institute Road, Worcester, Massachusetts 01609-2280

# Sangam<sup>†</sup>: Modeling Transformations For Integrating Now and Tomorrow

Kajal T. Claypool and Elke A. Rundensteiner

Department of Computer Science  
Worcester Polytechnic Institute  
Worcester, MA 01609-2280  
{kajal|rundenst}@cs.wpi.edu

## Abstract

Today many application engineers struggle to not only publish their relational, object or ascii file data on the Web but to also integrate information from diverse sources, often inventing and reinventing a suite of hard-wired integration tools. A model management system that supports the specification and manipulation of not only data models and schemata, but also mappings between the different models in a generic manner has the promise of solving these issues. However, support for modeling and managing such mappings as objects remains an unsolved challenge. In our work, we propose a powerful middle-ware tool that successfully tackles this challenge. For this, we propose a graph-theoretic framework that allows users to explicitly model mappings between different data models as well as re-structuring within one data model. Our map metamodel is based on a set of re-usable mapping constructs that can in principle be applied on any data model described in our framework. In our work, we have tested these operators for XML and relational model mappings. Using the description of maps at the model level, mappings between specific application schemas and transformations of associated application data can be automated by our framework. Our framework guarantees the correctness of the map, of the generated transformation code, of the output data model, and of the generated application schemas, based on the correctness criteria for the map metamodel. In this paper, we also introduce the model management system Sangam that we are developing to realize our proposed map modeling theory. With Sangam we show not only the feasibility of our approach but also demonstrate the re-usability and the ease of end-to-end development of modeling strategies. To further illustrate our ideas, we present a walk-through example of mapping an application DTD to a relational application schema using Sangam.

**Keywords:** Metamodel, Model Management, Integration, Schema Transformation, XML Mapping

## 1 Introduction

**The Problem of Integration.** Networked environments like the Internet have catalyzed a phenomenal growth in the publication of data bringing with it an increasing need to share and integrate information. Alas, the format of data to be integrated varies from company to company and sometimes from person to person. To accomplish tasks such as data sharing, exchange, and integration, we may need to map an XML schema to a relational schema to drive transformation of XML elements into relational data or map an XML schema

---

<sup>†</sup>Sangam (Sangam) is a Hindi word meaning Junction.

of one application to that of another. Data sources may need to be mapped into data warehouse tables, or a query posed against a high-level semantic model may need to be mapped into an equivalent query posed against a logical database schema [HMN<sup>+</sup>99].

Today with new data formats continuously emerging, research and industry likewise have had to visit and re-visit the issues of integration. And in the era of electronic information exchange, as current technology turns into legacy systems at an ever increasing pace, this data model impedance presents an increasingly critical challenge. To combat this recurring problem, researchers [BR00, AT96, GL98] have looked at model management as a possible solution. For over a decade, researchers have modeled data models and application schemas to facilitate the translation process between different models. The translation between the models, however, was and is restricted to hard-wired code interspersed with generic rules to allow for some flexibility [MR83, AT96, BR00, PR95]. Moreover, these systems are not designed to be extensible. Hence, the addition of a new data model typically requires integration engineers to write from scratch, the translation of the new model to all other data models as well as the re-structuring within the same data model. The only re-use here is perhaps that of a user's knowledge and experience.

**Modeling Maps.** The key thrust of our work now lies in the *modeling of maps* in a model management system. Our goal is to alleviate the necessity of writing hard-wired translation code and to reach a new era of tackling future legacy system integrations for now and tomorrow. We do so by using a graph-theoretic framework [AT96] to model not only the data models, but to now also explicitly model the maps between the data models and between the application schemas. For this, we provide a generic set of map metaconstructs to describe (1) transformation of one data model to another - for example, map a relation in the relational model to an element or sub-element in the XML model; (2) re-structuring within the same data model - for example, combine two elements in a DTD to produce a new merged element in an output DTD; (3) re-structuring across schema-data boundaries - for example, map an attribute in the input to a relation in the output. We guarantee the correctness of these mappings, i.e., we can establish the correctness of the output data model (or output application schema) when given a "correct" input data model (or application schema). In this paper, we limit our discussions to translation to and from the relational model and the XML model (DTD) and re-structuring within these models. However, mapping between other models can be done in a similar manner.

To show the feasibility of our approach, we present  $\mathcal{S}$ angam our prototype system [CRZS01]. To the best of our knowledge, this is the first model management system that models and manages mappings and treats them as first-class objects. We illustrate the ease with which a user can not only map an application DTD to an relational schema using  $\mathcal{S}$ angam but also have the system generate the corresponding code for the data transformation.

### **Advantages of Our Approach.**

- *Automatic Generation* -  $\mathcal{S}$ angam offers declarative support for handling transformation strategies from one data model to another. These strategies are declared once at the model level by a system administrator. A key advantage of  $\mathcal{S}$ angam lies in the fact that mapping between application schemas and the subsequent mapping between the application data can be almost completely automated, thereby

increasing user productivity.

- *Correctness* - Based on our theoretical framework, we can guarantee the correctness of the mapping strategies, the mapping code that is generated to transform one application schema to another, as well as the transformation of the application data. Correctness here is based on the conformation to the output data model.
- *Plug-n-Play* - An integration engineer can edit and customize a map, for example, by plugging in a new function to combine two attributes together in a domain-specific manner. The system would still guarantee the correctness of the output as stated above.
- *Extensibility* - As new, perhaps more optimized, mapping strategies are discovered between existing data models, they can be added in *ᄎangam* with relatively little effort (by composing new map models), when compared to writing a transformation program from scratch.
- *Map Evaluation* - In contrast to often ad-hoc translation code, our graph-theoretic representation of maps provides a solid basis for reasoning about maps, such as their capacity, their similarity or their performance. This allows us to evaluate different mappings and select the right mapping for a given task.
- *Generic Map Operators* - A generic powerful set of operators [BR00] can now be defined for maps, allowing easier maintainability and powerful management of maps compared to working directly with translation programs. This set would include, for example, operators to copy portions of a map, to compose more complex maps, and to propagate schema changes on either the input or the output application schemas.

**Overview.** In Section 2, we present our four-tier map modeling architecture. Section 3 reviews background on modeling of data models. Section 4 gives formal descriptions of our map metaconstructs and Section 5 describes how these metaconstructs can be put together to form maps between data models and application schemas. In Section 6, we walk-through a XML to relational model example to show the feasibility and realization of our theoretical framework via *ᄎangam* our prototype system. Section 7 presents some related work and we conclude in Section 8.

## 2 A Four Layer Architecture for Map Modeling

UML has introduced a four-tier UML Metamodel architecture [Boo94] for modeling of data models. We now put forth our novel extension of the architecture to also enable modeling of maps between data models and application schemas.

**Modeling Data Models and Application Schemas.** While our work focuses on *modeling the mappings* between data models and application schemas, declaratively describing the latter is a necessity as these are the inputs and outputs of our maps. Following the four layer UML Metamodel architecture [Boo94], the data models, application schemas and data are represented in a four-tier architecture depicted in Figure 1. The

metamodel layer describes the metaconstructs that are used in the model layer to define the data models. The data models in turn define the structure of an application schema in the application layer. The application data resides in the bottom layer, and its structure is defined by the corresponding application schema in the application layer above. Figure 2 (a) represents a DTD model in the data model layer; Figure 2 (b) represents an application DTD (conforming to the DTD model) in the application layer. A detailed description of these appears in Section 3.

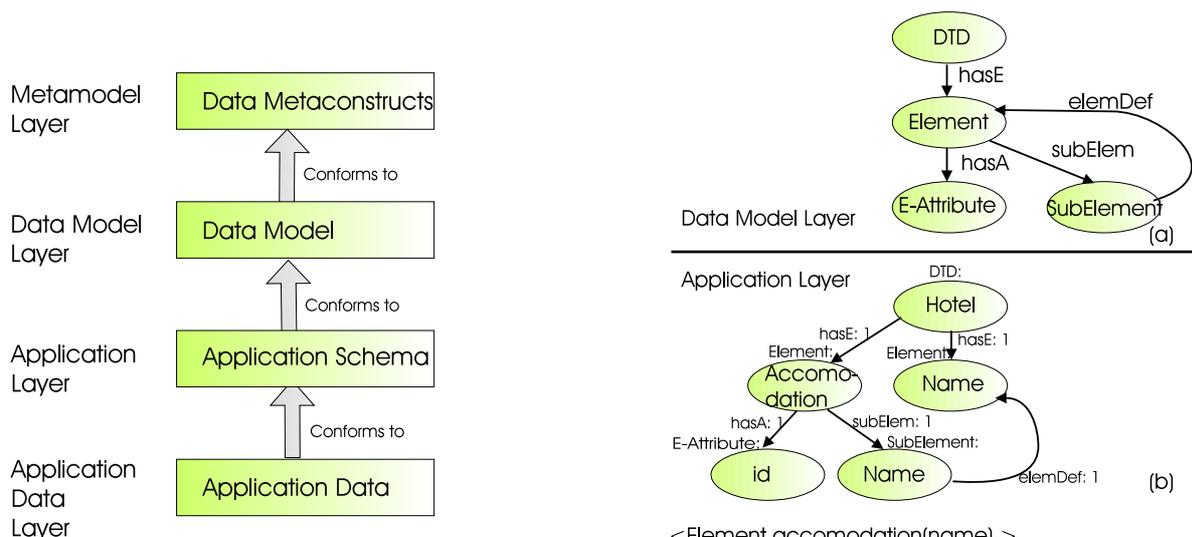


Figure 1: The Four Tier Architecture for Modeling Data Models and Application Schemas.

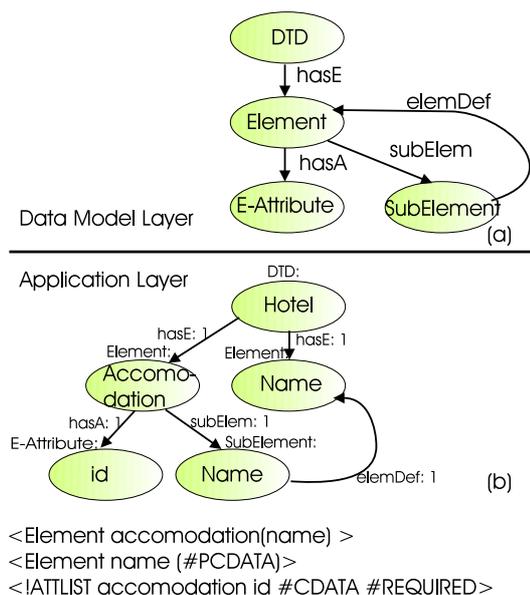


Figure 2: An Example Showing the DTD Model described in Data Model Layer and a DTD in the Application Layer.

**Modeling Model Maps and Application Maps.** We propose a novel extension to the familiar four tier architecture to now also model maps as shown in Figure 3. The map metaconstructs given in the top layer define the building blocks required for modeling maps. A map model composed of these map metaconstructs describes the translation of a given input data model to an output data model in the map model layer. For example, we can define a map model for translating the constructs of the XML (DTD) model (e.g., `element`) to constructs in the relational data model (e.g., `relation` or `attribute`).

In the application map layer, application maps conforming to the map model structure translate an input application schema to an output application schema as depicted in Figure 4. For example, an application map conforming to a given map model may map the element `Accommodation` to a relation `ACCOMMODATION`, or it may map the element `Name` to an attribute name in the relation `ACCOMMODATION`.

In the bottom layer, data maps describe the translation of the input application data to the output application data. This data transformation is guided by the information stored in the application map. Thus, for example, the values for element `Accommodation` for all XML documents conforming to the `Hotel` application DTD are translated to rows in the `ACCOMMODATION` relation; and the values for element `Name` are mapped to a value in a column of the relation `ACCOMMODATION`.

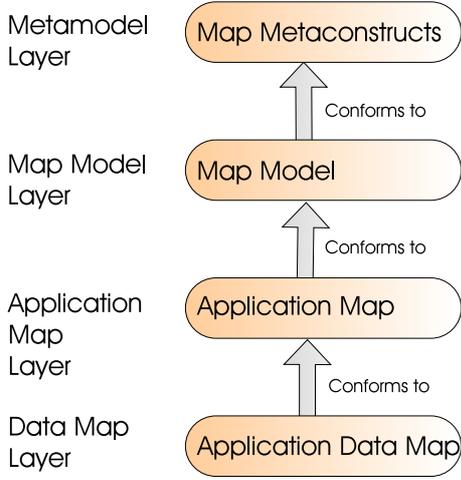


Figure 3: Our Proposed Four Tier Architecture for Modeling Map Models and Application Maps between Application Schemas.

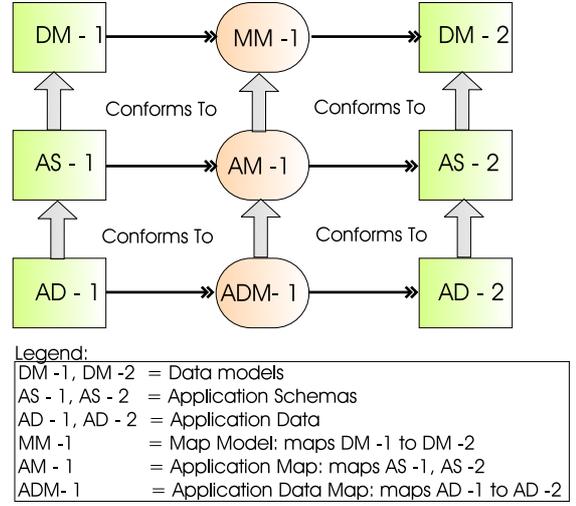


Figure 4: The Interaction between the Four-Tier Architecture for Data Models (Figure 1 and the Four-Tier Architecture for Map Modeling (Figure 3).

### 3 Background - Modeling Data Models and Application Schemas

In this section we briefly review the metamodel we use to express data models in the model layer and application schemas in the application layer. We utilize a subset of the graph-theoretic formalism as presented by Atzeni et al. [AT96]. Our metamodel is based on a fixed set of metaconstructs that are either nodes or edges. There are two node metaconstructs  $\mathcal{N}$ , namely *complex* ( $\square$ ) and *atomic* ( $\circ$ ) node types. Complex nodes have a set of outgoing edges while atomic nodes have no outgoing edges. There are two edge metaconstructs  $\mathcal{E}$ , namely *containment* ( $\rightarrow$ ) and *property* ( $--\rightarrow$ ) edges. A containment edge exists between two complex nodes while a property edge stems from a complex node and ends on an atomic node.

Throughout  $n \in \mathcal{N}$  refers to a node and  $e \in \mathcal{E}$  refers to an edge. The notation  $e: \langle n_1, n_2 \rangle$  refers to the two end-points of an edge. Here the edge  $e$  stems from  $n_1$  and ends on  $n_2$ . Given edges  $e_1: \langle n_1, n_2 \rangle$  and  $e_2: \langle n_2, n_3 \rangle$ , we denote the path from  $n_1$  to  $n_3$  by the path expression  $n_1 . e_1 . n_2 . e_2 . n_3$ .

**Patterns and Structures.** As in [AT96], two main notions are utilized to describe data models as well as application schemas. A *structure* is a directed acyclic graph whose nodes and edges are metaconstructs of our metamodel. A *pattern* is a rooted tree whose nodes and edges are metaconstructs of the metamodel and whose edges have *quantifiers* as labels. A quantifier given as a pair of integers  $[x:y]$ , with  $0 \leq x \leq y < \infty$ , specifies the minimum and maximum times an edge with the same label can appear in a structure.

**Definition 1 (Structure)** A structure is a triple  $S = (G, \mu, \epsilon)$  where  $G = (N, E)$  is a directed acyclic graph and  $\mu$  and  $\epsilon$  are typing functions:  $\mu: N \rightarrow \mathcal{N}$  and  $\epsilon: E \rightarrow \mathcal{E}$ .

Trees within the structure  $S$  composed of all the nodes and the edges reachable from the outer nodes<sup>1</sup> of  $S$  are referred to as *components* of  $S$ . One structure  $S_1 = ((N_1, E_1), \mu_1, \epsilon_1)$  is mapped to another structure  $S_2 = ((N_2, E_2), \mu_2, \epsilon_2)$  by a pair of functions  $\theta: N_1 \rightarrow N_2$  and  $\phi: E_1 \rightarrow E_2$ . Two structures are *isomorphic* if

<sup>1</sup>We use the term outer nodes to refer to nodes that have no edges incident on them.

for each  $n \in \mathbb{N}_1$ ,  $\mu_1(n) = \mu_2(\theta(n))$  and for each edge  $e: \langle n, n' \rangle \in E_1$ ,  $\phi(e) = e' : \langle \theta(n), \theta(n') \rangle \in E_2$  and  $\epsilon_1(e) = \epsilon_2(\phi(e))$ .

**Definition 2 (Pattern)** A pattern is a pair  $P = (S, \rho)$  where  $S = (G, \mu, \epsilon)$  is a structure such that  $G$  is a rooted tree, and  $\rho$  is a function that associates a quantifier with each edge of  $G$ .

A pattern describes a collection of structures that represent the same composition of the metaconstructs. The quantifier  $[x:y]$  associated with an edge in a pattern specifies the range, i.e., minimum ( $x$ ) and maximum ( $y$ ) number of times that the edge  $e$  can occur in a structure.

A structure  $S$  matches a pattern  $P = (S', \rho)$ , if for a node  $n \in S$  there is a set of edges in  $S$  denoted by  $E_n^e$  outgoing from  $n$  such  $\phi(e) = e'$  for some edge  $e' \in S'$ ; then, for each node  $n$  of  $S$  and for each set of edges  $E_n^e$ ,  $|E_n^e| \in \rho(e')$ .

A structure  $S$  is an instance of a set of patterns  $\mathcal{P}$ , if for each component  $S_i$  of a structure  $S$  there is a pattern  $P_i \in \mathcal{P}$  such that  $S_i$  matches  $P_i$ . Given a set of patterns  $\mathcal{P}$ ,  $Inst(\mathcal{P})$  denotes the set of structures that are instances of  $\mathcal{P}$ . Figure 5 represents a set of patterns and Figure 6 shows a structure that is an instance of patterns  $P_1$  and  $P_3$  in Figure 5.

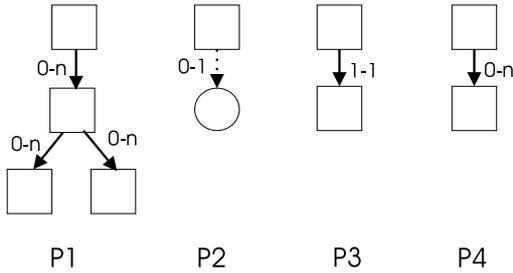


Figure 5: A Set of Patterns.

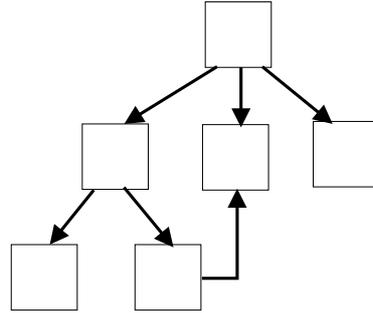


Figure 6: A Structure  $S_1$  - an Instance of Patterns  $P_1$  and  $P_3$  given in Figure 5.

Two patterns  $P_1 = (S_1, \rho_1)$  and  $P_2 = (S_2, \rho_2)$  are *isomorphic* if  $S_1$  and  $S_2$  are isomorphic and if for each edge  $e_1 \in P_1$ , there is  $e_2 \in P_2$  such that  $\phi(e_2) = e_1$  and  $\rho(e_1) = \rho(e_2)$ .

There is also a *subsumption* relationship between patterns, i.e., a set of patterns  $\mathcal{P}_1$  is subsumed by another set of patterns  $\mathcal{P}_2$ , if  $Inst(\mathcal{P}_1) \subseteq Inst(\mathcal{P}_2)$ .

**Data Model.** A data model  $DM = (\mathcal{P}, \gamma)$  is composed of a set of patterns  $\mathcal{P}$  with a labeling function  $\gamma$  that maps each element of  $\mathcal{P}$  of type  $\mathcal{N} \cup \mathcal{E}$  to a label  $l$  from a set of labels  $\mathcal{L}$ . These labels correspond to the names of constructs in a particular data model.

Figure 2 shows the segment of the DTD model obtained by assigning labels “*DTD*”, “*Element*”, “*E-Attribute*”, and “*SubElement*” to the first, second, third and fourth  $\square$  node respectively in the the pattern  $P_1$ ; and the labels “*SubElement*” and “*Element*” to the first and second  $\square$  node respectively in pattern  $P_3$  given in Figure 5. Similarly, labels are assigned to all the edges.

**Application Schema.** An application schema  $AS$  is a pair  $AS = (S, \lambda)$  composed of a structure  $S$  and a labeling function  $\lambda$  that maps each node and edge of  $S$  to a label, such that distinct labels are associated with

different elements of the structure. These labels correspond to real-world entities modeled by the application schema. A schema  $\mathcal{AS} = (S, \lambda)$  is an *instanceOf* a model  $\mathcal{DM} = (\mathcal{P}, \gamma)$  if the structure  $S$  is an *instanceOf*  $\mathcal{P}$ .

Figure 2 (b) shows a slice of the application DTD HOTEL that is obtained by assigning labels to the structure in Figure 6.

## 4 Metaconstructs for Modeling Maps

In this section we now describe our map metamodel and in particular its map metaconstructs. We describe the functionality and correctness of these map metaconstructs using the concept of patterns from Section 3. These (instances of) map metaconstructs can map between patterns as well as structures, and hence can map from one data model to another and also from one application schema to another as defined below. In this section we show how the map metaconstructs work for patterns and structures using examples of patterns and structures that are part of a data model or an application schema (i.e., they have labels). We also assume here that there is a global universe of discourse of all allowable patterns in the system denoted by  $\mathcal{P}$  and a set of all allowable structures in the system denoted by  $\text{Inst}(\mathcal{P})$ . This includes all sets of patterns that describe an individual data model and all structures that describe an application schema.

Our map metamodel is based on a fixed set of map node types  $\mathcal{M}$  and a fixed set of map edge types  $\xi$ . We define five types of map nodes: *cross* node ( $\otimes$ ), *merge* node ( $\oplus$ ), *ident* node ( $\ominus$ ), *project* node ( $\ominus$ ) and *connect* node ( $\ominus$ ). We also define three types of map edges: *input* ( $\rightarrow$ ), *output* ( $\rightarrow$ ) and *containment* ( $\rightarrow$ ).

### 4.1 Map Nodes

Each map node captures the semantics of some translation of a given input pattern to some output pattern, with their specific mappings and requirements for the input and output patterns as described below. We denote the mapping of a node  $n$  in an input pattern  $\mathcal{P}$  to a node  $n'$  in the output pattern  $\mathcal{P}'$  by  $\theta(n) = n'$  as given in Section 3. Similarly, an edge  $e \in \mathcal{P}$  is mapped to an edge  $e' \in \mathcal{P}'$  by  $\phi(e) = e'$ . For each map node we now give more precise definitions of  $\theta$  and  $\phi$ .

#### 4.1.1 Cross Node

**Cross Node for Patterns.** The cross map node is the simplest map node. It maps an input pattern consisting of one node (root) to an output pattern of exactly one node (root). Its functionality enables the mapping of one data model construct to another. For example, a cross node can map a node labeled *Element* in the XML model to a node labeled *Relation* in the relational model. Figure 7 (a) shows the cross mapping of a pattern  $\mathcal{P}$  to a pattern  $\mathcal{P}'$ .

**Definition 3 (Cross Node - Pattern)** *Given an input pattern  $\mathcal{P}$  with  $G = (N, E)$ ,  $|N| = 1$  and  $|E| = 0$ , a cross node, denoted by  $\otimes$ , produces as output a pattern  $\mathcal{P}'$  with  $G' = (N', E')$ ,  $|N'| = 1$  and  $|E'| = 0$  and function  $\theta: N \rightarrow N'$ .*

**Pattern Correctness of Cross Node.** Given an input pattern  $\mathcal{P}_1 \in \mathcal{P}$  with  $G_1 = (N_1, E_1)$  a cross node  $\otimes$  is said to be correct if it produces as output some pattern  $\mathcal{P}_2$  with  $G_2 = (N_2, E_2)$ , such that  $\mathcal{P}_2 \in \mathcal{P}$ .

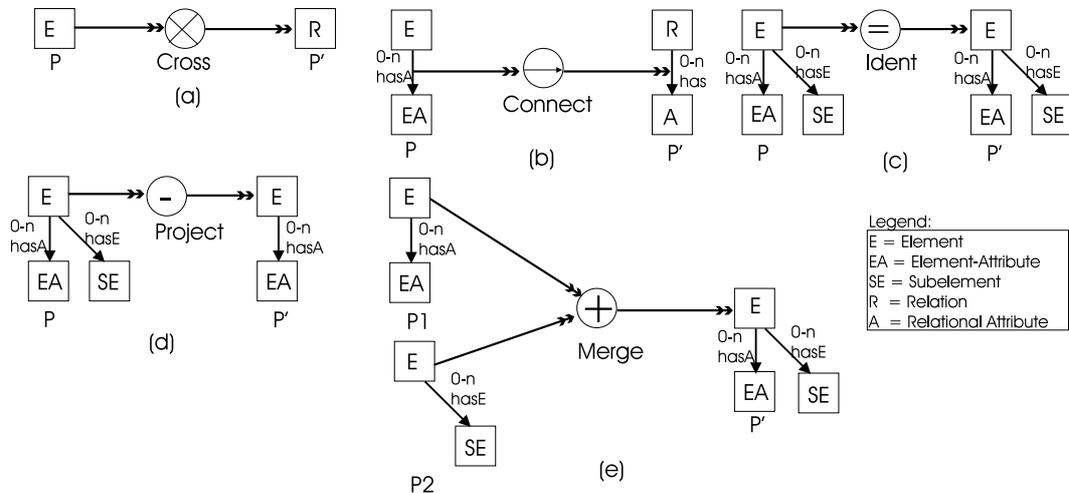


Figure 7: The Five Map Metaconstructs For Patterns.

**Cross Node for Structures.** At the application map layer, (an instance of) a cross node maps an input structure to produce an output structure, thereby allowing the mapping of one application schema to another. For example, Figure 8 (a) shows the mapping of an XML element `Accomodation` to a relation `ACCOMODATION`.

**Definition 4 (Cross Node - Structure)** Given an input structure  $S$  with  $G = (N, E)$ ,  $|N| = 1$  and  $|E| = 0$ , a cross node, denoted by  $\otimes$ , produces as output a structure  $S'$  with  $G' = (N', E')$ ,  $|N'| = 1$  and  $|E'| = 0$  and function  $\theta: N \rightarrow N'$ .

**Structure Correctness of Cross Node.** Given an input structure  $S_1 \in \text{Inst}(\mathcal{P})$  with  $G_1 = (N_1, E_1)$ , a cross node  $\otimes$  is said to be correct if it produces as output some structure  $S_2$  with  $G_2 = (N_2, E_2)$  such that  $S_2 \in \text{Inst}(\mathcal{P})$ .

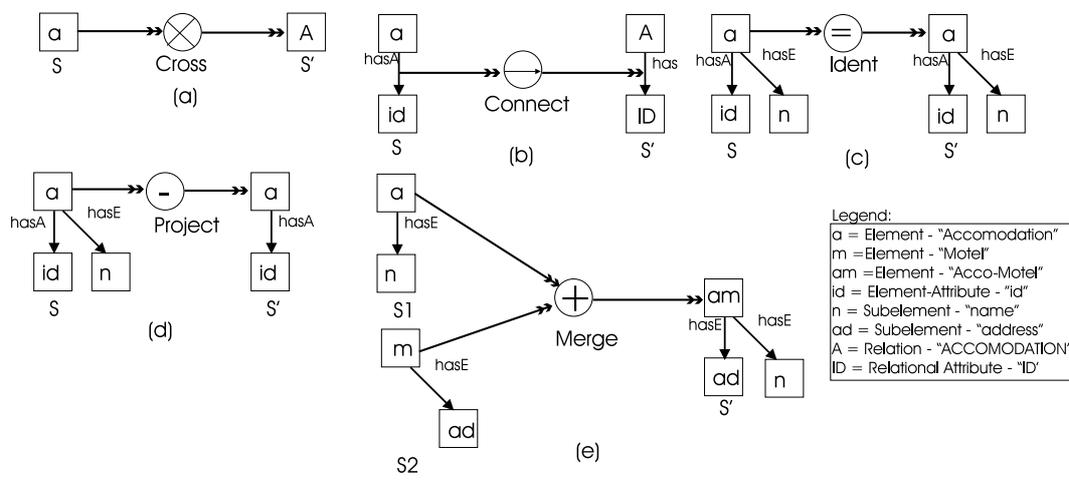


Figure 8: The Five Map Metaconstructs For Structures.

### 4.1.2 Connect Node

In general, the connect node enables the mapping of an edge  $e_1: \langle n_1, n_2 \rangle$  in an input pattern or structure to another edge  $e_2: \langle n_3, n_4 \rangle$  in the output pattern or structure respectively.

**Connect Node for Patterns.** The connect node enables the mapping of an edge  $e_1: \langle n_1, n_2 \rangle$  in an input pattern to an edge  $e_2: \langle n_3, n_4 \rangle$  in the output pattern. For example, a connect node can map the *Subelem* relationship between an *Element* and its *Subelement* in the XML model to the *ForeignKey* relationship between two *Relations* in the relational model. Input to a connect node can be a direct edge between two nodes such as *has*, *hasA*, or *Subelem*; or a path expression such as *Element.SubElem.SubElement.ElemDef* that refers to the definition of an *Element* that defines a *Subelement*. Output of an edge map node is restricted to be a direct edge<sup>2</sup>.

For correctness of the connect node, we require that the end-points of the edge  $e_1$  be mapped to the end-points of the output edge  $e_1'$  to avoid dangling edges. Thus, node  $n_1$  must be mapped to either  $n_1'$  or  $n_2'$ . If  $n_1$  is mapped to  $n_1'$ , then  $n_2$  must be mapped to  $n_2'$ ; or vice versa. Hence, it is not meaningful to map the *has* edge between the *Relation* and the *Attribute* to the *hasA* edge between an *Element* and *E-Attribute* without also mapping the nodes appropriately.

**Definition 5 (Connect Node - Pattern)** Given a pattern  $\mathcal{P}$  with  $G = (N, E)$ ,  $|N| = 2$ ,  $|E| = 1$ , and edge  $e: \langle n_1, n_2 \rangle \in \mathcal{P}$ , such that  $n_1, n_2 \in N$  and  $n_1 \text{ noteq } n_2$ , a connect node, denoted by  $\oplus$ , maps  $\mathcal{P}$  to an output pattern  $\mathcal{P}'$ , with  $G' = (N', E')$ , and  $|N'| = 2$ ,  $|E'| = 1$ , via function  $\phi: E \rightarrow E'$  such that  $\phi(e) = e': \langle n_1', n_2' \rangle$  and  $e' \in E'$ ,  $n_1', n_2' \in N'$ . Also we require  $\theta(n_1) = n_1'$  and  $\theta(n_2) = n_2'$  or vice versa.

**Pattern Correctness of Connect Node.** Given an input pattern  $\mathcal{P} \in \mathcal{P}$  with edge  $e: \langle n_1, n_2 \rangle \in \mathcal{P}$ , the connect node is said to be correct if the output pattern  $\mathcal{P}'$  it produces is in  $\mathcal{P}$ .

**Connect Node for Structures.** The connect node to map between two structures is similar in functionality to that between two patterns. Hence, a connect node maps edge  $e: \langle n_1, n_2 \rangle$  in an input structure to another edge  $e': \langle n_1', n_2' \rangle$  in the output structure. Figure 8 shows a connect node that maps the edge between the XML element *Accomodation* and its attribute *id* to the edge between the relation *ACCOMODATION* and its attribute *ID*. Similar to the connect node between patterns, a direct edge or a path expression in the input structure is valid input for the (instance of) connect node in the application layer.

For correctness of the connect node here, we again stipulate as for the connect node between patterns, that the end-points of the edge in the input structure must be mapped to the end-points of the edge in the output structure.

**Definition 6 (Connect Node - Structure)** Given a structure  $\mathcal{S}$  with  $G = (N, E)$ ,  $|N| = 2$ ,  $|E| = 1$ , and edge  $e: \langle n_1, n_2 \rangle \in \mathcal{S}$ , such that  $n_1, n_2 \in N$ , a connect node, denoted by  $\oplus$ , maps  $\mathcal{S}$  to an output structure  $\mathcal{S}'$ , with  $G' = (N', E')$ , and  $|N'| = 2$ ,  $|E'| = 1$ , via function  $\phi: E \rightarrow E'$  such that  $\phi(e) = e': \langle n_1', n_2' \rangle$  and  $e' \in E'$ ,  $n_1', n_2' \in N'$ . Also we require  $\theta(n_1) = n_1'$  and  $\theta(n_2) = n_2'$  or vice versa.

<sup>2</sup>Path expressions on the output node can be simulated by a combination of map nodes. Thus, for example it is possible to map a direct input edge to an output path expression (splitting of an edge) by combining several map nodes.

**Structure Correctness of Connect Node.** Given an input structure  $S \in \text{Inst}(\mathcal{P})$  with edge  $e: \langle n_1, n_2 \rangle \in S$ , the connect node is said to be correct if the output structure  $S'$  it produces is also in  $\text{Inst}(\mathcal{P})$ .

### 4.1.3 Ident Node

The ident ( $\ominus$ ) takes as input a pattern or a structure and produces as output an exact copy of the input pattern or structure respectively.

**Ident Node for Patterns.** The ident ( $\ominus$ ) takes as input a pattern  $P$  with a complex root node and produces as output a pattern  $P'$  such that  $P$  and  $P'$  are isomorphic.

**Definition 7 (Ident Node - Pattern)** *Given an input pattern  $P$  with  $G = (N, E)$ ,  $|N| > 1$  and  $nr \in N$  is the root, and output pattern  $P'$  with  $G' = (N', E')$ , an ident node, denoted by  $\ominus$ , maps  $P$  to  $P'$  via a bijection  $\theta: N \rightarrow N'$  such that for every  $e: \langle n_1, n_2 \rangle \in E$  there exists an edge  $e': (\theta(n_1), \theta(n_2)) \in E'$  and  $nr' \in N'$  is root of pattern  $P'$  with  $\theta: nr \rightarrow nr'$  and quantifier  $\rho(e) = \rho(e')$ .*

**Pattern Correctness of Ident Node.** Given a pattern  $P \in \mathcal{P}$  as input, an ident node  $\ominus$  is said to be correct if it maps to an output pattern  $P'$  such that  $P$  and  $P'$  are isomorphic<sup>3</sup>.

**Ident Node for Structures.** The ident ( $\ominus$ ) takes as input a structure  $S$  with a complex root node and produces as output a pattern  $S'$  such that  $S$  and  $S'$  are isomorphic.

**Definition 8 (Ident Node - Structure)** *Given an input structure  $S$  with  $G = (N, E)$ ,  $|N| > 1$  and  $nr \in N$  is the root, and output structure  $S'$  with  $G' = (N', E')$ , an ident node, denoted by  $\ominus$ , maps  $S$  to  $S'$  via a bijection  $\theta: N \rightarrow N'$  such that for every  $e: \langle n_1, n_2 \rangle \in E$  there exists an edge  $e': (\theta(n_1), \theta(n_2)) \in E'$  and  $nr' \in N'$  is root of structure  $S'$  with  $\theta: nr \rightarrow nr'$ .*

**Structure Correctness of Ident Node.** Given a structure  $S \in \text{Inst}(\mathcal{P})$  as input, an ident node  $\ominus$  is said to be correct if it maps to an output structure  $S'$  such that  $S$  and  $S'$  are isomorphic.

### 4.1.4 Merge Node

The merge node ( $\oplus$ ) maps either two input patterns or two input structures to produce as output one pattern or structure respectively.

**Merge Node for Patterns.** The merge node ( $\oplus$ ) maps two input patterns  $P_1$  and  $P_2$  with roots  $nr_1$  and  $nr_2$  respectively to one output pattern  $P'$  with root  $nr'$ . The set of patterns  $\mathcal{P}_1$  and  $\mathcal{P}_2$  reachable from the roots  $nr_1$  and  $nr_2$  and subsumed by the input patterns  $P_1$  and  $P_2$  respectively are mapped to the output pattern  $P'$  such that  $\mathcal{P}_1 \cup \mathcal{P}_2$  are subsumed by  $P'$  and are reachable from its root  $nr'$ . Consider the example given in Figure 7 (e). Here the merge node maps the patterns  $P_1$  and  $P_2$  to produce pattern  $P'$  such that all components stemming from the roots of  $P_1$  and  $P_2$  are subsumed in the output structure and are reachable from the root of the new pattern  $P'$ .

---

<sup>3</sup>Isomorphism for patterns and structures is defined in Section 3.

**Definition 9 (Merge Node - Pattern)** <sup>4</sup> Let  $P_1$  and  $P_2$  be two patterns. Here  $P_1$  has  $G_1 = (N_1, E_1)$ ,  $|N_1| \geq 1$  and  $nr_1 \in N_1$  the root. Let  $\mathcal{P}_1$  be the set of disjoint patterns rooted in  $nr_1$  and subsumed by  $P_1$ . Let  $N_1$  and  $E_1$  represent the set of nodes and edges of  $\mathcal{P}_1$ . Let  $P_2$  and  $\mathcal{P}_2$  be defined similar to  $P_1$  and  $\mathcal{P}_1$ .

Given the patterns  $P_1$  and  $P_2$  as input, a merge node, denoted by  $\oplus$ , maps  $P_1$  and  $P_2$  to  $P'$ , with  $G' = (N', E')$  and root node  $nr'$ , by  $\theta: N_1 \cup N_2 \rightarrow N'$  such that  $\theta(nr_1) = nr'$ ;  $\theta(nr_2) = nr'$ ; and for every  $e: (nr, n_2) \in E_1 \cup E_2$ ,  $nr = nr_1$  or  $nr = nr_2$ , there exists an edge  $e': (nr', \theta(n_2)) \in E'$ ; and for every  $e: (n_1, n_2) \in E_1 \cup E_2$ , there exists an edge  $e': (\theta(n_1), \theta(n_2)) \in E'$ . If  $e_1 \in E_1$  and  $e_2 \in E_2$  such that  $\phi(e_1) = \phi(e_2) = e'$ ,  $e' \in E'$ , then  $\rho(e') = \max(\rho(e_1)) + \max(\rho(e_2))$ <sup>5</sup>.

**Pattern Correctness of the Merge Node.** Let  $P_1, P_2 \in \mathcal{P}$  with roots  $nr_1$  and  $nr_2$  respectively be two patterns. Let  $P_{r1}$  and  $P_{r2}$  be patterns subsumed by  $P_1$  and  $P_2$  respectively such that  $|N_{r1}| = 1$  and  $nr_1 \in N_{r1}$ , and  $|N_{r2}| = 1$  and  $nr_2 \in N_{r2}$ . Given two patterns  $P_1$  and  $P_2$ , a merge node is said to be correct if it produces an output pattern  $P'$ , such that  $P' \in \mathcal{P}$  and  $P_{r1}$  and  $P_{r2}$  are subsumed in  $P'$ .

**Merge Node for Structures.** The merge node ( $\oplus$ ) maps two input structure  $S_1$  and  $S_2$  with roots  $nr_1$  and  $nr_2$  respectively to one output structure  $S_o$  with root  $nr_o$ . The set of components  $\mathcal{S}_1$  and  $\mathcal{S}_2$  reachable from the outer nodes (roots)  $nr_1$  and  $nr_2$  are mapped to the output structure  $S_o$  such that  $\mathcal{S}_1 \cup \mathcal{S}_2$  are components reachable from the outer node (root)  $nr_o$  of the output structure  $S_o$ . Consider the example given in Figure 8 (e). Here the merge node maps the structures  $S_1$  and  $S_2$  to produce structure  $S'$  such that all components stemming from the roots of  $S_1$  and  $S_2$  are components in the output structure  $S'$ .

**Definition 10 (Merge Node - Structure)** Let  $S_1$  and  $S_2$  be two structures. Here  $S_1$  has  $G_1 = (N_1, E_1)$ ,  $|N_1| \geq 1$  and  $nr_1 \in N_1$  the root. Let  $\mathcal{S}_1$  be the set of components rooted in  $nr_1$ . Let  $N_1$  and  $E_1$  represent the set of nodes and edges of  $\mathcal{S}_1$ . Let  $S_2$  and  $\mathcal{S}_2$  be defined similar to  $S_1$  and  $\mathcal{S}_1$ .

Given the structures  $S_1$  and  $S_2$  as input, a merge node, denoted by  $\oplus$ , maps  $S_1$  and  $S_2$  to  $S'$ , with  $G' = (N', E')$  and root node  $nr'$ , by  $\theta: N_1 \cup N_2 \rightarrow N'$  such that  $\theta(nr_1) = nr'$ ;  $\theta(nr_2) = nr'$ ; and for every  $e: (nr, n_2) \in E_1 \cup E_2$ ,  $nr = nr_1$  or  $nr = nr_2$ , there exists an edge  $e': (nr', \theta(n_2)) \in E'$ ; and for every  $e: (n_1, n_2) \in E_1 \cup E_2$ , there exists an edge  $e': (\theta(n_1), \theta(n_2)) \in E'$ .

**Structure Correctness of the Merge Node.** Let  $S_1, S_2$  with roots  $nr_1$  and  $nr_2$  respectively be two structures such that  $S_1$  and  $S_2 \in \text{Inst}(P)$ , i.e., they are instances of the same pattern  $P$ . Let  $S_{r1}$  and  $S_{r2}$  be components of  $S_1$  and  $S_2$  respectively such that  $|N_{r1}| = 1$  and  $nr_1 \in N_{r1}$ , and  $|N_{r2}| = 1$  and  $nr_2 \in N_{r2}$ . Given the two structures  $S_1$  and  $S_2$ , a merge node is said to be correct if it produces an output structure  $S'$ , such that  $S'$  is also an instance of the same pattern  $P$  and the set of components  $S_{r1}$  and  $S_{r2}$  are also components of  $S'$ .

#### 4.1.5 Project Node

A project node maps a set of patterns subsumed in a pattern to an output pattern or a set of components of an input structure to an output structure.

<sup>4</sup>This is a very general definition of the merge node. More specific merge nodes can be described similarly.

<sup>5</sup> $\max$  returns the maximum quantifier from the quantifier  $[\min:\max]$  pair.

**Project Node for Patterns.** Let  $\mathcal{P}$  be a pattern with root  $nr$ . The project node maps a set of patterns  $\mathcal{P}_1$  subsumed in pattern  $\mathcal{P}$  such that all patterns  $p_1 \in \mathcal{P}_1$  are rooted at  $nr$  to an output pattern  $\mathcal{P}'$ . Consider again the pattern in Figure 2 (a) where *Element* has two edges *hasA* and *Subelem*. A project can select an edge, for example *hasA*, and map it to the output pattern that now includes the root *Element*, the edge *hasA* and the sub-pattern on which edge *hasA* is incident, eliminating the edge *Subelem*. The project node is depicted in Figure 7 (d).

**Definition 11 (Project Node - Pattern)** Let  $\mathcal{P}$  be a pattern with  $G = (N, E)$ ,  $|N| \geq 1$ , and  $nr \in N$  be the root of  $\mathcal{P}$ . Let  $\mathcal{P}_1$  be the set of disjoint patterns rooted in  $nr$  and subsumed in  $\mathcal{P}$ . Let  $\mathbb{N}_1$  and  $\mathbb{E}_1$  represent the set of nodes and edges respectively for  $\mathcal{P}_1$ . Given  $\mathcal{P}$  as the input pattern, a project node, denoted by  $\ominus$ , maps  $\mathcal{P}$  to an output pattern  $\mathcal{P}'$ , with  $G' = (N', E')$ , by the bijection  $\theta: \mathbb{N}_1 \rightarrow N'$  such that for every  $e: (n_1, n_2) \in \mathbb{E}_1$  there exists an edge  $e': (\theta(n_1), \theta(n_2)) \in E'$ , and  $\rho(e) = \rho(e')$ ;  $nr' \in N'$  is root of pattern  $\mathcal{P}'$  and  $\theta(nr) = nr'$ .

**Pattern Correctness of Project Node.** Given a pattern  $\mathcal{P} \in \mathcal{P}$  as input, a project node is said to be correct if it maps to an output pattern  $\mathcal{P}'$  such that  $\mathcal{P}'$  is subsumed by  $\mathcal{P}$  and  $\mathcal{P}'$  in  $\mathcal{P}$ .

**Project Node for Structures.** Let  $S$  be a structure with root  $nr$ . The project node maps a set of components  $\mathcal{S}_1$  of structure  $S$  to an output structure  $S'$  such that the set of components  $\mathcal{S}'_1$  of  $S'$  is a subset of  $\mathcal{S}_1$ . The project node is depicted in Figure 8 (d).

**Definition 12 (Project Node - Structure)** Let  $S$  be a structure with  $G = (N, E)$ ,  $|N| \geq 1$ , and  $nr \in N$  be the root of  $S$ . Let  $\mathcal{S}_1$  be the set of components rooted in  $nr$ . Let  $\mathbb{N}_1$  and  $\mathbb{E}_1$  represent the set of nodes and edges respectively for  $\mathcal{S}_1$ . Given  $S$  as the input structure, a project node, denoted by  $\ominus$ , maps  $S$  to an output structure  $S'$ , with  $G' = (N', E')$ , by the bijection  $\theta: \mathbb{N}_1 \rightarrow N'$  such that for every  $e: (n_1, n_2) \in \mathbb{E}_1$  there exists an edge  $e': (\theta(n_1), \theta(n_2)) \in E'$ ;  $nr' \in N'$  is root of structure  $S'$  and  $\theta(nr) = nr'$ .

**Structure Correctness of Project Node.** Given a structure  $S \in \text{Inst}(\mathcal{P})$  as input, a project node is said to be correct if it maps to an output structure  $S'$  such that  $S'$  is subsumed by  $S$  and  $S'$  in  $\text{Inst}(\mathcal{P})$ .

## 4.2 Map Edges

**Input and Output Edges.** Each map node is required to always have at least one input edge ( $\rightarrow$ ) and one output edge ( $\rightarrow$ ) that connect to patterns or structures of a data model or an application schema respectively. For example, when a cross node maps an input pattern  $\mathcal{P}_1$  with a node labeled *Relation* to a pattern  $\mathcal{P}_2$  with a node labeled *Element*, then the cross node has one input edge that ends at (node of)  $\mathcal{P}_1$  and an output edge that is incident on (a node of)  $\mathcal{P}_2$ . The input and the output edges of a map denote here that a *Relation* is mapped to an *Element*.

**Containment Edges.** A containment edge from map node  $mn_1$  to map node  $mn_2$  indicates that the map node  $mn_2$  is part of the map node  $mn_1$ . Containment edges between map nodes often (but not necessarily always) reflect containment edges between two nodes in the input pattern  $\mathcal{P}_1$ . For example, consider two

cross map nodes such that  $mn_1$  maps an *Element* to a *Relation* ( $Element \rightarrow \otimes \rightarrow Relation$ ); and  $mn_2$  maps *E-Attribute* to *Attribute* ( $E-Attribute \rightarrow \otimes \rightarrow Attribute$ ). In that case, we may also have a connect node  $mn_3$  to map the edge *hasA* to *has* ( $hasA \rightarrow \ominus \rightarrow has$ ). Here  $mn_1 \rightarrow mn_2$  and  $mn_1 \rightarrow mn_3$ , that is the nodes  $mn_2$  and  $mn_3$  are contained in map node  $mn_1$ .

## 5 Map Model and Application Maps

In this section we give a brief overview of the graph-based map metamodel used to express map models and the application maps. Throughout we use  $mn$  to refer to a map node and  $me$  to refer to a map edge.

### 5.1 Map Patterns and Map Structures

#### 5.1.1 Map Structure

We now extend and apply the concepts from Section 3 to maps, that is, we introduce the notion of *map patterns* and *map structures* for map modeling. A *map structure* is a directed acyclic graph whose nodes and edges are elements of our set of map metaconstructs  $\mathcal{M} \cup \xi$ . Roughly speaking, a map structure maps an input structure to an output structure. Figure 9 (b) depicts a map structure.

**Definition 13 (Map Structure)** *A structure is a five-tuple  $MS = (MG, \nu, \varepsilon, S_1, S_2)$  where  $MG = (MN, ME)$  is a directed acyclic graph, and  $\nu$  and  $\varepsilon$  are typing functions:  $\nu: MN \rightarrow \mathcal{M}$  and  $\varepsilon: ME \rightarrow \xi$  and  $S_1$  and  $S_2$  are the input and output structures as defined below.*

Let  $ME_{mn}^i$  be the set of input ( $\rightarrow$ ) edges for map node  $mn \in MN$ . A structure  $S_1 = ((N_1, E_1), \mu_1, \epsilon_1)$  is defined as an *input structure* for a map structure  $MS = (MG, \nu, \varepsilon, S_1, S_2)$  if for every map node  $mn$  each input edge in  $ME_{mn}^i$  is incident on either a node  $n \in N_1$  or an edge  $e \in E_1$ . Moreover, every node and edge in  $S_1$  has at least one input edge from some map node  $mn \in MS$  incident on it.

Let  $ME_{mn}^o$  be the set of output edges for map node  $mn \in MN$ . A structure  $S_2 = ((N_2, E_2), \mu_2, \epsilon_2)$  is defined as an *output structure* for a map structure  $MS = (MG, \nu, \varepsilon, S_1, S_2)$  if for every map node  $mn$  each output edge in  $ME_{mn}^o$  is incident on either a node  $n \in N_2$  or an edge  $e \in E_2$ . Moreover, every node and edge in  $P_2$  has at least one output edge from some map node  $mn \in MS$  incident on it.

We use the same notion of *components* for map structures as defined for structures in Section 3.

**Correctness of Map Structures.** A map structure  $MS$  is said to be correct if for a given input structure  $S_1 \in Inst(\mathcal{P})$ , it maps to an output structure  $S_2 \in Inst(\mathcal{P})$  and all map nodes  $mn \in MN$  are correct by Section 4.

**Properties of Map Structures.** Consider two map structures  $MS_1 = (MG_1, \nu_1, \varepsilon_1, S_1, S_2)$  and  $MS_2 = (MG_2, \nu_2, \varepsilon_2, S_3, S_4)$ , then  $MS_1$  can be *mapped* to  $MS_2$  via a pair of functions  $\alpha: MN_1 \rightarrow MN_2$  and  $\beta: ME_1 \rightarrow ME_2$ . Two map structures are *equivalent* if for each  $mn \in MN_1$ ,  $\nu_1(mn) = \nu_2(\alpha(mn))$  and for each edge  $me: (mn, mn') \in ME_1$ ,  $\varepsilon_1(me) = \varepsilon_2(\beta(me))$  and  $\beta(me) = (\alpha(mn), \alpha(mn'))$ . Two map structures  $MS_1$  and  $MS_2$  are *isomorphic* if  $MS_1$  and  $MS_2$  are equivalent and if the two inputs  $S_1$  and  $S_3$  of the structures  $MS_1$  and  $MS_2$  respectively are isomorphic.

### 5.1.2 Map Pattern

Like a pattern defined in Definition 14, a map pattern describes a collection of structures that represents a specific composition of the map metaconstructs. The quantifier in a map pattern specifies the range, i.e., minimum and maximum number of times, an edge can occur in the corresponding map structure. Input and output of the map pattern are patterns as described in Section 3 implying that input and output edges from map nodes in the map pattern must be incident on nodes in the patterns. Figure 9 (a) depicts a map pattern.

**Definition 14 (Map Pattern)** A map pattern is a four-tuple  $MP = (MS, \varrho, P_1, P_2)$  where  $MS = (MG, \nu, \varepsilon, S_1, S_2)$  is a map structure such that  $MG$  is a rooted tree, and  $\varrho$  is a function that associates a quantifier with each edge of  $MG$ ; and  $P_1$  and  $P_2$  are input and output patterns, respectively.

Let  $ME_{mn}^i$  be the set of input edges for map node  $mn \in MN$ . A pattern  $P_1 = (S_1, \rho_1)$  is defined as an *input pattern* for a map pattern  $MP$  if for every map node  $mn \in MN$  each input edge in  $ME_{mn}^i$  is incident on either a node  $n \in N_1$  or an edge  $e \in E_1$ . Every node and edge in the input pattern must have at least one input edge from some map node  $mn \in MP$  incident on it.

Let  $ME_{mn}^o$  be the set of output edges for map node  $mn \in MN$ . A pattern  $P_2 = (S_2, \rho_2)$  is defined as an *output pattern* for a map pattern  $MP$  if for every map node  $mn$ , each output edge in  $ME_{mn}^o$  is incident on either a node  $n \in N_2$  or an edge  $e \in E_2$ . Moreover, every node  $n$  in  $N_2$  and every edge  $e \in E_2$  has at least one output edge from some map node  $mn \in MP$ .

**Correctness of Map Patterns.** A map pattern  $MS$  is correct if for a given input pattern  $P_1 \in \mathcal{P}$ , it outputs a pattern  $P_2$  such that  $P_2 \in \mathcal{P}$  and all map nodes  $mn \in MN$  are correct by Section 4.

**Properties of Map Patterns.** Two map patterns  $MP_1$  and  $MP_2$  are *isomorphic* if  $MS_1$  and  $MS_2$  are isomorphic and if each map edge  $me_1$  in  $MP_1$ ,  $me_1 = \beta(me_2)$  and  $me_2 \in MP_2$  and  $\varrho(me_1) = \varrho(me_2)$ .

**Map Structures and Map Patterns.** The *instanceOf* relationship between a map structure and a map pattern is defined similar to that between a structure and a pattern. We first need to be able to show that a map structure  $MS_i$  can be represented by some pattern  $MP_i$  in a given set of patterns  $\mathcal{MP}$ . For this we first define a *map match*.

**Definition 15 (Map Match)** A map structure  $MS = ((MN_s, ME_s), \nu, \varepsilon, S_1, S_2)$  **matches** a map pattern  $MP = (MS', \varrho', P_1', P_2')$  where  $MS' = (MG', \nu', \varepsilon', S_1', S_2')$  if for a map node  $mn \in MN_s$  there is a set of edges  $ME_{mn}^{me}$  stemming from the  $mn$  such that they map to the same edge  $me' \in ME'$  i.e.,  $\beta(me) = me'$  for  $me' \in ME'$ , and  $|ME_{mn}^{me}| \in \varrho(me')$ ;  $S_1$  is *instanceOf*  $P_1$  and  $S_2$  is *instanceOf*  $P_2$ .

A map structure  $MS$  is an *instanceOf* a set of map patterns  $\mathcal{MP}_1$  if for each component  $MS_i$  of a map structure  $MS$ , there is a pattern  $MP_i \in \mathcal{MP}_1$  such that  $MS_i$  matches  $MP_i$  by Definition 15.

We also define a *subsumption* relationship between map patterns, i.e., a set of map patterns  $\mathcal{MP}_1$  is subsumed by another set of map patterns  $\mathcal{MP}_2$ , if for each map pattern  $MP_1 \in \mathcal{MP}_1$  there is a map pattern  $MP_2 \in \mathcal{MP}_2$  such that  $MP_1$  and  $MP_2$  are isomorphic; and the input pattern  $P_{i1}$  for  $MP_1$  is subsumed by the input pattern  $P_{i2}$  for map pattern  $MP_2$ ; and the output pattern  $P_{o1}$  for  $MP_1$  is subsumed by the output pattern  $P_{o2}$  for map pattern  $MP_2$ .

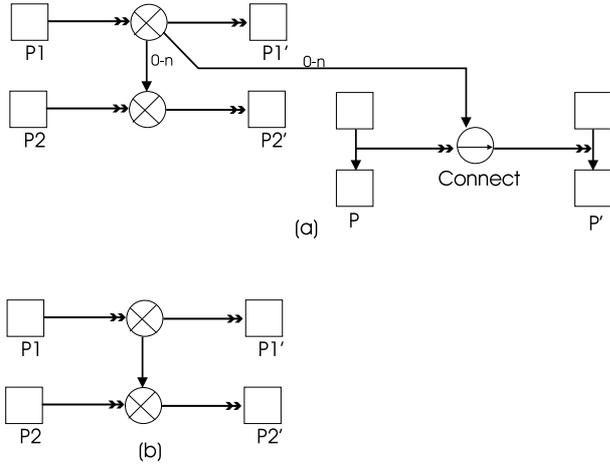


Figure 9: (a) An Example Map Pattern and (b) Map Structure.

## 5.2 Map Model and Application Model

**Map Model.** A map model  $MM = (\mathcal{MP}, \eta, \delta)$  is composed of a set of map patterns  $\mathcal{MP}$  with a labeling function  $\eta$  that maps the elements of  $\mathcal{MP}$  which are of types  $\mathcal{M} \cup \xi$  to a set of labels; and a constraint function  $\delta$  that maps the map edges  $me \in \mathcal{MP}$  to a set of constraints. Constraints are conditions applied primarily for the input and output edges of the map nodes. A constraint is satisfied if the condition holds true. The input and output patterns for each map pattern are also mapped via the labeling function  $\gamma$  to a set of labels (see the definition of a data model). Intuitively, a map model  $MM$  maps (a subset of) a data model  $DM_1$  to (a subset of) another data model  $DM_2$ . Figure 10 depicts a map model that assigns labels to the map pattern represented in Figure 9 (a). This map model maps from DTD to the relational model.

**Correctness of Map Model.** A map model  $MM$  is correct from  $DM_1$  and  $DM_2$ , if (1) all map patterns  $MP \in \mathcal{MP}$  are correct; (2) all constraints mapped via the constraint function  $\delta$  are satisfied; and (3) for a given input model  $DM_1$  the map model produces a set of patterns  $\mathcal{P}_1$  that is a subset of the set of patterns  $\mathcal{P}_{DM_2}$  for  $DM_2$ .

**Application Map.** An application map  $AM = (AS, \kappa, \zeta)$  is defined by a map structure  $AS$ , a labeling function  $\kappa$  that maps each map node and map edge of  $AS$  to a label, a constraint function  $\zeta$  that maps each map edge  $me \in AM$  to a constraint, and a labeling function  $\lambda$  that maps each node and edge of  $S_1$  and  $S_2$  to a label. In other words, an application map defines a mapping between two application schemas. Figure 11 shows a segment of a mapping between an application DTD with an *Element-ACCOMMODATION* and the relational schema with *Relation-ACCOMMODATION*.

An application map  $AM$  is an *instanceOf* a map model  $MM$  if the structure  $MS$  that defines  $AM$  is an *instanceOf*  $\mathcal{MP}$  that defines  $MM$ .

**Correctness of Application Map.** An application map  $AM$  is correct if (1) for a given correct input application schema  $AS_1$ , the application map produces a structure  $S$  that defines a correct application schema  $AS_2$ ; (2) the application map  $AM$  is an *instanceOf* a map model  $MM$ ; and (3) all constraints mapped via the constraint function  $\zeta$  are satisfied.

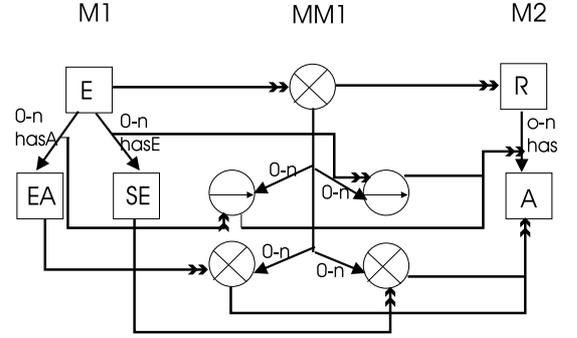


Figure 10: An Example Map Model to Map DTD to Relational Model.

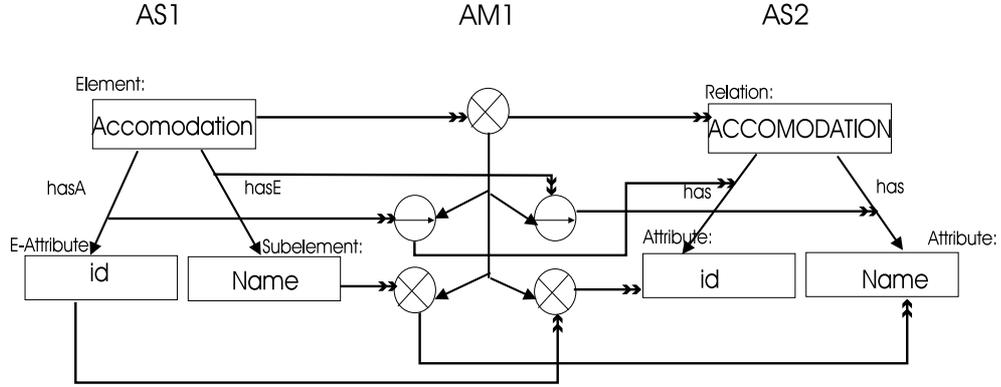


Figure 11: Example Application Map to Transform Application DTD to Relational Application Schema.

## 6 $\mathcal{G}$ angam - Our Model Management System

In the previous sections we have presented a theoretical framework that allows us to model maps between data models as well as between application schemas; and allows us to ensure that the outputs for the maps are correct. However, this is a far cry from the actual realization of the system that must now not only allow the modeling of these maps and ensure their correctness (as per the theory) but must also perform the data transformation and all the while making it easier for users to actually perform these tasks. Towards this goal, we now show how the proposed model management theory can be realized using a relational engine as data manager. Below we show the steps for defining a map between data models and between application schemas to illustrate the mapping theory in concrete terms. However, we would like to point out that once a map model is defined in  $\mathcal{G}$ angam, most of the mapping steps illustrated here can be automated. Hence a user of  $\mathcal{G}$ angam while reaping benefits from its formal basis need not be aware of its details. In this section we describe the data model and map model layers and work our way down to the application data and data map layer (Figure 4). For simplicity, we use relational tables to represent the storage structures to capture models and maps, however an object-relational model is utilized in our implementation.

### 6.1 Setting up Data Models and Map Models

In this section we describe how a system administrator of  $\mathcal{G}$ angam would initialize the system by first entering data models and map models relevant for the domain.

**Modeling Data Models.** Table 1 describes the storage structure for explicitly modeling the data models in  $\mathcal{G}$ angam. The `NODE` and `EDGE` tables represent the nodes and the edges respectively of patterns that define a data model (Section 3). Table 2 shows a subset of the DTD model captured as data in our metamodel tables in the data model layer (Figure 1). Other data models such as the relational model can be represented in a similar fashion.

**Modeling Map Models.** Similarly Table 3 stores a map model. The tables `MAP-NODE` and `MAP-EDGE` store the map nodes and map edges contained in the set of patterns that define a map model. Properties such

NODE		EDGE	
Field	Description	Field	Description
Node Name	Label of node.	From-Node	Id of node that edge stems at.
Node Type	Complex (C) or atomic (A).	Edge-Label	Label of edge.
NodeID	Internal unique identifier.	To-Node	Id of node edge is incident on.
		Edge-Type	Co - containment, P - property.
		Quantifier	Range of edge.
		EdgeID	Internal unique identifier.

Table 1: The Storage Structure for a Data Model in  $\mathcal{G}$ angam.

NODE			EDGE					
Node Name	Node Type	NodeID	From-Node	Edge-Label	To-Node	Edge-Type	Quant.	EdgeID
DTD	C	N-1	DTD	Name	String	P	1-1	E-1
Element	C	N-2	DTD	FullPath	String	P	1-1	E-2
E-Attribute	C	N-3	DTD	HasE	Element	Co	1-n	E-3
SubElement	C	N-4	Element	HasA	Attribute	Co	0-n	E-4
Group	C	N-5	Element	Name	String	P	1-1	E-5
Bool	A	N-6	Element	Grps	Group	Co	0-n	E-6
String	A	N-7	Element	ID	String	P	1-1	E-7
Integer	A	N-8	Element	Subelem	SubElement	Co	0-n	E-9
			E-Attribute	Required	Bool	P	0-1	E-10
			E-Attribute	Name	String	P	1-1	E-11
			E-Attribute	Type	String	P	1-1	E-12
			E-Attribute	Fixed	String	P	0-1	E-13
			Group	Required	Bool	P	0-1	E-14
			Group	CanRepeat	Bool	P	0-1	E-15
			Group	Order	Integer	P	0-1	E-16
			Group	Subelem	SubElement	Co	1-n	E-17
			SubElement	Required	Bool	P	0-1	E-18
			SubElement	CanRepeat	Bool	P	0-1	E-19
			SubElement	Order	Integer	P	0-1	E-20
			SubElement	ElemDef	Element	Co	1-1	E-8

Table 2: The DTD Model Described Using the Structure in Table 1 in the Data Model Layer. We use the Name of the Node for Clarity. Here the Columns To-Node and From-Node are References (Foreign Keys) to the `NODE.NodeID`.

as label, type, quantifier, constraints, etc. are stored along with the map nodes and map edges in these tables. Table 4 captures the map pattern given in Figure 10. This map pattern maps a fragment of the DTD model (Table 2) to a fragment of the relational model.

## 6.2 Application Schema and Application Map Layer

In this section we describe how an application schema and a map between two application schemas are stored in the application layer of  $\mathcal{G}$ angam. Storage structures for the application schemas and the application maps are generated by our system directly from the data model (Table 1) and the map model (Table 3) respectively.

**Application Layer.** Data models described in the data model layer (such as Table 2) are used by the system to generate storage structures in the application layer. The Storage Structure Generator, SSG for short, in

MAP-NODE		MAP-EDGE	
Field	Description	Field	Description
Node Name	Label of map node.	From-Map-Node	Id of map node edge stems from.
Node Type	Ident (I), Cross (C), Merge (M), Project (P), Edge (E).	Map-Edge-Label	Label of map edge.
NodeID	Internal unique identifier.	To-Map-Node	Id of map node edge is incident on.
		Edge-Type	Containment (C), Input (I), Output (O)
		Quantifier	Range for map edge.
		MapConstraint	Constraint on map edge
		EdgeID	Internal unique identifier.

Table 3: The Storage Structure for a Map Model.

MAP-NODE			MAP-EDGE						
Node Name	Node Type	NodeID	From-M-Node	M-Edge-Label	To-M-Node	M-Edge-Type	M-Quant.	MapC	M-EdgeID
Cross-1	C	M-1	Cross-1	elem-in	Element	I	1-1	-	ME-1
Cross-2	C	M-2	Cross-1	rel-out	Relation	O	1-1	-	ME-2
Cross-3	C	M-3	Cross-1	att-con	Cross-2	C	0-n	-	ME-3
Edge-1	E	M-4	Cross-1	subE-con	Cross-3	C	0-n	-	ME-4
Edge-2	E	M-5	Cross-1	hasA-con	Edge-1	E	0-n	-	ME-5
			Cross-1	se-con	Edge-2	E	0-n	-	ME-6
			Cross-2	att-in	E-Attribute	I	1-1	-	ME-7
			Cross-2	att-out	Attribute	O	1-1	-	ME-8
			Cross-3	sube-in	SubElement	I	1-1	-	ME-9
			Cross-3	att-out	Attribute	O	1-1	*	ME-10
			Edge-1	hasA-in	hasA	I	1-1	-	ME-11
			Edge-1	has-out	has	O	1-1	-	ME-12
			Edge-2	sube-in	Subelem	I	1-1	-	ME-13
			Edge-2	has-out	has	O	1-1	-	ME-14

Table 4: A Map Pattern describing the mapping of a DTD pattern to a relational pattern. This map pattern is part of the map model that maps an entire DTD to a relational model. We use the Name of the node for clarity purpose. \* = If  $|\text{SubElement}.\text{ElemDef}.\text{Element}.\text{hasA}| = 0$

Figure 13 is a two-pass algorithm. In the first pass, COLLECT (Figure 12) gathers the set of nodes that belong to a given data model by following the containment edges from a given root. In the second pass, we generate the application schema storage structure. Here we use the intuition that most nodes can have multiple instantiations. Thus given a data model DM, for every row in the NODE table representing a node  $n \in \text{DM}$  we generate a table in the application layer with the table name equal to the node's label. All edges between the nodes are represented as attributes of the table generated for either the From-Node or the To-Node using the Edge-Label as the attribute name. Figure 13 gives the details of the algorithm.

Table 5 is gives a subset of tables generated from the Table 2 using the algorithms given in Figure 12 and 13 in the application layer (Figure 1). Table 6 depicts a segment of the sample application DTD (shown in Figure 2). The data in Table 6 is validated at loading time and thus ensures that it is an instanceOf the DTD model in Table 2. It thereby conforms to the structure of the DTD model and preserves the constraints and quantifiers specified in the DTD model.

```

COLLECT (Node: Root,
         Set<String>: TablesToCreate)
Root: root node (starting point)
TablesToCreate : Set of Node Labels
e: an edge, a row in the EDGE table
{
  if Root ∈ TablesToCreate:
    return;
  for all e in EDGE:
    {
      if (e.From-Node = Root)
        if (e.EdgeType) = Co
          TablesToCreate.add(e.To-Node)
          COLLECT(e.To-Node,
                  TablesToCreate)
    }
}

```

Figure 12: First Pass: Collecting the names of all Types that must be Created in the Application Layer from the NODE and EDGE Tables.

```

SSG (Set<Strings>: TablesToCreate)
TablesToCreate : Node Labels to create tables
e: an edge, a row in the EDGE table
CreateDef(): creates empty types with name
addAttribute(): schema function to add an attribute
                to existing type
getType(): Helper function - returns handle on
                the type given a name.
{
  CreateDef(TablesToCreate)
  for all e in EDGE
    where e.From-Node ∈ TablesToCreate
    {
      if (e.EdgeType = P)
        addAttribute (e.EdgeLabel,
                     getType(e.From-Node), String)
      else if ((e.EdgeType = Co) &&
              (e.Quantifier = 0-n))
        {
          addAttribute (e.EdgeLabel,
                       getType(e.To-Node),
                       REF(getType(From-Node) )
        }
      else if ((e.EdgeType = Co) &&
              (e.Quantifier = 0-1))
        {
          addAttribute (e.EdgeLabel,
                       getType(e.From-Node),
                       REF(getType(To-Node))
        }
    }
}

```

Figure 13: Second Pass: Creating Table Definitions Using Data in NODE and EDGE Tables.

DTD	
Field	Description
Name	Name of the DTD
FullPath	Full path of the DTD
UID	Internal unique identifier

ELEMENT	
Field	Description
Name	Name of the Element
ID	User assigned id for element
UID	Internal unique identifier
hasE	reference to DTD

E-ATTRIBUTE	
Field	Description
Name	Name of the Attribute
Type	Type of the Attribute (String)
Required	Boolean
Fixed	Boolean
hasA	reference to Element.
ID	Internal unique identifier

SUBELEMENT	
Field	Description
Required	Boolean
CanRepeat	Boolean
Order	Integer
ElemDef	reference to Element (element definition)
Subelem	reference to Element (parent element)
UID	Internal unique identifier

Table 5: Storage Structure in Application Layer for Application DTDs.

**Application Map Layer.** Similar to the data model, the map model is used to generate storage structures for the application map. Table 7 depicts the set of tables generated from the map model in Table 4. To generate this structure we use the map storage structure generator called MSSG. MSSG is also a two-pass algorithm. The first pass of the algorithm, similar to the COLLECT algorithm, collects a set of tables that

DTD		
Name	FullPath	ID
Acc	url	D1

E-ATTRIBUTE					
Name	Type	Required	Fixed	hasA	UID
id	CDATA	False	False	E1	D5

ELEMENT			
Name	ID	hasE	UID
accommodation	E1	D1	D2
name	E2	D1	D3
PCDATA	E3	D1	D4

SUBELEMENT					
Required	Can-Repeat	Order	Elem-Def	hasE	UID
True	False	1	E2	E1	D6
True	False	1	E3	E2	D7

Table 6: Application DTD Stored in the Application Layer.

need to be created. The second pass of the algorithm (MSSG) creates the table definitions as represented by the map model based on the algorithm given in Figure 14. Here each map node is represented by a table with the table name the same as the node's label. Different kinds of edges are treated differently. The details of this algorithm are given in Figure 14.

```

MSSG (Set<Strings>: MTablesToCreate)
MTablesToCreate : Map Node Labels to create types
e: an edge, a row in the EDGE table
CreateDef(): creates empty types with name
addAttribute(): schema function to add an attribute to existing type
getType(): returns handle on type given a name.
{
  CreateDef(MTablesToCreate)
  for all m in MTablesToCreate
    if m.NodeType = I or M or P
      createType((m.Name)-SECONDARY)

  for all e in EDGE    where e.From-Node ∈ MTablesToCreate
  {
    if (e.EdgeType = I or O)
      addAttribute(e.EdgeLabel, getType(e.From-Node), REF(getType(e.To-Node)))
    else if ((e.EdgeType = C) && (e.Quantifier = 0-n))
    {
      addAttribute(e.EdgeLabel, getType(e.To-Node), REF(getType(e.From-Node)))
    }
    else if ((e.EdgeType = C) && (e.Quantifier = 0-1))
    {
      addAttribute(e.EdgeLabel, getType(e.From-Node),
        REF(getType(To-Node)))
    }
    else return;
  } } }

```

Figure 14: Second Pass MSSG: Creating Table Definitions Using Data in MAP-NODE and MAP-EDGE Tables.

Given the data in Table 4, the MSSG algorithm generates the set of tables for the application map layer given in Table 7. Data in this storage structure maps an application DTD to produce as output a relational application schema. Each table in this layer has an additional attribute which refers to the code segment

(function) used for the corresponding data transformation as described next. Table 8 represents an application map, depicted in Figure 11, whose input is the application DTD given in Table 6. The application map in Table 8 is an instance of the map model in Table 4 and conforms to the structure and the constraints of the map.

It should be noted here that it is possible to automate the creation of the application map as depicted in Table 8. To automate this process a Sngam user needs to specify the map model as well as the input for the application map. It is a fairly straightforward process to then create the application map using this information.

CROSS-1		CROSS-2	
Field	Description	Field	Description
UID	Internal unique identifier	UID	Internal unique identifier
elem-in	Input for the Cross Node	att-in	E-attribute input
rel-out	Output of the Cross Node	att-out	Attribute output
func	Data transformation function	att-con	Containment, Reference to Cross-1
		func	Data transformation function

CROSS-3		EDGE-1	
Field	Description	Field	Description
UID	Unique Identifier	UID	Internal unique identifier
subelem-in	Subelement input	hasA-in	edge input
att-out	Attribute output	has	edge output
sube-con	Containment - reference to Cross-1	hasA-con	reference to Cross-1
func	Data transformation function	func	Data transformation function

EDGE-2	
Field	Description
UID	Unique Identifier
subelem-in	Subelement edge input
has-out	Attribute output
se-con	Containment - reference to Cross-1
func	Data transformation function

Table 7: Storage Structure for an Application Map Conforming to the Map Model in Table 4.

CROSS-1				CROSS-2				
UID	elem-in	rel-out	func	UID	eatt-in	att-out	att-con	func
am1	accommodation	ACCOMMODATION	elem-rel	am2	id	id	am1	ea-a

CROSS-3					EDGE-1				
UID	subelem-in	att-out	sube-con	func	UID	hasA-in	has-out	hasA-con	func
am3	name	name	am1	se-a	am5	hasA-id	ACC-has	am1	edgeT
am4	PCDATA	name	am3	tostring					

EDGE-2				
UID	subelem-in	has-out	sube-con	func
am6	se-name	ACC-has	am1	edgeT

Table 8: Application Map to Map the Application DTD in Table 6 to a Relational Structure.

### 6.3 Application Data and Map Data Transformers

Once the map between application schemas is in place, the next, perhaps the most crucial, step is to transform the input application data in the data layer to data in the desired output format. In *Sangam* we automate this process by encoding atomic units of transformation code as a part of each map node metaconstruct. For some map nodes there may be multiple default functions. A user or the system can select from these different choices at map generation time in the application layer. *Data transformers* in the data layer are generated using these atomic units of code and are composed as per the application map. While we are able to generate these transformers in a generic fashion, we rely on these pre-coded transformation functions to for instance extract data from XML documents or to load it into relational tables. The default functions may be overridden by a user-defined function such as a function that allows the user to merge two attributes, for example `First-Name` and `Last-Name` into one attribute `Full-Name` by concatenating the string values for the two attributes.

## 7 Related Work

**Meta Data Modeling.** Meta-modeling has been utilized to model many different types of information, data, roles and business rules [AT96, BR00, GL98]. In particular, it has been looked at as a middle-ware medium to handle schema integration and data transformation over the last twenty years [MR83, AT96, BR00, PR95]. Papazoglou et al. [PR95] propose a middle-layer meta-model to accomplish transformations between the OO and relational data models. The transformations are accomplished by a set of pre-defined translation rules that can convert the OO or relational data models to and from the middle-layer meta-model. Using translations as basic building blocks, they aim to automatically generate mappings from one given model to another at run-time. Atzeni et al. [AT96] have presented a framework to describe data models and application schemas. They focus on discovering translations between data models and hence application schemas. We make use of their graph model to express the data models and application schemas in our system. Commercially, Microsoft Repository [Ber99] and Rochade Information Model [Roc00] are meta-repositories that generically describe mainly data models for system integration purposes. To the best of our knowledge, these works all place the primary focus on the modeling of the data models and not of transformations at the meta level. In our work we now focus on modeling of maps between data models and between application schemas and the subsequent automatic generation of code to perform the data transformation.

**Schema Integration and Data Transformation.** The crux of MMS, the mapping, has been dealt with in the literature under the umbrella of schema transformation and integration [HMN<sup>+</sup>99, MZ98, FK99, MIR93, CJR98]. However, this work is typically specific to either the application domain or to a particular data model and does not deal with meta-modeling [MR83, AT96, BR00, PR95]. Recent work related to ours are Clio [HMN<sup>+</sup>99] a research project at IBM's Almaden Research Center and work by Milo and Zohar [MZ98]. Clio, a tool for creating mappings between two data representations semi-automatically (i.e., with user input) focuses on supporting querying of data in either the source or the target representation and on just in time cleansing and transformation of data. Milo et al. [MZ98] have looked at the problem of data translations based on schema-matching. They follow an approach similar to Atzeni et al. [AT96] and Papazoglou et

al. [PR95], but not at the meta-level, in that they define a set of translation rules to enable discovery of relationships between two application schemas. We can directly make use of translation algorithms from the literature, such as the algorithms for translating between an XML-DTD and relational schema [FK99] or mapping rules [MZ98], and model them as mappings in the MMS. However, our focus is not discovering such algorithms for mapping but rather on the generic expressibility of any possible (future) mapping and its management. Work on equivalence of the translations between models [MIR93] is of particular importance as such properties of maps can also be established.

## 8 Conclusions

In this paper we have presented a framework to allow modeling of mapping between different data models as well as to allow re-structuring within a data model. We are currently in the process of implementing *Gangam* in Java using Oracle 8i (and its object extensions) as the MMS data store. A version of this system will be demonstrated at SIGMOD 2001 [CRZS01].

In conclusion, we would like to point out that while a model management approach may not be the quickest solution to mapping between data models, it still offers some immediate advantages in terms of flexibility and extensibility. Beyond its immediate advantages, we believe our approach has a more far reaching impact, for as today's data models become legacy models our framework does not become obsolete. But rather it is extensible in that we can, not only add other data models but we can also easily describe or re-target existing maps between such new data models using the map metaconstructs, thus getting actual transformation code for free with very little effort.

## References

- [AT96] P. Atzeni and R. Torlone. Management of Multiple Models in an Extensible Database Design Tool. In Peter M. G. Apers and et al., editors, *Advances in Database Technology - EDBT'96, Avignon, France, March 25-29*, LNCS. Springer, 1996.
- [Ber99] Bernstein, P. A. and Bergstraesser, T. et al. Microsoft Repository Version 2 and the Open Information Model. *Information Systems*, 2(24):71–98, 1999.
- [Boo94] G. Booch. *Object-Oriented Analysis and Design*. Benjamin Cummings Pub., 1994.
- [BR00] P. A. Bernstein and E. Rahm. Data Warehouse Scenarios for Model Management. In *International Conference on Conceptual Modeling*, 2000.
- [CJR98] K.T. Claypool, J. Jin, and E.A. Rundensteiner. SERF: Schema Evolution through an Extensible, Re-usable and Flexible Framework. In *Int. Conf. on Information and Knowledge Management*, pages 314–321, November 1998.
- [CRZS01] K.T. Claypool, E.A. Rundensteiner, X. Zhang, and H. et al. Su. Model Management - A Solution to Support Multiple Data Models, Their Mappings and Maintenance. In *Demo Session Proceedings of SIGMOD'01, to appear*, 2001.
- [FK99] D. Florescu and D. Kossmann. Storing and Querying XML Data Using an RDBMS. In *Bulletin of the Technical Committee on Data Engineering*, pages 27–34, Sept. 1999.
- [GL98] S. Gbel and K. Lutze. Development of Meta Databases for Geospatial Data in the WWW. In *ACM-GIS*, pages 94–99, 1998.

- [HMN<sup>+</sup>99] L.M. Haas, R.J. Miller, B. Niswonger, M.T. Roth, P. Schwarz, and E.L. Wimmers. Transforming Heterogeneous Data with Database Middleware: Beyond Integration. *IEEE Data Engineering Bulletin*, 22(1):31–36, 1999.
- [MIR93] R. J. Miller, Y. E. Ioannidis, and R. Ramakrishnan. The Use of Information Capacity in Schema Integration and Translation. In *Int. Conference on Very Large Data Bases*, pages 120–133, 1993.
- [MR83] L. Mark and N. Roussopoulos. Integration of Data, Schema and Meta-Schema in the Context of Self-Documenting Data Models. In Carl G. Davis, Sushil Jajodia, Peter A. Ng, and Raymond T. Yeh, editors, *Proceedings of the 3rd Int. Conf. on Entity-Relationship Approach (ER'83)*, pages 585–602. North Holland, 1983.
- [MZ98] T. Milo and S. Zohar. Using Schema Matching to Simplify Heterogeneous Data Translation. In Ashish Gupta, Oded Shmueli, and Jennifer Widom, editors, *VLDB'98, Proceedings of 24th International Conference on Very Large Data Bases, August 24-27, 1998, New York City, New York, USA*, pages 122–133. Morgan Kaufmann, 1998.
- [PR95] M.P. Papazoglou and N. Russell. A Semantic Meta-Modeling Approach to Schema Transformation. In *CIKM '95*, pages 113–121. ACM, 1995.
- [Roc00] Rochade Information Model. Rochade Information Model. <http://support.viasoft.com/html/2RochDocDownload.html>, 2000.