

# Developing Metadata-Intensive Applications with Rondo

Sergey Melnik, Erhard Rahm  
University of Leipzig, Germany  
melnik@db.stanford.edu, rahm@informatik.uni-leipzig.de

Philip A. Bernstein  
Microsoft Research, Redmond  
philbe@microsoft.com

## ABSTRACT

The future of the Semantic Web depends on whether or not we succeed to integrate reliably thousands of online applications, services, and databases. These systems are tied together using mediators, mappings, database views, and transformation scripts. Model management aims at reducing the amount of programming needed for the development of such integrated applications. We present a first complete prototype of a generic model-management system, in which high-level operators are used to manipulate models and mappings between models. We define the key operators and conceptual structures and describe their use and implementation. We examine the solutions for three model-management tasks: change propagation, view reuse, and reintegration.

**Keywords:** Generic Model Management

## 1. INTRODUCTION

The future of the Semantic Web depends on whether or not we succeed to integrate reliably thousands of online applications, services, and databases. These systems are tied together using mediators, mappings, database views, and transformation scripts of various kinds, whose development is extremely costly. Once up and running, the maintenance of the plumbing used to connect heterogeneous systems becomes a second crucial issue. Even minor changes in database schemas or interfaces of online services break components that rely on them. Thus, the development of applications that help bridge disparate systems lies on the critical path that leads to the Semantic Web.

Such applications deal with the tasks that arise in the context of database design, data integration, data translation, model-driven website management, data warehousing, etc. They manipulate a variety of metadata artifacts that are called *models*, such as relational and XML schemas, interface definitions, mediator specifications, or website layouts, and *mappings* between models, such as SQL views or XSLT transformations. Reducing the amount of programming required for the development of such metadata-intensive tasks is the subject of model management research. In fact, many of today's metadata-intensive tasks are still solved manually, because an automated approach requires too much implementation effort due to the lack of a common programming platform.

Database and software engineering researchers have been studying the individual aspects of model management in depth for decades. However, factoring out the common aspects of model management has only recently come to the fore of active research [7]. A major goal of this recent research has been to develop a set of algebraic operators, such as Compose, Match and Merge, that generalize the transformation operations utilized across various metadata applications. These operators are applied to models and mappings as a

whole, rather than to their individual elements, and simplify the programming of metadata applications by raising the level of abstraction. Moreover, the operators are *generic* in the sense that they can be utilized for different kinds of models and scenarios. Although many model-management tasks can be automated, there remain critical places where human decision-making is needed, e.g., to address semantic heterogeneity. Thus, some of the operations are inherently semiautomatic and require feedback of a human engineer before, during, or after operator execution.

Our goal is to investigate whether metadata management can be done in a generic fashion, the key question raised in [7]. Detailed walkthroughs of various model-management problems have been examined to address this question (e.g., in [5,9]). Our contribution is that we succeeded in making such abstract programs executable. In this article, we present a prototype of a programming platform for model management and describe the conceptual structures and operators that we developed (a short version of this article appeared in [18]). Primarily, our prototype supports the developers of model-management solutions, by providing a high-level programming environment. However, it also addresses the needs of the engineers who deploy these solutions by offering a graphical user interface (GUI) to receive their feedback in semiautomatic operations.

In designing and implementing our prototype, we consciously focus on simplicity. We investigate how far we can go with a comparatively weak representation of models and mappings that can be used to solve an interesting class of problems. We also determine how much code is needed for the most basic, but still useful, model management system.

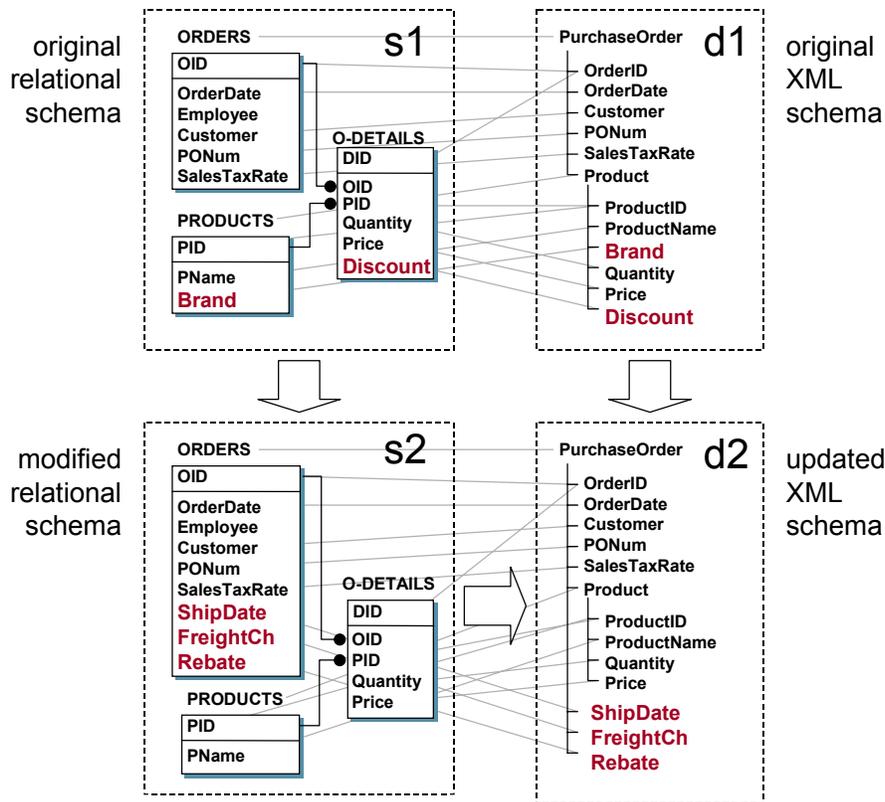
The key contributions of this article are as follows:

- We introduce conceptual structures used for representing models and mappings. We explore a simple class of mappings between models that we call morphisms and suggest a new structure called selector.
- We define the semantics of the key model-management operators on the conceptual structures that we introduce, and suggest several new generic operators.
- We present new algorithms used for implementing the operators Extract and Merge.
- We examine the solutions for three important model-management tasks that involve manipulations of relational schemas, XML schemas, and SQL views.
- Finally, we describe the first complete prototype implementation of model management and demonstrate how it can be extended to embrace new kinds of models.

This article is organized as follows. In Section 2 we walk through a model-management scenario to motivate the conceptual structures and operator definitions that we present in Sections 3 and 4. Section 5 is devoted to the implementation and the algorithms that we developed. Section 6 describes our prototype in more detail. In Sections 7 and 8, we discuss two further model-management scenarios, view reuse and reintegration. Related work is reviewed in Section 9. We outline some preliminary ideas on structural vs. state-based semantics of operators and scripts in Section 10, and conclude in Section 11.

## 2. MOTIVATING SCENARIO

To motivate the operator definitions that we give in this article, we will use a scenario that is illustrated in Figure 1 and exemplifies one of the patterns that can be found in many metadata-intensive applications. Consider an e-commerce company that needs to supply its purchase order data to a business partner. The data is stored in a relational database according to a relational schema  $s1$ . For the purpose of data exchange, both companies agree to use a common XML schema  $d1$ . (The correspondences between the



**Figure 1: Scenario illustrating propagation of changes from a relational to an XML schema**

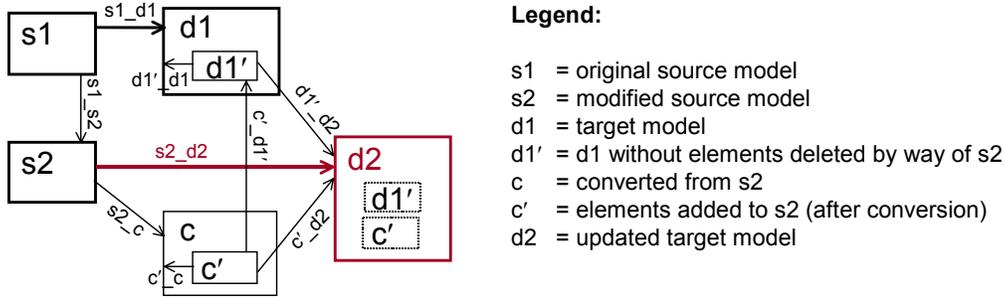
elements of schema  $s1$  and  $d1$  are depicted as light gray lines.) Schema  $d1$  differs from  $s1$  in terms of structure and naming.

The relational schema undergoes periodic changes due to the dynamic nature of the business. Assume that  $s2$  is a new version of  $s1$ , in which columns “Brand” and “Discount” have been deleted, and columns “ShipDate”, “FreightCh” (freight charge), and “Rebate” have been added. These changes (highlighted in bold in Figure 1) need to be propagated to the XML schema, so that  $d1$  becomes  $d2$ .

The change propagation described above can be done as follows. First, the changes introduced by  $s2$  need to be detected, i.e.,  $s1$  and  $s2$  need to be matched. Then, the  $d1$  images of the elements deleted in  $s1$  need to be removed from  $d1$ . Finally, the XML schema counterparts of the added and renamed columns in  $s1$  need to be merged into  $d1$  to obtain  $d2$ . During these steps, intervention of a human engineer may be required, for example, to decide whether the new column “Rebate” should indeed be added to the exchange schema or is not part of the exchanged data and should be omitted. Still, a major portion of the work is mechanical and can be automated.

Notice that the procedure sketched above could be applied in the reverse case, when the XML schema  $d1$  is the one that has been modified and the changes are to be propagated back to the relational schema  $s1$ . Another instance of the same pattern is round-tripping the modifications from a relational schema like  $s1$  to an existing conceptual schema of the data, which may be expressed as an ER diagram. A key idea of generic model management is to solve such tasks at a high level of abstraction using a concise generic script.

Below we present an actual model-management script that implements the above solution for our change propagation scenario, and is directly executable by our prototype.



**Figure 2: Schematic representation of a solution for change propagation scenario of Figure 1**

We will use the script to introduce the major model-management operators, which we define in the subsequent sections. To explain the individual steps of the script, we use a schematic representation of the solution shown in Figure 2. The rectangles labeled  $s1$ ,  $s2$ ,  $d1$ , and  $d2$  represent the four schemas of Figure 1. The arcs between the rectangles denote the *mappings* between the schemas. For example, the correspondences between schemas  $s1$  and  $d1$  in Figure 1 are shown as a single arc from rectangle  $s1$  to  $d1$  in Figure 2.

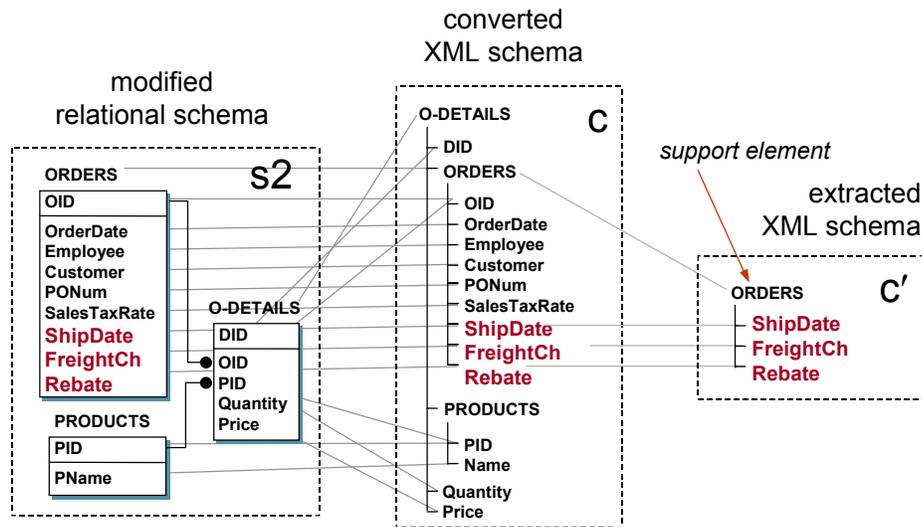
At the bottom of Figure 2, there is a schema  $c$ , which does not appear in Figure 1. To see why it is needed, recall that  $s1$  and  $d1$  are expressed using two different schema languages. The new schema elements added to  $s1$  by way of  $s2$  have no counterparts in schema  $d1$ . That is, the new elements need to be converted from the source schema language to the target language. For example, the attribute “ShipDate” added to relation “ORDERS” needs to be converted to a subelement of the complex type “PurchaseOrder” in the XML schema. This step is often referred to as schema translation in the literature. In our solution, we assume that such a translation tool is available as an operator, say SQL2XSD, which takes as input a relational schema and produces as output an XML schema and a mapping between the original and converted schema elements. Thus, the schema  $c$  and the mapping  $s2\_c$  between  $s2$  and  $c$  shown in Figure 2 are obtained as  $\langle c, s2\_c \rangle = \text{SQL2XSD}(s2)$ . Schema  $c$  is illustrated in Figure 3. Note that  $c$  is not yet the desired result  $d2$ ; for example,  $c$  contains an unneeded complex type O-DETAILS, and differs from  $d2$  structurally.

Now, our solution for the change propagation scenario can be expressed as the following script:

```
operator PropagateChanges(s1, d1, s1_d1, s2, c, s2_c)
1. s1_s2 = Match(s1, s2);
2.  $\langle d1', d1\_d1 \rangle = \text{Delete}(d1, \text{Traverse}(\text{All}(s1) - \text{Domain}(s1\_s2)), s1\_d1)$ ;
3.  $\langle c', c\_c \rangle = \text{Extract}(c, \text{Traverse}(\text{All}(s2) - \text{Range}(s1\_s2)), s2\_c)$ ;
4.  $c\_d1' = c\_c * \text{Invert}(s2\_c) * \text{Invert}(s1\_s2) * s1\_d1 * \text{Invert}(d1\_d1)$ ;
5.  $\langle d2, c\_d2, d1\_d2 \rangle = \text{Merge}(c', d1', c\_d1')$ ;
6.  $s2\_d2 = s2\_c * \text{Invert}(c\_c) * c\_d2 +$ 
    $\text{Invert}(s1\_s2) * s1\_d1 * \text{Invert}(d1\_d1) * d1\_d2$ ;
7. return  $\langle d2, s2\_d2 \rangle$ ;
```

The script defines a generic operator PropagateChanges, which takes six parameters as input (including the converted schema  $c$ ), and produces two return values  $\langle d2, s2\_d2 \rangle$  as output. Below, we explain the script line by line.

1. In line 1, schemas  $s1$  and  $s2$  are “matched” to detect the changes. The result is a mapping  $s1\_s2$  shown schematically in Figure 2. Speaking informally, the mapping



**Figure 3: Converted schema  $c$  and support element  $ORDERS$  in  $c'$**

connects the equivalent elements of  $s1$  and  $s2$ . The new elements of  $s2$  (e.g., “ShipDate”) and deleted elements of  $s1$  (e.g., “Brand”) have no matching counterparts, so they remain unconnected.

2. Line 2 illustrates how operators can be combined. First, the deleted elements of  $s1$  are identified using the expression  $All(s1) - Domain(s1\_s2)$ , i.e., all elements of  $s1$  without the matched (and thus not deleted) elements. Then, these elements are used to “traverse” the mapping  $s1\_d1$ . For example, the deleted relational attribute “Brand” traverses  $s1\_d1$  and yields the XML schema element “Brand” of  $d1$ . Finally, these  $d1$  images of the deleted elements are removed from  $d1$  using the operator Delete. The result is a new schema  $d1'$  (a “subschema” of  $d1$ ), and a mapping  $d1'\_d1$ , which describes how  $d1'$  relates to  $d1$ .
3. Line 3 is quite similar to line 2. The new elements of  $s2$ , i.e., those missing from the range of  $s1\_s2$ , traverse  $s2\_c$  into the converted model  $c$  (see Figure 3). For example, the image of relational attribute “ShipDate” is an XML schema element “ShipDate” obtained by conversion. A “subschema”  $c'$  containing the images of the new elements is then extracted from  $c$  using the operator Extract, which also returns the mapping  $c'\_c$ . In addition to the elements obtained by traversal like “ShipDate”,  $c'$  contains an extra element of  $c$ , the complex type “ORDERS” that encloses “ShipDate”. Such extra elements are called “support” elements [5]. Support elements may have to be extracted to make  $c'$  a well-formed XML schema.
4. At this point,  $d1'$  is a subschema of  $d1$  without the deleted elements, and  $c'$  contains the added elements and their support elements. Schemas  $d1'$  and  $c'$  need to be merged to obtain the final result  $d2$  (line 5). As we explain in Section 4.5, the merging of two schemas is driven by a mapping that tells how elements of the two schemas, specifically the support elements of  $c'$ , correspond to each other. The mapping between  $d1'$  and  $c'$  is shown in Figure 2 as an arc connecting the two enclosed rectangles. This mapping can be obtained by “composing” the existing mappings between  $c'$ ,  $c$ ,  $s1$ ,  $s2$ ,  $d1$ , and  $d1'$  as  $c'\_c * Invert(s2\_c) * Invert(s1\_s2) * s1\_d1 * Invert(d1'\_d1)$ . To get the composition right, mappings  $s2\_c$ ,  $s1\_s2$ , and  $d1'\_d1$  need to be “inverted”, i.e., the domains and ranges of the mappings need to be swapped.

Thus, we determine by composition that the support element “ORDERS” in  $c'$  corresponds to the element “PurchaseOrder” in  $d1'$ .

5. The final result of change propagation, schema  $d2$ , is computed by the Merge operator. Among other things, the operator Merge creates a single complex type definition from complex type “ORDERS” from  $c'$  and “PurchaseOrder” from  $d1'$ . Additionally, the operator returns two mappings,  $c'_d2$  and  $d1'_d2$ , which describe how  $d2$  relates to the inputs to Merge,  $c'$  and  $d1'$ .
6. As a last step, we compute  $s2_d2$ , a new version of the mapping  $s1_d1$  given as part of the input. We need  $s2_d2$  to ensure that our change propagation script can be re-applied if the source schema evolves again. Since  $d2$  is obtained by merging  $d1'$  and  $c'$ , the mapping  $s2_d2$  is essentially a union of two mappings, the one between  $s2$  and the  $d1'$ -portion of  $d2$ , and the one between the  $s2$  and  $c'$ -portion of  $d2$ . These two mappings can be obtained by composition as  $s2_c * \text{Invert}(c'_c) * c'_d2$  and  $\text{Invert}(s1_s2) * s1_d1 * \text{Invert}(d1'_d1) * d1'_d2$ , respectively. Their union is denoted using the plus sign (+). To illustrate, the first mapping establishes the correspondences between the added elements “ShipDate”, “FreightCh”, “Rebate” in  $s2$  and their  $d2$  counterparts. The second mapping in the union tells us that “OID” in  $s2$  corresponds to “OrderID” in  $d2$ , etc.

Notice that the above script is not limited to propagating changes from relational schemas to XML schemas. In fact, the reverse propagation problem can be solved using the same script by assigning the original and modified XML schemas to  $s1$  and  $s2$ , and the relational schema to  $d1$ . Of course, the input parameters  $c$  and  $s2_c$  need to be obtained using a different converter, e.g., as  $\langle c, s2_c \rangle = \text{XSD2SQL}(s2)$ .

In our implementation, every intermediate result of a script such as the one above can be examined and adjusted by a human engineer using a graphical tool. Specifically, the result of Match in line 1 can be post-processed to remove incorrectly suggested matches and add missing ones. Similarly, the merging in line 5 is in general a semiautomatic process, which requires human feedback. Finally, by adjusting the intermediate results of operator compositions in lines 2 and 3 the engineer can decide which additions and deletions should not be propagated.

In the above discussion, we introduced several operators informally. To make these operators effective and usable by developers, their semantics needs to be specified precisely. Our goal is to make the semantics as “generic” as possible, so the operators can serve a broad range of model-management tasks. In the next two sections we describe this semantics, first by defining the structures on which they operate, and then by describing the operators themselves.

### 3. CONCEPTUAL STRUCTURES

Model-management applications deal with a wide range of metadata artifacts, which include not only schemas, such as the relational and XML schemas in our motivating scenario, but also view definitions, interface specifications, etc. We represent the formal descriptions, or *models*, of these artifacts as directed labeled graphs. This graph representation is quite flexible and can accommodate virtually any type of models.

We also introduce two additional structures, called *morphisms* and *selectors*. Morphisms are binary relationships that establish n:m correspondences between the elements of two models (i.e., nodes of two graphs). For example, in our motivating scenario morphisms are used for keeping track of the XML counterparts of the relational schema elements. Two morphisms, one between  $s1$  and  $d1$  and another between  $s2$  and  $d2$ , are shown in Figure 1 using light gray lines. The third conceptual structure, selector, is a set of elements used in models. A major benefit of using selectors is that various

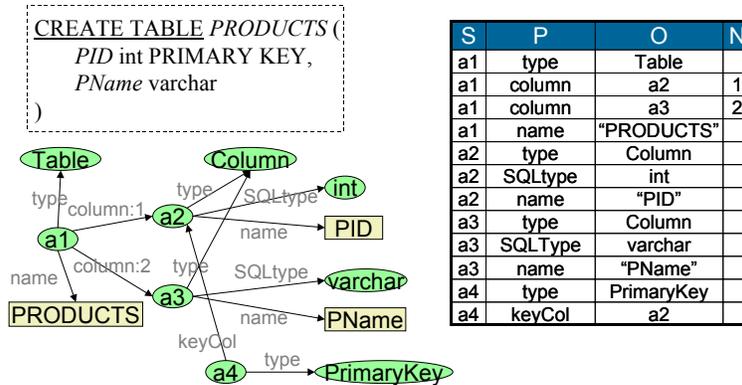


Figure 4: Sample model shown as graph and 4-tuples

operations, in particular the set operations, which would typically produce non-well-formed models if used directly, can be applied to selectors safely.

In the following subsections, we define models, morphisms, and selectors as abstract graph and set structures. We also describe them in an equivalent representation as relations. The latter will make it easier to define the semantics of the operators, which follow later.

### 3.1 Models

We represent models as directed labeled graphs. The nodes of such graphs denote *model elements*, such as relations and attributes in relational schemas, type definitions in XML schemas, clauses of SQL statements, etc. We assume that each element is uniquely identified by an object identifier (OID). A directed labeled graph is a set of edges  $\langle s, p, o \rangle$  where  $s$  is the source node,  $p$  is the edge label, and  $o$  is the target node<sup>1</sup>. For a given source  $s$  and label  $p$ , the target nodes may be sequentially ordered. Their order can be captured by an ordinal property on edges. Thus, conceptually a graph can be viewed as a relation  $M$  with four attributes,  $M(S: \text{OID}, P: \text{OID}, O: \text{OID} \cup \text{Literal}, N: \text{integer})$ , where  $N$  is an optional attribute used for ordering and  $S, P, O$  form a unique key. The node identifiers and edge labels are drawn from the set of OIDs, which can be implemented as integers, pointers, URIs, etc. The literals include strings, integers, floats, and other data types. The type of attribute  $O$  is defined as a union type of OIDs and literals.

Consider the example in Figure 4. It illustrates how a relational table PRODUCTS defined in SQL DDL (top left) is represented as a graph (bottom left) and as a corresponding set of 4-tuples (on the right). The ovals in the graph denote OIDs, and rectangles denote literals. Nodes  $a1$ ,  $a2$ ,  $a3$  represent the table PRODUCTS and its columns PID and PName, respectively. Node  $a4$  represents the primary key constraint on PID. For readability, the identifiers such as Table or Column are spelled out as names rather than opaque IDs.

The order of the columns identified by the nodes  $a2$  and  $a3$  is determined by the values 1 and 2 of attribute  $N$  (fourth attribute of the table with 4-tuples). In general, the node ordering with respect to a given  $\{\text{src node}\}$  and  $\{\text{edge label}\}$  is determined by the SQL query: `SELECT M.O FROM M WHERE M.S={src node} AND M.P={edge label} ORDER BY M.N`. In the example, we have  $M.S=a1$  AND  $M.P=column$ .

A formal specification of the rules for encoding a model as a graph is called a *meta-model*. A model is *well-formed* if it conforms to its meta-model. For example, Figure 4 illustrates a graph encoding of relational schemas that uses specific edge labels, such as

<sup>1</sup> The notation  $(s, p, o)$  stands for (subject, predicate, object).

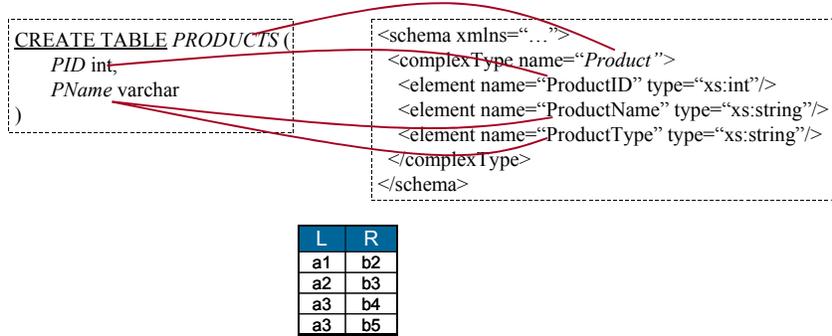


Figure 5: Morphism between relational and XML schema

SQLtype or name, and auxiliary nodes, such as Table, varchar, or PrimaryKey. If we know the relational meta-model, we can tell whether or not a given graph represents a well-formed relational schema. For example, if we know that each column must have an SQL type, then removing the edge ⟨a2, SQLtype, int⟩ from the graph in Figure 4 yields a model that is not well-formed. For the purposes of this article, it is unimportant how a meta-model is represented and how one checks that a model conforms to its meta-model. The details of the graph representation of models remain opaque to the developer of model management applications. Of course, the representation is visible to developers of model management operators. So, a developer must be aware of the representation to implement a custom, non-generic operator, e.g., an operator to normalize relational schemas.

### 3.2 Morphisms

Many metadata-intensive applications, such as data integration and warehousing tools, use a graphical metaphor like the one shown in Figure 1 for representing schema mappings. These mappings are shown to the engineer as sets of lines connecting the elements of two schemas. We call such mappings (*schema*) *morphisms*. Thus, a morphism is a binary relation over two (possibly overlapping) sets of OIDs, i.e., a set of pairs ⟨l, r⟩ drawn from  $OID \times OID$ .

Clearly, a morphism is a weaker representation of a transformation between two models than an SQL view or the mapping languages and expressions suggested in [3,5,14,19,20]. In particular, a morphism carries no semantics about the transformation of instances that conform to the models (e.g., no SQL WHERE-clause). Still, we have found that many mappings can be expressed in this way such as in our change propagation scenario of Section 2. The morphisms have several other advantages. Given our graph representation of models, a morphism can represent a mapping between different kinds of models, e.g., between a relational and XML schema. A morphism can always be inverted

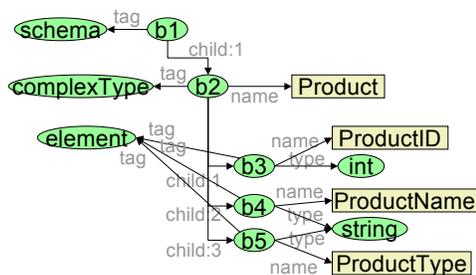


Figure 6: Graph representation of XML schema in Figure 5

V
a1
a2
a3
a4
Table
Column
PrimaryKey
int
varchar

Figure 7: Example of a selector

and composed. (In contrast, an SQL view cannot be composed with an XSLT transformation in an obvious way.) And since morphisms can be expressed as binary relations, they can be implemented and manipulated easily.

Consider the example in Figure 5. The top part of the figure shows the relational schema of Figure 4 and an XML schema. A morphism between the two schemas is depicted graphically as four arcs that connect the elements of the schemas. The bottom part of the figure shows the same morphism represented as a relation. The node identifiers  $a_1$ ,  $a_2$ ,  $a_3$  correspond to those of Figure 4. The nodes  $b_2$ ,  $b_3$ ,  $b_4$ ,  $b_5$  denote respectively the complex type “Product” and the elements “ProductID”, “ProductName”, and “ProductType” defined in the XML schema (its graph representation is illustrated in Figure 6). Notice that a node can be connected to multiple nodes; e.g.  $a_3$  (“PName”) is connected to  $b_4$  (“ProductName”) and  $b_5$  (“ProductType”). Moreover, various kinds of model elements, such as relations or attributes, can participate in a morphism.

In an implementation, it may be convenient to annotate the pairs  $\langle l, r \rangle$  with additional properties. For example, most implementations of the Match operator compute similarity values between the elements of two models. These values can be returned conveniently using a morphism in which each pair has an additional similarity property. Hence, although we define a morphism conceptually as a binary relation  $H(L: \text{OID}, R: \text{OID})$ , it may contain additional attributes, as required by the individual operators. Typically, the  $L$  elements originate from one model, and the  $R$  elements from another.

### 3.3 Selectors

A selector is a set of node identifiers, which may originate from a single or multiple models. It can be represented as a relation with a single attribute,  $S(V: \text{OID})$ , where  $V$  is a unique key. Figure 7 shows an example of a selector that contains all OIDs used in the model depicted in Figure 4.

## 4. OPERATORS

In our motivating scenario, we introduced high-level operators whose inputs and outputs are models, morphisms and selectors, such as Match, Delete, Traverse, Extract, and Invert. Such operators raise the level of abstraction of manipulating metadata structures by considering whole models and morphisms at a time, as opposed to node-at-a-time primitives. In this section, we define the precise semantics of these operators on the structures defined in Section 3. Their implementation is covered in Section 5.

We start our presentation of operator semantics in Section 4.1 with what we call *primitive* operators. These are generic operators whose semantics can be defined formally using the relational algebraic manipulation of the relational representations of Section 3. For notational convenience, we express this manipulation in SQL. After that, we introduce the other more powerful operators: such as Extract, Delete, Match, and Merge, whose semantics is more subtle and still a subject of ongoing research.

As we will see, some operators, such as Subgraph or Copy, are agnostic about the kind of models passed as input, whereas the semantics of others depends on the underlying meta-model. The GUI operators EditMap and EditSelector allow arbitrary transformations of morphisms and selectors by an engineer. Thus, their semantics cannot be constrained any further.

### 4.1 Primitive operators

Table 1 lists the definitions of seven primitive operators. The left column contains the operator definitions expressed in SQL. Variables  $m$ ,  $s$ , and  $map$  hold a model, a selector, and a morphism, respectively. The right column illustrates the application of the operators using simple examples. All primitive operators defined in the table are standard

set-theoretic operators. Notice that their definitions are expressed declaratively, i.e., the implementation of these operators, or functional combinations thereof, can be optimized using standard query optimization techniques.

The operator Domain extracts the “left” elements from a morphism and returns a selector that holds the result. The operator RestrictDomain restricts a morphism to a smaller element domain, which is specified by the selector passed as a second parameter of the operator. The Invert operator swaps the left and right elements of a morphism. The Compose (\*) operator is defined as the natural join of two morphisms, yielding another morphism. The TransitiveClosure operator on morphisms is specified using a recursive SQL definition. The Id operator creates an identity morphism over a given selector.

The operator Subgraph( $m, s$ ) extracts from model  $m$  a subgraph induced by the nodes referenced in  $s$ . The literals attached to the nodes in  $s$  are also extracted from  $m$ . In the example of Table 1, the literal “PID” is not contained in the input selector  $s$ , but the edge  $\langle a2, name, \text{“PID”} \rangle$  is nevertheless returned as part of the result. The extracted subgraph may not be a well-formed model. That is, it may not be fully connected and may not conform to its meta-model.

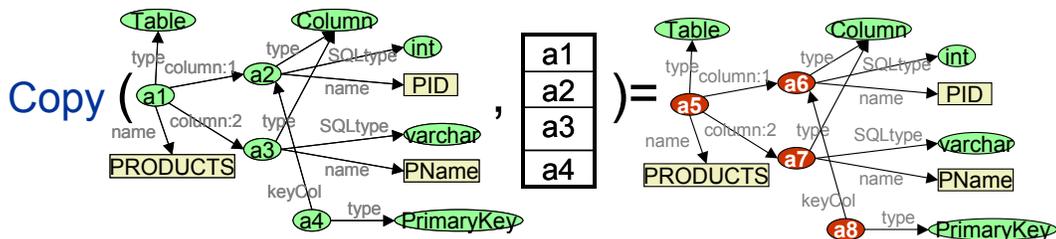
**Table 1: Definitions of primitive operators**

Definition	Example
<b>Domain</b> (map) := SELECT DISTINCT map.L AS V FROM map	$\text{Domain}\left(\begin{bmatrix} a1 & b1 \\ a2 & b2 \end{bmatrix}\right) = \begin{bmatrix} a1 \\ a2 \end{bmatrix}$
<b>RestrictDomain</b> (map, s) := SELECT * FROM map WHERE map.L IN s	$\text{RestrictDomain}\left(\begin{bmatrix} a1 & b1 \\ a2 & b2 \end{bmatrix}, [a1]\right) = \begin{bmatrix} a1 & b1 \end{bmatrix}$
<b>Invert</b> (map) := SELECT map.R AS L, map.L AS R FROM map	$\text{Invert}\left(\begin{bmatrix} a1 & b1 \\ a2 & b2 \end{bmatrix}\right) = \begin{bmatrix} b1 & a1 \\ b2 & a2 \end{bmatrix}$
<b>Compose</b> (map1, map2) := SELECT DISTINCT map1.L, map2.R FROM map1, map2 WHERE map1.R = map2.L	$\text{Compose}\left(\begin{bmatrix} a1 & b1 \\ a2 & b2 \end{bmatrix}, \begin{bmatrix} b1 & c1 \end{bmatrix}\right) = \begin{bmatrix} a1 & c1 \end{bmatrix}$
<b>TransitiveClosure</b> (map) := WITH RECURSIVE TC(L, R) AS (map UNION SELECT DISTINCT TC.L, map.R FROM TC, map WHERE TC.R = map.L) SELECT * FROM TC	$\text{TransitiveClosure}\left(\begin{bmatrix} a & b \\ b & c \end{bmatrix}\right) = \begin{bmatrix} a & b \\ b & c \\ a & c \end{bmatrix}$
<b>Id</b> (s) := SELECT s.V AS L, s.V AS R FROM s	$\text{Id}\left(\begin{bmatrix} a1 \\ a2 \end{bmatrix}\right) = \begin{bmatrix} a1 & a1 \\ a2 & a2 \end{bmatrix}$
<b>Subgraph</b> ( $m, s$ ) := SELECT * FROM m WHERE m.S IN s AND (m.O IN s OR isLiteral(m.O))	$\text{Subgraph}(M, \begin{bmatrix} a2 \\ \text{Column} \\ \text{int} \end{bmatrix}) =$ where $M$ = model of Figure 4

The set operators Union (+), Difference (−), and Intersection ( $\cap$ ) are another three important primitive operators. We define these on models, morphisms, and selectors by the corresponding set operations on their representation as relations. For example,

**Union**( $x, y$ ) := SELECT \* FROM x UNION SELECT \* FROM y

Note that applying the set operations to well-formed models may produce a model that is not well-formed.



**Figure 8: Examples of copying the model of Figure 4 using selector {a1, a2, a3, a4}**

The last two primitive operators are All and Copy. The operator  $All(m)$  returns a selector that contains only those nodes of  $m$  that denote the model elements of the model's meta-model, such as tables or columns in the relational meta-model. For example, for the model of Figure 4 the operator All yields the selector {a1, a2, a3, a4} and filters out all auxiliary nodes, such as Table or PrimaryKey, that are used in the graph encoding.

Frequently, it is important to ensure that a given node identifier is used in exactly one model. Furthermore, unique node IDs make it possible to refer to model elements across model boundaries. For these reasons, we use the operator Copy to create a copy of a model  $m$  in which the selected node IDs are replaced by new, uniquely created IDs. In the following definition of Copy, the function uniqueOID() generates a unique OID on each call, and the function ifNULL( $x, y, z$ ) returns  $y$  whenever  $x$  is a NULL value,  $z$  otherwise. If  $s=All(m)$ , the output morphism  $m'_m$  is a bijection between  $All(m')$  and  $All(m)$ .

```
Copy(m, s) :=
  m'_m = SELECT uniqueOID(), s.V FROM s;
  m' = SELECT ifNULL(T1.L, m.S, T1.L), m.P,
            ifNULL(T2.L, m.O, T2.L)
  FROM m, m'_m as T1, m'_m as T2
  LEFT OUTER JOIN ON m.S=T1.R, m.O=T2.R;
  return (m', m'_m);
```

Figure 8 illustrates the operator Copy. The operator takes as input the model  $m$  of Figure 4 and selector {a1, a2, a3, a4} =  $All(m)$ . As a result of copying, a new model has been created (on the right), in which the nodes IDs a1, a2, a3, a4 have been replaced by the generated unique IDs a5, a6, a7, a8, respectively.

## 4.2 Derived operators

The derived operators are functional combinations of other operators. For example, consider the definitions shown below.

```
operator Range(map)
  return Domain(Invert(map));

operator RestrictRange(map, selector)
  return Invert(RestrictDomain(Invert(map), selector));

operator Traverse(selector, map)
  return Range(RestrictDomain(map, selector));

operator Restrict(map, m1, m2)
  return RestrictRange(RestrictDomain(map, All(m1)), All(m2));
```

The Range of a morphism is obtained as the domain of an inverted morphism, by combining the primitive operators Domain and Invert of Table 1. Similarly, RestrictRange is specified in terms of the operator RestrictDomain by first inverting the

input morphism, then applying `RestrictDomain`, and finally inverting the resulting morphism once again.

The third operator, `Traverse`, was used in our motivating scenario for locating the *dl* images of the elements deleted from the relational schema *sl*. To “traverse” the morphism, it is first domain-restricted by the selector, and the range of the restricted morphism is returned as output.

The last operator, `Restrict`, confines the domain and range of a morphism to the elements of two models *m1* and *m2*. Notice that the definitions of the derived operators above are expressed declaratively, allowing the implementations to be optimized.

### 4.3 Extract and Delete

Extracting and deleting portions of models are operations that are heavily deployed in metadata applications. To perform these operations, we propose the generic operators `Extract` and `Delete`. The operator `Extract` is applied as follows:  $\langle m', m'_m \rangle = \text{Extract}(m, s)$ . The inputs are a well-formed model *m* and a selector *s* that identifies the set of nodes to be extracted. The output model *m'* satisfies the following properties: (i) *m'* contains all selected nodes, (ii) *m'* is a well-formed model, (iii) *m'* is an equally or less expressive model than *m*, i.e., *m* can represent all information of *m'*, and (iv) *m'* is a “minimal” model that satisfies (i)–(iii). Condition (ii) may require that unselected “support” elements be included in *m'*. Condition (iii) can be characterized formally in terms of dominance and information capacity as suggested in [15,19]. The morphism *m'\_m* is an injective function from  $\text{All}(m')$  to  $\text{All}(m)$ , i.e., each model element of *m'* has at most one counterpart in *m*.

In general, a model may contain implicit information, such as transitive relationships between model elements. In such cases, the result of `Extract` may need to make such information explicit. For example, consider a class diagram with three classes A, B, C, and two explicit subclass definitions: A is a subclass of B, and B is a subclass of C. Due to condition (iii), `Extract(m, {A, C})` should return a class diagram in which A is defined as a subclass of C. This example illustrates that extraction is a rich operation, whose semantics and implementation may be non-trivial.

Conceptually, the semantics of the operator `Extract(m, s)` can be realized using the following algorithm:

1. Create a “closure” of *m*, i.e., a model *m'* in which all implicit information of *m* is represented explicitly.
2. Assign  $s' = s$ , where *s'* is a temporary selector.
3. For each *x* in *s'*, extend *s'* with elements needed to satisfy conditions (ii) and (iii).
4. Apply 3 until a fixpoint is reached, i.e., *s'* does not change.
5. Extract subgraph *t'* induced by *s'* as  $t' = \text{Subgraph}(m', s')$ .
6. Obtain a “cover” of *t'*, i.e., a minimal model *t* that is semantically equivalent to *t'*.
7. Return `Copy(t, All(t))` as result of extraction. Notice that the operator `Copy` (Section 4.1) returns a model and a mapping.

Deleting a selected portion of a model can be defined as extraction of the unselected portion. Thus, we define

```
operator Delete(m, s)
  return Extract(m, All(m) – s);
```

Note that the nodes of *s* that do not represent the model elements of *m*, i.e., are not members of  $\text{All}(m)$ , have no impact on the result of deletion due to applying  $\text{All}(m) - s$ .

#### 4.4 Match

The purpose of Match is to uncover how two models “correspond” to each other. It takes two models as input and returns a morphism between them. Match is inherently heuristic. So like the previous literature on Match [24], we do not offer a formal definition of what constitutes a correct output morphism. In general, matching two schemas requires information that is not present in the schemas and cannot be fully automated. Hence, a human engineer needs to review and adjust the suggestions produced by an automatic procedure, either in a post-processing step or iteratively.

#### 4.5 Merge

To combine two models into one, we utilize the operator Merge, applied as  $\langle m, m1\_m, m2\_m \rangle = \text{Merge}(m1, m2, map)$ . If the input models  $m1$  and  $m2$  are well-formed, Merge should produce a well-formed model  $m$  that (i) is at least as expressive as each of the input models, i.e., capable of representing the information contained in both models, and (ii) is “minimal”, i.e., deleting any element makes the model less expressive than one of the input models. The third parameter to Merge is a morphism  $map$  that describes elements of  $m1$  and  $m2$  that are equivalent and should be “merged” into a single element in  $m$ . The output morphisms  $m1\_m$  and  $m2\_m$  identify the counterparts of the elements of  $m1$  and  $m2$  in the merged model  $m$ .

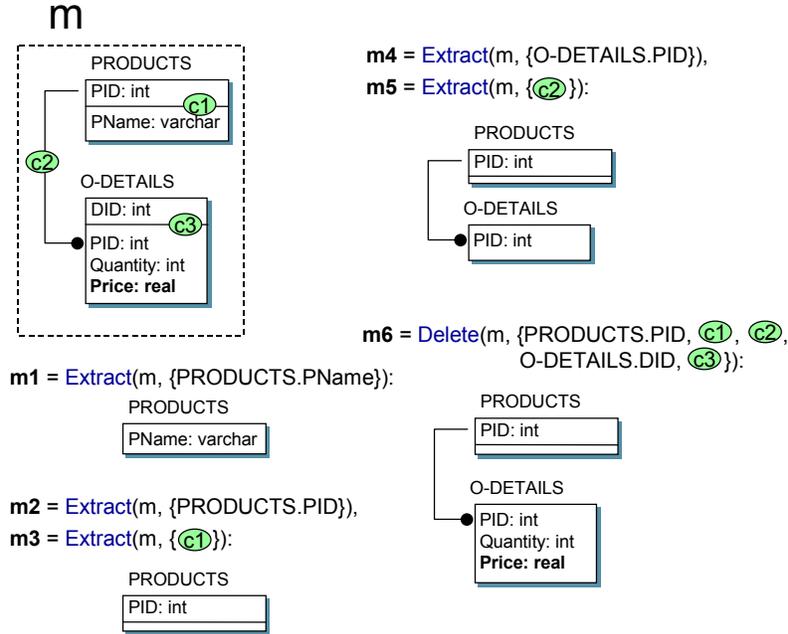
The conceptual definition of Merge given above does not say anything about the naming and ordering of model elements. For example, it does not prescribe that the attribute names of  $m1$  take precedence over those of  $m2$ , or the other way around. These details are not considered to be part of the semantics of Merge because they inherently involve end-user decision making. They are discussed in Section 5.7.

### 5. IMPLEMENTATION

In this section we discuss our implementation of the conceptual structures and operators presented above. We have found that the relations that were used in Section 3 as standard mathematical representation of graphs actually are a convenient implementation structure too. Our graph representation is based on the classical relational data model, in which node identifiers are constants that can be shared across models. We chose a relational approach instead of an object-oriented one (e.g., the one in [5]) to simplify the implementation and specification of the operators, which can often be done using SQL. Our relational graph model is based on the W3C’s Resource Description Framework (RDF).

For encoding relational schemas, XML schemas, and SQL views as graphs we use the following approach. Our meta-model for relational schemas is based on OIM [8]. For example, the model elements of a relational schema comprise tables, columns, and constraints; a table contains an ordered list of columns, each of which has a type; tables and columns carry names; the constraints are specialized into primary key, unique key, non-null, or referential constraints; a referential constraint refers to two columns, one of which is a foreign key and the other is a primary key; etc. Our graph representation of XML schemas builds on XML DOM. The graph representation of SQL views that we deploy is comparable to a parse tree produced by an SQL processor (see Figure 16 in Section 7). All clauses, statements, alias definitions, functional terms, etc. are represented as separate nodes. A view graph does not replicate the names of attributes and relations used in schemas, but refers directly to the respective nodes in the schema graphs.

The output of the primitive operators is defined uniquely in Section 4, except for the operator All, which is implemented differently for each meta-model. For example, for relational schemas the implementation of All is specified as follows:



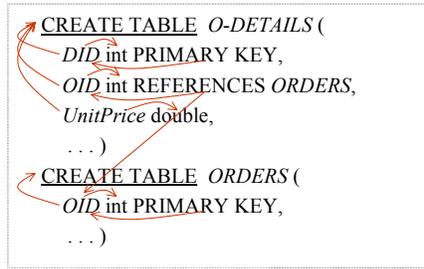
**Figure 9: Examples of extraction and deletion from a relational schema  $m$  (output morphisms not shown)**

$All(m, s) := SELECT m.S FROM m WHERE m.P=type AND m.O IN \{Table, Column, PrimaryKey, UniqueKey, NonNull, ReferentialConstraint\}$

## 5.1 Extract and Delete

To describe our implementation of the Extract and Delete operators we focus on the relational schemas. Consider the schema  $m$  shown on the left of Figure 9. The primary key constraints on PID and DID are depicted as horizontal bars underlining the respective attributes. The referential constraint is shown as a line connecting PRODUCTS.PID and O-DETAILS.PID. Assume that in the graph representation of  $m$  the three constraints are denoted by the nodes  $c1$ ,  $c2$ , and  $c3$ , respectively. For brevity, we henceforth refer to the graph nodes representing the attributes of  $m$  simply by using their names.

Figure 9 illustrates six examples of extraction and deletion. The output morphisms  $m1_m, \dots, m6_m$  are omitted in the figure for compactness. The first example demonstrates extraction of the attribute PName yielding schema  $m1$ . Condition (ii) of Section 4.3 ensures that  $m1$  is a well-formed relational schema, i.e., attribute PName belongs to a relation and has a type specification. Applied to relational schemas, condition (iii) requires that the extracted schema contain all constraints present in the original schema that affect the selected elements. For example, extracting the attribute PRODUCTS.PID from  $m$  causes the primary key constraint  $c1$  to be extracted as well, yielding the schema  $m2$ . Dropping  $c1$  would violate (iii), since it would allow the attribute PID to contain duplicates and thus the original schema  $m$  could not represent all information of  $m2$ . Analogously, extracting O-DETAILS.PID from  $m$  (as schema  $m4$ ) needs to preserve the referential constraint  $c2$ , which in turn requires the presence of PRODUCTS.PID and its primary key constraint  $c3$ . Condition (iv) prevents any other attributes from appearing in  $m4$ .



**Figure 10: Example of existential dependencies in a relational schema**

In our prototype, the implementation of operator  $\text{Extract}(m, s)$  for relational schemas is based on the conceptual algorithm of Section 4.3. Steps 1 (“closure”) and 6 (“cover”) are equality assignments. Step 3 of the algorithm is implemented as follows:

- If  $s'$  contains constraint  $x$ , add to  $s'$  all attributes that participate in the constraint definition.
- If  $s'$  contains attribute  $x$ ,  $s'$  is extended to include (a) the enclosing relation of  $x$ , (b) the type definition of  $x$ , (c) the referential constraint or non-null constraint for  $x$ , (d) the primary key or unique key definition for  $x$ , but only when all attributes participating in the key definition are contained in  $x$ .

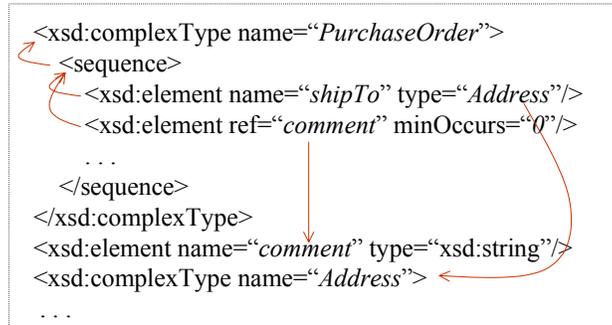
In Figure 9, schemas  $m3$  and  $m5$  illustrate the extraction of nodes that denote constraints. To illustrate case (d), consider a relation  $P(\text{Name}, \text{DOB}, \text{Addr})$  with a unique key constraint on  $(\text{Name}, \text{DOB})$ . According to the algorithm,  $\text{Extract}(m, \{P.\text{Name}\})$  yields  $P(\text{Name})$ . The unique key constraint is not included since  $P.\text{DOB}$  is not selected.

Notice that condition (iii) of  $\text{Extract}$  makes it impossible to delete a constraint on a relational attribute without deleting the attribute definition, or to delete the primary key attribute participating in a referential constraint without deleting its foreign key attribute. For example, consider schema  $m6$  in Figure 9. Selecting  $\text{PRODUCTS.PID}$  and the constraints  $c1$  and  $c2$  is not sufficient for deleting this attribute: the attribute  $\text{O-DETAILS.PID}$ , which is a foreign key on  $\text{PRODUCTS.PID}$ , is not selected; therefore, dropping  $\text{PRODUCTS.PID}$  would extend the set of possible values that  $\text{O-DETAILS.PID}$  may take beyond those contained in  $\text{PRODUCTS.PID}$  and hence violate condition (iii). In Sections 5.3 and 5.4, we present more flexible operators  $\text{ExtractMin}$ ,  $\text{DeleteHard}$ , and  $\text{DeleteSoft}$ , which allow such deletions by providing fewer consistency guarantees than  $\text{Extract}$  and  $\text{Delete}$ .

Extraction from XML schemas is implemented analogously to the above algorithm. Type references in XML schemas are treated similarly to the referential constraints in relational schemas. Currently, derived types are not supported.

## 5.2 Dependencies

As we observed above, the operators  $\text{Extract}$  and  $\text{Delete}$  disallow semantically questionable transformations on schemas, such as dropping arbitrary constraints, and are defined for schemas only. In general, deletion on models, which may or may not be schemas, needs to be done in a careful way to ensure that the consistency of the resulting model with respect to its meta-model is not violated. For example, consider the relation  $\text{ORDERS}$  shown at the bottom of Figure 10. If we were to delete just the definition of the table  $\text{ORDERS}$ , we risk getting an inconsistent model, in which fields like  $\text{OID}$  do not belong to any table. Or, if we delete the field  $\text{ORDERS.OID}$ , we might get a malformed referential constraint for  $\text{O-DETAILS.OID}$ , whose target key definition is now missing.



**Figure 11: Example of existential dependencies in an XML schema**

To deal with such consistency issues in a more general way, we exploit the concept of existential dependencies between model elements.

Figures 10 and 11 show examples of dependencies that hold between the elements of a relational schema, and between the elements of an XML schema. Each of the arcs specifies that the source element of the arc is existentially dependent on the target element. For example, in the relational schema of Figure 10, the attribute “UnitPrice” cannot exist without its type definition (arc from “UnitPrice” to *double*). Similarly, the primary key constraint in table O-DETAILS is malformed if the constrained field “DID” is missing. The referential constraint between the fields O-DETAILS.OID and ORDERS.OID spans two tables, and requires both a foreign key and a primary key. Analogously, in the XML schema of Figure 11, the definition of the element “shipTo” depends on the existence of the complex type “Address” as well as on the enclosing sequence element, etc.

As illustrated in Figures 10 and 11, dependencies are binary relations over the elements of a single model. Thus, we represent dependencies as intra-model morphisms, whose left elements are dependent on the right ones. To obtain the dependencies for a given model, we use the operator *Dependencies*, which invokes a non-generic implementation to compute the dependency morphism for the given model. For each supported model type, one such non-generic implementation is provided (one for relational schemas, another one for XML schemas, etc.). In our implementation, the operator *Dependencies* uses the arc types defined in the meta-model to determine what arcs are dependency arcs. For example, the arcs *column* and *SQLtype* of Figure 4 are marked as dependency arcs in our representation of the meta-model for relational schemas; the target of an arc of type *SQLtype* depends on the source, and the source of arc of type *column* depends on its target.

### 5.3 ExtractMin

A general intuition behind extraction is that we want to obtain a minimal model that contains the nodes in the selector and all those nodes and edges that are necessary to make the resulting subgraph a ‘complete’, well-formed model. Obviously, such model has to contain at least those nodes that are existentially required for the nodes in the selector. This minimalist subgraph can be obtained using the operator *ExtractMin* defined below, which uses an auxiliary derived operator *Reachable*.

```

operator ExtractMin(M, selector, dependencies)
  T = Subgraph(M, selector + Reachable(selector, dependencies));
  return Copy(T, All(T));

```

```
operator Reachable(selector, map)
  return Range(RestrictDomain(TransitiveClosure(map), selector));
```

The operator `ExtractMin` takes three parameters as input, a source model  $M$ , a selector that identifies the elements to be selected, and the dependency morphism for  $M$ . The operator returns the subgraph of  $M$  induced by the union of the nodes in the selector and all nodes that are required to satisfy the existential dependencies of the selected nodes. These required nodes are obtained using the operator `Reachable`.

To illustrate how `Reachable` works, imagine that it is called with parameters  $\{a,d\}$  as selector and  $\{\langle a,b \rangle, \langle b,c \rangle\}$  as the dependency morphism of model  $M$ . We get:  $\text{Reachable}(\{a,d\}, \{\langle a,b \rangle, \langle b,c \rangle\}) = \text{Range}(\text{RestrictDomain}(\{\langle a,b \rangle, \langle b,c \rangle, \langle a,c \rangle\}, \{a,d\})) = \text{Range}(\{\langle a,b \rangle, \langle a,c \rangle\}) = \{b,c\}$ . Thus, selecting  $\{a,d\}$  from model  $M$  yields  $\text{Subgraph}(M, \{a,d\} + \{b,c\}) = \text{Subgraph}(M, \{a,b,c,d\})$ . The resulting subgraph contains by definition all edges between  $\{a,b,c,d\}$  and their incident literals. Notice that the operator `Reachable` can be executed by the optimizer efficiently, without materializing the transitive closure. This observation is important, since the dependency closures of even moderately-sized models may contain hundreds of thousands of entries.

As another example, consider selecting a single node denoting the attribute “UnitPrice” from the model of the relational schema of Figure 10 using `ExtractMin`. As shown in the figure, the type definition of “UnitPrice” and the relation “O-DETAILS” are required for the attribute definition, so that the operator `Extract` returns a subgraph of the model that represents the relational schema

```
CREATE TABLE O-DETAILS (UnitPrice double)
```

Similarly, if a single node denoting the primary key of table `ORDERS` is selected, we get

```
CREATE TABLE ORDERS (OID int PRIMARY KEY)
```

In this case, the node identifying the table `ORDERS` is pulled out due to the transitive dependency of the primary key on the table definition via the attribute definition.

## 5.4 DeleteHard and DeleteSoft

As noted in Section 4.3, extracting a selected portion of a model can be viewed as deletion of the unselected portion. To support a broader range of model management scenarios, we define additional two variants of deletion, `DeleteHard` and `DeleteSoft`. Both operators remove a portion of a model referenced by a selector. The intuition behind `DeleteHard` is that we want to obtain a maximal consistent submodel without the selected nodes. It is defined as follows.

```
operator DeleteHard(M, selector, dep)
  toDelete = selector + Reachable(selector, Invert(dep));
  toKeep = All(M) – toDelete;
  return ExtractMin(M, toKeep, dep);
```

Essentially, the operator `DeleteHard` takes `All(M)` elements of  $M$ , subtracts from this set the elements to be deleted, and applies `ExtractMin` to extract the unselected portion of the model. To take the existential dependencies into account, `DeleteHard` extends the selector passed as input to include all elements of  $M$  that would become “dangling”, i.e., elements that are existentially dependent on the elements to be deleted. Such would-be dangling elements are obtained by passing the selector and the inverted dependency morphism to the operator `Reachable`. That is, the dependencies are traversed in the reverse direction.

Consider again the example in Figure 10. Imagine that we `DeleteHard` the nodes representing the attribute `O-DETAILS.UnitPrice` and the table `ORDERS`. The set of elements `Reachable` from these selected elements over the inverted dependency morphism are the foreign key constraint on `O-DETAILS.UnitPrice` and all attributes of `ORDERS` (to see that, the arcs in the figure need to be traversed in the reverse direction).

That is, the constraint and the table ORDERS with all its attributes will be removed, and we get the schema

```
CREATE TABLE O-DETAILS (DID int PRIMARY KEY, OID int)
```

In contrast to DeleteHard, the operator DeleteSoft removes each selected element only if it has no unselected dependent elements. That is, in the above example, the table ORDERS would not be deleted since it is referenced by the unselected foreign key on O-DETAILS.OID. The result of applying DeleteSoft for the same input parameters is shown below. Only O-DETAILS.UnitPrice has been removed.

```
CREATE TABLE O-DETAILS (
  DID int PRIMARY KEY,
  OID int REFERENCES ORDERS)
CREATE TABLE ORDERS (OID int PRIMARY KEY, ...)
```

The operator DeleteSoft is defined below. Instead of extending the selector to cover the would-be dangling elements, it is restricted to make sure that no unselected elements are removed. The selector that keeps the elements that cannot be deleted (cannotBeDeleted) is first obtained by collecting all elements which the unselected elements depend on. Now, the input selector is adjusted to eliminate all these undeletable elements. Finally, the operator ExtractMin is applied, just as in the operator DeleteHard.

```
operator DeleteSoft(M, selector, dep)
  cannotBeDeleted = Reachable(All(M) – selector, dep);
  toDelete = selector – cannotBeDeleted;
  toKeep = All(M) – toDelete;
  return ExtractMin(M, toKeep, dep);
```

Table 2 summarizes the differences between the operators discussed above and illustrates them using a single characteristic example for relational schemas:

**Table 2: Comparison of variants of extraction and deletion**

Operator	Example
Extract	Cannot extract a field without the constraints defined for the field.
ExtractMin	Can extract a field without the constraints defined for the field.
Delete	Cannot delete a constraint defined on a field without deleting the field.
DeleteSoft	Can delete a constraint defined on a field without deleting the field. Cannot delete fields referenced by unselected fields.
DeleteHard	Can delete fields even if they are referenced by unselected fields. In this case, dangling references would be deleted, too.

## 5.5 Diff

The Diff operator computes the difference between a model M and another model that is connected to M using a mapping map. Intuitively, the difference between two models is a sub-model of M that does not participate in the mapping map. In other words, to obtain the difference we eliminate from M all elements that do have matching counterparts in the other model. Thus, we define the operator Diff as shown below:

```
operator Diff(M, map)
  return Delete(M, Range(map));
```

Similarly to the operators DeleteSoft and DeleteHard, we provide additional two versions of the Diff operator: DiffSoft and DiffHard.

```
operator DiffSoft(M, map)
  return DeleteSoft(M, Range(map));

operator DiffHard(M, map)
  return DeleteHard(M, Range(map));
```

Notice that given the the differencing operators, we could define deletion as derived operations. For example, the operator Delete could be defined based on Diff as

```
operator Delete(M, s)
  return Diff(M, Id(s));
```

## 5.6 Match

In our prototype, the Match operator takes as input two models of the same kind, e.g., two relational schemas, and returns as output a morphism. We implemented Match using the Similarity Flooding (SF) algorithm, a graph-matching algorithm presented in [17]. The SF algorithm exploits the structure of the graphs to be matched and performs especially well for detecting the differences between two versions of a schema, which is the case in our motivating scenario and many other metadata applications.

The SF algorithm takes as input two graphs  $m1$  and  $m2$ , and a set of initial similarity values between the nodes of the graphs, expressed as a weighted binary relation *seed*. Each pair  $\langle l, r \rangle$  of *seed* carries a similarity value between zero and one. In a fixpoint computation, the algorithm iteratively propagates the initial similarity of nodes to the surrounding nodes, using the intuition that neighbors of similar nodes are similar. The output of the algorithm is another weighted binary relation.

In Section 3.2 we defined a morphism as a binary relation. To include weights in a morphism, we add to it a third attribute Sim that holds a similarity value for each pair of nodes. The primitive operators in Section 4.1 ignore this extra information. We implement the operator Match as

```
operator Match(m1, m2, seed)
  multimap = SFJoin(m1, m2, seed);
  multimap = Restrict(multimap, m1, m2);
  map = FilterBest(multimap);
  return <map, multimap>;
```

The operator SFJoin encapsulates the SF algorithm. As explained in [17], the *multimap* returned by the algorithm may contain a large fraction of the cross product of the nodes in  $m1$  and  $m2$ , and needs to be filtered. The operator FilterBest implements the filter suggested in [17], which exploits the stable-marriage property. In addition to filtering, we restrict the result of the SFJoin operator to the nodes that represent the model elements of  $m1$  and  $m2$  using the operator Restrict (Section 0). The input morphism *seed* is typically obtained using another auxiliary operator NGramMatch( $m1, m2$ ), which computes the similarities of literals in  $m1$  and  $m2$  based on the number of n-grams that they have in common. Alternatively, *seed* can be obtained by composition of morphisms. If *seed* is omitted, NGramMatch is invoked in SFJoin by default.

The above Match implementation returns both the filtered morphism *map*, and the unfiltered *multimap*. The morphism *map* can be adjusted by the engineer using a graphical tool by invoking the operator EditMap on the outputs of Match, e.g., as  $map = \text{EditMap}(map, multimap)$ . The graphical tool allows the engineer to inspect all candidate matches suggested in *multimap*.

The script used above for implementing the Match operator can be easily adapted to call other external schema matchers, which may deploy thesauri, analyze schema

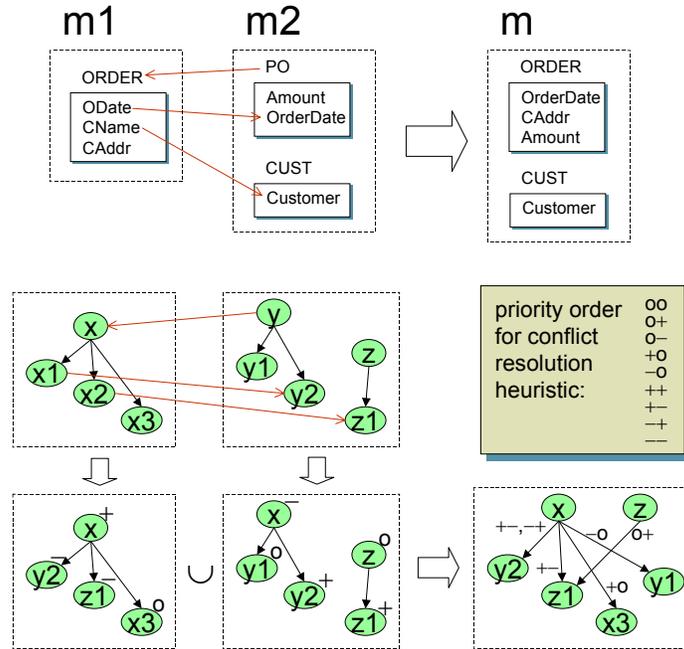


Figure 12: Merging two sample schemas

annotations, mine samples of instance data, reuse previous match results, etc., to reduce the manual post-processing effort.

## 5.7 Merge

We discuss our implementation of the Merge operator using the example in Figure 12. On the top, two sample models  $m1$  and  $m2$  get merged into  $m$  (the output morphisms are omitted). The morphism  $map$  is depicted using directed arcs. The direction of each arc establishes a preference between two model elements; when collapsing the two elements, the target element is kept in the output  $m$ , whereas the source element is discarded. For example, the attribute  $PO.OrderDate$  is kept and  $ORDER.ODate$  is discarded. Such preferences are not part of the semantics of the Merge operator (Section 4.5), but are essential for practical deployment. The input morphism  $map$  contains an extra attribute  $Dir$  to hold the direction of the arcs ( $\rightarrow$  or  $\leftarrow$ ). Before Merge is executed, a human engineer has a chance to specify the arc direction in a graphical tool by invoking the operator `EditMap`.

The bottom of Figure 12 depicts  $m1$  and  $m2$  as graphs. For brevity, the arc labels, type edges, and literals are omitted (compare to Figure 4). Node  $x$  corresponds to relation  $ORDER$ ,  $x1$  denotes  $ORDER.ODate$ , etc. The morphism  $map$  is  $\{\langle x, y, \leftarrow \rangle, \langle x1, y2, \rightarrow \rangle, \langle x2, z1, \rightarrow \rangle\}$ .

To implement the Merge operator, we developed an algorithm called `GraphMerge`, which we describe below. Similar to [11,23], the algorithm consists of three conceptual steps: node renaming, graph union, and conflict resolution.

1. In the first step, the graph nodes at the blunt ends of  $map$  are renamed to their targets at sharp ends, in both graphs  $m1$  and  $m2$ . The result of renaming is shown on the bottom left of Figure 12. Nodes  $y$ ,  $x1$ , and  $x2$  of both graphs have been renamed respectively to  $x$ ,  $y2$ , and  $z1$ .
2. In the second step, we do a graph union, i.e., a set union of two sets of edges, and obtain the graph depicted on the bottom right of the figure. This graph is not a well-

formed model, because the node  $z1$ , which used to represent the attribute CUST.Customer in  $m2$ , has now become an attribute of two different relations,  $x$  (ORDER) and  $z$  (CUST).

3. Such conflicts are resolved in the third and final step of the GraphMerge algorithm. The above conflict is eliminated by deleting either the edge between  $x$  and  $z1$ , or the edge between  $z$  and  $z1$ , effectively making Customer an attribute of either relation CUST or relation ORDER in the merged schema. The choice is made by a human engineer.

Step 3 is the costliest step of the algorithm, since it requires human feedback. To partially automate conflict resolution, we developed the following heuristic. Observe that in Figure 12 it seems more “natural” to keep the attribute Customer in relation CUST than to move it to ORDER. To generalize this observation, we track the origin of each edge in the merged graph, and assign to each edge a tag, such as  $+−$  or  $o+$ , which indicates whether each of the nodes incident at the edge was a source node of *map* ( $−$ ), a target node ( $+$ ) of *map*, or none of the two ( $o$ ) (these are the only three possible cases assuming that source and target nodes of *map* are disjoint). For example, the edge  $\langle x, z1 \rangle$  obtained by renaming from  $\langle x, x2 \rangle$  is tagged with  $+−$ , since  $x$  is a target node and  $x2$  is a source node of *map*. Analogously, the edge  $\langle z, z1 \rangle$  is tagged with  $o+$ , since  $z$  does not appear in *map* at all.

If we knew that  $o+$  edges are always preferred over  $+−$  edges, then, in a conflict  $\langle x, z1 \rangle$  could be eliminated without asking the engineer. We examined a variety of merge problems in the context of relational schemas, XML schemas, and SQL views, and established empirically a total order among all tag variations, which helps resolve many conflicts automatically in a way that matches human intuition. This order is shown in the middle right of Figure 12. Intuitively, edges between unchanged nodes ( $oo$ ) are least likely to be rejected in a conflict, and thus have the highest priority. Similarly, edges incident at  $+$  seem more likely to be preferred than those incident at  $−$ . Thus, Steps 2 and 3 are realized as follows. First, all edges in the merged graph are sorted by decreasing priority. Then, iteratively, each edge is taken off the top of the sorted list and is appended to an (initially empty) graph  $G$ . If appending the edge violates model consistency, it is rejected. Once all edges have been appended, the engineer examines the result and the choices made heuristically, and makes any necessary adjustments.

In the above description of the algorithm, we factored out an important aspect, the ordering of nodes within parent. To illustrate how we reestablish a correct order in the merged schema, consider Figure 12. Node  $y$  denoting the relation PO is renamed to  $x$ . Thus, when merging this node with the original  $x$  in  $m1$ , we move attributes  $y1$  (Amount) and  $y2$  (OrderDate) to the last position in the merged schema  $m$ . However, OrderDate “overrides” ODate, the first attribute in relation ORDER, and should remain at the first position. Hence, in schema  $m$ , the resulting order of attributes is OrderDate, CAddr, Amount.

The GraphMerge algorithm is summarized below:

**Algorithm** GraphMerge( $m1, m2, map$ )

$M := m1 \cup m2$ ;  $L :=$  empty list;  $G :=$  empty graph

**for each** edge  $e$  in  $M$  **do**

    rename nodes of  $e$  using *map*; assign tag to  $e$ ; append  $e$  to  $L$ ;

**end for**

sort edges in  $L$  by decreasing tag priority;

$maxN :=$  SELECT max( $M.N$ ) FROM  $M$ ;

**while**  $L$  not empty **do**

```

take edge  $e = \langle s, p, o, n \rangle$  off top of  $L$ ;
if tag( $e$ ) one of {"-o", "-+", "--"} then
   $n := n + \text{max}N$ ;
  if  $o$  is literal then continue loop end if
end if
if exists  $e' = \langle s, p, o, n' \rangle$  in  $G$  then
  replace  $e'$  in  $G$  by  $\langle s, p, o, \min\{n, n'\} \rangle$ ;
else if not conflictsWith( $\langle s, p, o, n \rangle$ ,  $G$ ) then
  append  $\langle s, p, o, n \rangle$  to  $G$ ; end if
end if
end while
return  $G$ 

```

The number  $\text{max}N$  is obtained as the highest existing value of the ordinal property  $N$  in  $m1$  and  $m2$  (compare Section 3.1). It is used to move the nodes hanging off renamed nodes to the last positions. To test for renamed nodes, we check whether the corresponding edge tag starts with  $-$ , i.e., is one of  $-o$ ,  $-+$ , or  $--$ . The literals belonging to such renamed nodes are removed, to ensure that, e.g., the relation corresponding to node  $x$  in the merged graph of Figure 12 will be named "ORDER" and not "PO". The function `conflictsWith()` checks whether appending a new edge to  $G$  causes a conflict.

The `GraphMerge` algorithm can be used for various kinds of models by implementing the function `conflictsWith()` appropriately. In our prototype, we deploy the algorithm for merging relational schemas, XML schemas, and SQL views. For example, conflict detection for relational schemas checks that relations cannot contain relations instead of attributes, or that attributes cannot be shared among relations, etc.

The Merge operator is implemented as follows:

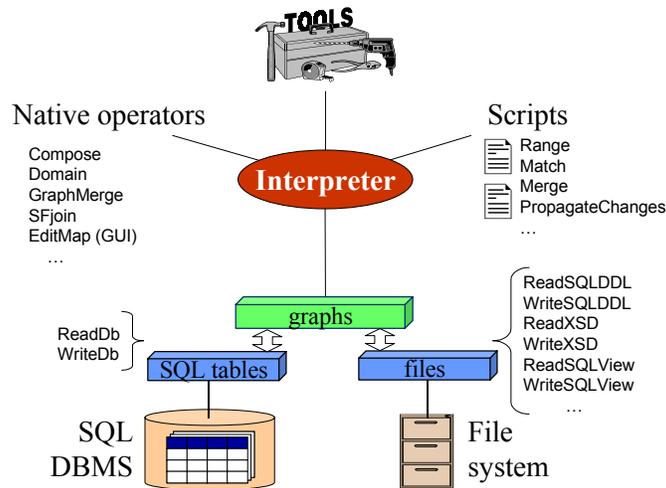
```

operator Merge( $m1, m2, \text{map}$ )
 $G = \text{GraphMerge}(m1, m2, \text{map})$ ;
 $s = \text{SELECT } L \text{ FROM } \text{map} \text{ WHERE Dir} = \text{"\to"} \text{ UNION}$ 
   $\text{SELECT } R \text{ FROM } \text{map} \text{ WHERE Dir} = \text{"\leftarrow"};$ 
 $m1\_G = \text{RestrictDomain}(\text{map}, \text{All}(m1) \cap s) + \text{Id}(\text{All}(m1) - s)$ ;
 $m2\_G = \text{RestrictDomain}(\text{map}, \text{All}(m2) \cap s) + \text{Id}(\text{All}(m2) - s)$ ;
 $\langle m, m\_G \rangle = \text{Copy}(G, \text{All}(G))$ ;
return  $\langle m, m1\_G * \text{Invert}(m\_G), m2\_G * \text{Invert}(m\_G) \rangle$ ;

```

Recall that `Merge` must also return morphisms from each of its input models to its output model. Thus, after applying `GraphMerge` to obtain the merged model  $G$ , we compute the morphisms  $m1\_G$  and  $m2\_G$ . The selector  $s$  contains all source nodes of  $\text{map}$ . For the example of Figure 12, we obtain  $m1\_G$  as the union of domain-restricted  $\text{map}$ ,  $\{\langle x1, y2 \rangle, \langle x2, z1 \rangle\}$ , which maps each renamed  $m1$  node to its new name, and the identity morphism on not renamed nodes,  $\{\langle x, x \rangle, \langle x3, x3 \rangle\}$ . Finally,  $G$  is copied to make the node IDs of the output model  $m$  unique, and the morphisms  $m1\_G$  and  $m2\_G$  are composed with `Invert`( $m\_G$ ), so they range over  $m$  instead of  $G$ .

The `GraphMerge` algorithm does not "invent" new model elements or establish new relationships between the existing elements. Therefore, the operator `Merge` as implemented above cannot reorganize schemas to resolve structural conflicts. For example, consider two XML schemas,  $S1$  with element `FullName` and  $S2$  with elements `FirstName` and `LastName`. Merging  $S1$  and  $S2$  should ideally create a new complex type `Name` with subordinate elements `FirstName` and `LastName`. Currently, we are working on addressing such structural conflicts by using  $n$ -way merges, in which intermediate schemas  $S_j$  are used for describing the desired structural transformations.



**Figure 13: Architecture of the prototype**

In Section 4.5 we postulated two “semantic” conditions that Merge should satisfy. Our implementation does not automatically ensure that condition (i) holds. For example, the engineer might decide to “override” a non-null constraint on an attribute in one schema S1 by a primary key constraint of the other schema S2, in which case the output model would be less expressive (i.e. more constrained) than S1. Although this flexibility is often desirable in practice, we are working on a more restrictive version of Merge that always guarantees to satisfy (i) and (ii).

## 6. PROTOTYPE

In this section, we describe our prototype, called Rondo<sup>2</sup>, in more detail. Its architecture is shown in Figure 13. Its central component is an interpreter that executes scripts. The interpreter can be run from the command line, or invoked programmatically by external applications and tools. Its main task is to orchestrate the data flow between the operators. The operators can be defined either by providing a native implementation, or by means of scripts. For example, a native operator like ReadSQLDDL reads a text document containing the definition of a relational database and creates its graph representation, whereas WriteSQLDDL exports the graph back as text. Similarly, two native operators ReadDb and WriteDb load and store arbitrary graphs in an SQL DBMS. Native operators are defined in scripts using statements like

```
alias ReadSQLDDL <Java class name>;
```

Other operators that have been implemented natively include all primitive operators of Section 4.1, operators that launch GUIs for editing morphisms and selectors, such as EditMap or EditSelector, schema translation and conversion operators, and the operators SFjoin and GraphMerge. All other operators, such as Range, Match, or Merge, are implemented by scripts presented in the previous sections. The specification of the commonly used native or derived operators can be grouped in a single script and utilized in other scripts using include statements.

The interpreter provides a debugging facility that allows examining the execution traces of complex scripts, and supports flexible handling of the input and output parameters of operators. For example, if an operator returns more than one argument (as

<sup>2</sup> Rondo: a musical work in which the main theme returns a number of times. A demo of the prototype is available for download at <http://www-db.stanford.edu/~melnik/mm/rondo/>

does our implementation of the operator Match), some of which are not used subsequently (as in script PropagateChanges in Section 2), they can be tacitly ignored.

For minimizing the amount of GUI programming needed for visualizing various kinds of models, we used the following technique. We require an operator like WriteSQLDDL to output not only the textual representation of the model, but also a data structure that describes how the terms in the text relate to the model elements, or graph nodes. In this way the schema elements shown in Figure 15 enclosed in boxes are associated with the graph nodes representing those elements, and the GUI operators EditMap and EditSelector can be used in exactly the same way for relational schemas (Figure 15) or SQL views (Figure 16).

At the current stage, our prototype supports the basic features of SQL DDL, XML Schema, RDF Schema, and SQL views, and, in preliminary form, UML. To introduce a new modeling language in the prototype, two steps are required. First, the import/export operators need to be provided, which ensure lossless round-tripping from the native format to graphs and back. Second, several callbacks need to be implemented for supporting the operators All, Extract, and GraphMerge.

The code breakdown of the prototype is shown in Figure 14. A large share of the implementation effort was due to the graph APIs responsible for in-memory representation and manipulation of graphs and morphisms, and the database support. The key generic model-management functionality comprises less than 7K lines of code. It includes the interpreter (2050), primitive operators (660), SFjoin (1760) and GraphMerge (700) implementations, as well as the generic GUI operators (1400). The non-generic part is essentially divided among the code needed to support SDL DDL, XML schemas, and SQL views. The smallest portion of code is due to converters: XSD2SQL (260), SQL2XSD (250), View2Morphism (90), and Morphism2View (200). The compactness of the converters is mostly due to the fact that they operate on the internal graph representation using expressive queries. The total amount of code in the prototype is below 24K lines. The total scripting code developed so far is measured in hundreds of lines. The scenarios shown in the article run in a few seconds on a 600 MHz laptop with 256 MB of memory.

Further scenarios that we implemented include a reintegration scenario from the context of version management, iterative merge, a warehousing scenario, in which we extract a subset of the schema that is sufficient to answer a given set queries, and a view reuse scenario. Due to space limitations, we cannot present all of them in this article. The view reuse scenario is in Section 7. Among other aspects, it illustrates how views can be merged, presents the GUIs used in our prototype, and demonstrates the use of the operators Morphism2View and View2Morphism. The reintegration scenario is covered in Section 8.

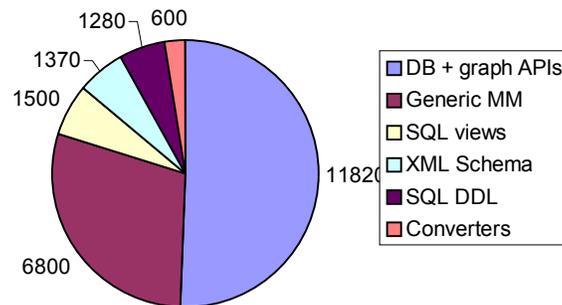


Figure 14: Code breakdown in prototype (in lines of code)

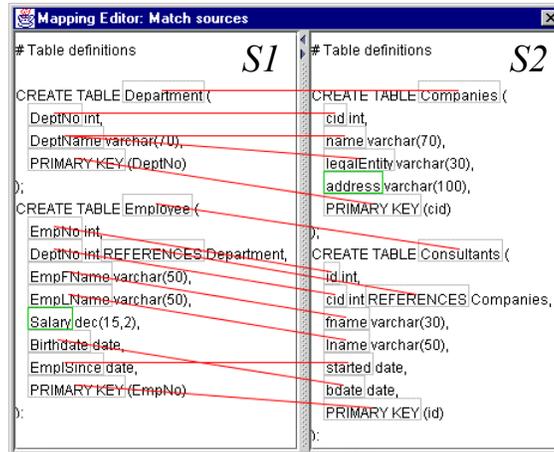


Figure 15: Morphism between sources  $S1$  and  $S2$

## 7. VIEW-REUSE SCENARIO

In this section, we examine another scenario, which illustrates the use of the operators presented in this article for addressing a typical data warehousing task. Consider adding a new source  $S2$  to a data warehouse  $D$ . Assume that  $S2$  is similar to an existing source  $S1$ . The morphism  $S1\_S2$  between the two source schemas is shown in Figure 15. Let an existing SQL view  $vS1\_D$  describe how the instances of  $S1$  populate  $D$ . The view  $vS1\_D$  is depicted in the middle of Figure 16 (the relevant portion of the warehouse schema can be seen in the CREATE VIEW clause). Our goal is to reuse the view  $vS1\_D$  for importing  $S2$  data into  $D$ , i.e., creating the view  $vS2\_D$ . Conventionally, this problem is solved manually involving a tiresome and error-prone renaming of the attribute and relation names of  $vS1\_D$  based on the similarities between  $S1$  and  $S2$ . In our prototype, we obtain  $vS2\_D$  using the following script:

1.  $S1\_S2 = \text{Match}(S1, S2)$ ;
2.  $S1\_D = \text{View2Morphism}(vS1\_D)$ ;
3.  $S2\_D = \text{Invert}(S1\_S2) * S1\_D$ ;
4.  $vS2\_D' = \text{Morphism2View}(S2\_D)$ ;
5.  $\text{map} = \text{Match}(vS2\_D', vS1\_D, \text{Invert}(S1\_S2))$ ;
6.  $vS2\_D = \text{Merge}(vS2\_D', vS1\_D, \text{map} + S1\_S2)$ ;

First, we match  $S1$  and  $S2$  to determine the correspondences between the schemas. As can be seen in Figure 15, some of the elements of  $S1$  and  $S2$  remain unmatched, whereas others, such as  $\text{Department.DeptName}$  are matched to two elements,  $\text{Companies.name}$  and  $\text{Companies.legalEntity}$ . In Step 2, we extract the morphism  $S1\_D$  from the view definition  $vS1\_D$  using a non-generic operator  $\text{View2Morphism}$ . For example, the morphism  $S1\_D$ , which is omitted in the figures for brevity, associates the attribute  $\text{Personnel.Pname}$  with two attributes,  $\text{Employee.EmpFName}$  and  $\text{Employee.EmpLName}$ , etc. Next, we compute the morphism  $S2\_D$  by composition. In Step 4, a “template” view definition  $vS2\_D'$  is generated from  $S2\_D$  using another non-generic operator  $\text{Morphism2View}$ . It is shown on the left of Figure 16. Morphism  $S2\_D$  contains no information as to how the values of the attribute  $\text{Personnel.Affiliation}$  are obtained from  $\text{Companies.name}$  and  $\text{Companies.legalEntity}$ . Therefore, a functional term  $\text{fct1}$  is generated in  $vS2\_D'$  as a placeholder.

In Step 5, the template  $vS2\_D'$  and the existing view  $vS1\_D$  are matched, using as a seed the morphism between  $S1$  and  $S2$ . The resulting morphism, after minor manual corrections, is depicted in Figure 16. Finally, in Step 6 both view definitions are merged

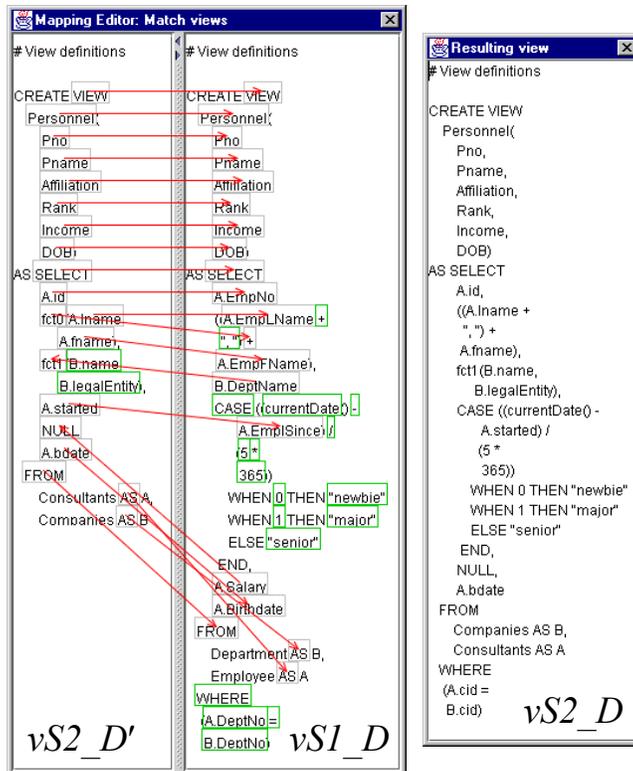


Figure 16: Merging two SQL views

to obtain  $vS2\_D$ , shown on the right. Notice that the function symbol `fct0` has been correctly replaced by the nested concatenation, whereas `fct1` was left as is. The unmatched `WHERE` clause was borrowed from  $vS1\_D$ ; the attribute references have however been correctly replaced by `Companies.cid` and `Consultants.cid`. To achieve that, the morphism *map* passed to `Merge` is extended to include  $S1\_S2$ . The heuristic deployed in the `GraphMerge` algorithm produces  $vS2\_D$  fully automatically, due to relative simplicity of the input views.

## 8. REINTEGRATION SCENARIO

In this section, we illustrate another scenario called reintegration, or 3-way merge. The reintegration problem arises when a model is modified independently by several engineers or tools. We focus on the case when there are two such independent modifications. Assume that model  $m$  was changed independently into  $m1$  by Ann and into  $m2$  by Bob. Our goal is to obtain the reconciled model  $m3$  that incorporates the changes done by Ann and Bob, and the mappings  $m\_m3$ ,  $m1\_m3$  and  $m2\_m3$  that describe how the models  $m$ ,  $m1$ , and  $m2$  relate to the reconciled version  $m3$ .

Consider the example in Figure 17. The original (relational) schema  $m$  is depicted on the top of the figure. In table `ORDERS` in schema  $m$ , employees are represented by an opaque identifier. To store employees' names, Ann creates the table `EMPLOYEEES` and makes `ORDERS.EID` a foreign key into the new table. Also, she deletes `ORDERS.PONum` and `O-DETAILS.UnitPrice` and adds `PRODUCTS.PDesc`. Meanwhile, Bob creates the table `BRANDS` and replaces the attribute `PRODUCTS.Brand` by a foreign key pointing to the new table. In addition, he adds a new attribute `PRODUCTS.ISIC` that holds the classification description of products. He deletes `DETAILS.UnitPrice`, just as Ann, and in addition he also deletes `DETAILS.Discount`.

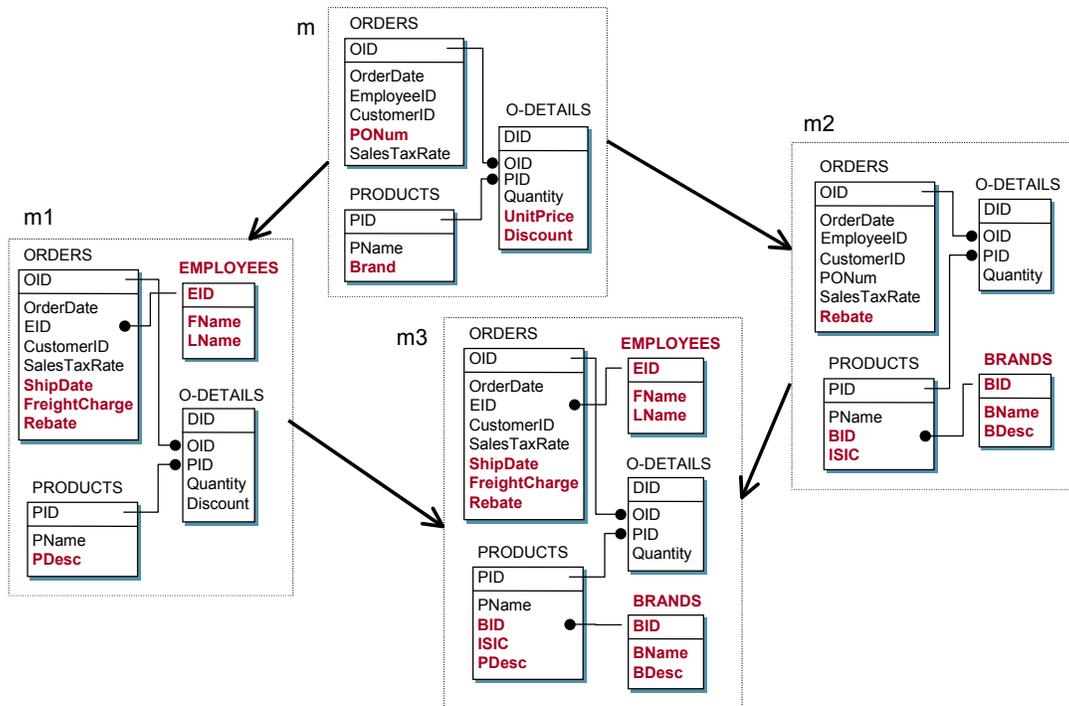


Figure 17: Reintegration scenario (3-way merge)

One way of obtaining  $m_3$  is to simply merge  $m_1$  and  $m_2$ . That is, in the script shown below, we first match  $m_1$  and  $m_2$  (line 1) and apply the Merge operator (line 2). To compute the mapping  $m\_m_3$ , we need to know how  $m$  corresponds to each of  $m_1$  and  $m_2$ . So, we match them in lines 3-4. Now, each of the compositions  $m\_m_1 * m_1\_m_3$  and  $m\_m_2 * m_2\_m_3$  describes a part of the mapping from  $m$  to  $m_3$ . To obtain  $m\_m_3$ , we combine both compositions in line 5.

```

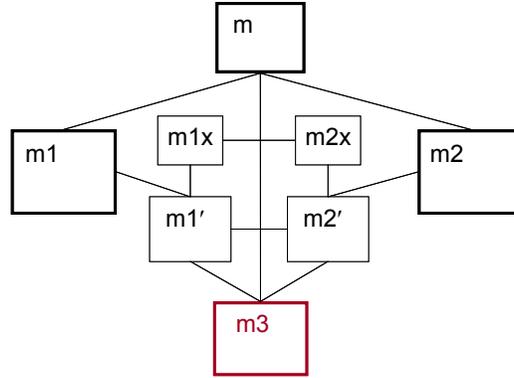
operator ReintegrateFirstCut(m, m1, m2)
1. m1_m2 = Match(m1, m2);
2. <m3, m1_m3, m2_m3> = Merge(m1, m2, m1_m2);
3. m_m1 = Match(m, m1); // or given
4. m_m2 = Match(m, m2); // or given
5. m_m3 = m_m1 * m1_m3 + m_m2 * m2_m3;

```

6. return <m3, m\_m3, m1\_m3, m2\_m3>;

The above approach has two major weaknesses. First, we have to apply the Match operator three times, each potentially requiring expensive human intervention. In practice,  $m\_m_1$  and  $m\_m_2$  could be tracked automatically by the schema editing tool used by Ann and Bob. Still, matching  $m_1$  and  $m_2$  from scratch can be costly. Second, the above script discards all deletions done exclusively by either Ann or Bob. That is, `ORDERS.PONum` and `O-DETAILS.Discount` would appear in  $m_3$  albeit both have been deleted. `O-DETAILS.UnitPrice` would, however, be correctly removed.

To address the first problem, we could modify the above script by moving lines 3-4 to the top and obtaining  $m_1\_m_2$  as the composition  $m_1\_m_2 = \text{Invert}(m\_m_1) * m\_m_2$ . By doing so, however, we duplicate the equivalent additions done by both Ann and Bob, since the added equivalent elements have no counterparts in  $m$  and hence their correspondences get lost upon composition. That is, after executing such modified script, `ORDERS.Rebate` would appear in  $m_3$  twice. And yet, we could use  $m_1\_m_2$  computed by composition to drive the match between  $m_1$  and  $m_2$ , as in  $m_1\_m_2 = \text{Match}(m_1, m_2,$



**Figure 18: Schematic representation of the reintegration scenario**

$\text{Invert}(m\_m1) * m\_m2 + \text{NGramMatch}(m1, m2)$ ). Moreover, when  $m1$  and  $m2$  are large, it may be more effective to extract only the new portions of  $m1$  and  $m2$  and match those.

To address the second problem, which is due to losing deletions done exclusively by Ann or Bob, we could apply to  $m1$  all deletions done in  $m2$ , and likewise apply to  $m2$  all deletions of  $m1$ . We incorporate both ideas in the script below:

```

operator Reintegrate(m, m1, m2)
1. m_m1 = Match(m, m1); // or given
2. m_m2 = Match(m, m2); // or given
3. <m1', m1'_m1> = Delete(m1, Traverse(All(m) - Domain(m_m2), m_m1));
4. <m2', m2'_m2> = Delete(m2, Traverse(All(m) - Domain(m_m1), m_m2));
5. <m1x, m1x_m1'> = Extract(m1', Traverse(All(m1) - Range(m_m1), Invert(m1'_m1)));
6. <m2x, m2x_m2'> = Extract(m2', Traverse(All(m2) - Range(m_m2), Invert(m2'_m2)));
7. m1x_m2x_core = m1x_m1' * m1'_m1 * Invert(m_m1) *
    m_m2 * Invert(m2'_m2) * Invert(m2x_m2');
8. m1x_m2x = Match(m1x, m2x, m1x_m2x_core + NGramMatch(m1x, m2x));
9. m1'_m2' = Invert(m1x_m1') * m1x_m2x * m2x_m2' +
    m1'_m1 * Invert(m_m1) * m_m2 * Invert(m2'_m2);
10. <m3, m1'_m3, m2'_m3> = Merge(m1', m2', m1'_m2');
11. m1_m3 = Invert(m1'_m1) * m1'_m3;
12. m2_m3 = Invert(m2'_m2) * m2'_m3;
13. m_m3 = m_m1 * m1_m3 + m_m2 * m2_m3;
14. return <m3, m_m3, m1_m3, m2_m3>;

```

To illustrate the script, consider the schematic representation in Figure 18. In line 3, we obtain the model  $m1'$  that contains all of  $m1$ , i.e., the model produced by Ann, without the elements deleted by Bob by way of  $m2$  (DETAILS.Discount). The expression  $\text{All}(m) - \text{Domain}(m\_m2)$  produces a selector that holds the elements of  $m$  that do not appear in  $m2$ . The images of these elements obtained by traversing  $m\_m1$  into  $m1$  are then deleted. Analogously,  $m2'$  contains all of  $m2$  without the elements deleted by way of  $m1$ , such as ORDERS.PONum.

In line 5, we extract a portion  $m1x$  of  $m1'$  that comprises only the elements added by Ann (e.g., PRODUCTS.PDesc) and their support elements (e.g., PRODUCTS). We achieve this by traversing the added elements  $\text{All}(m1) - \text{Range}(m\_m1)$  from  $m1$  to  $m1'$ . Line 6 does a similar job for  $m2x$ . Notice that line 5 could be realized alternatively as  $\langle m1x, m1x\_m1' \rangle = \text{Diff}(m1', m\_m1 * \text{Invert}(m1'_m1))$ ;

In line 7, we compute the mapping  $m1x\_m2x\_core$  between  $m1x$  and  $m2x$  to establish the correspondences between the support elements of  $m1x$  and  $m2x$ . This mapping is then used to drive the Match between the added portions in line 8. Here, the engineer

executing the script has a chance to state that `ORDERS.Rebate` added by Ann is equivalent to `ORDERS.Rebate` added by Bob. Notice that this Match is relatively inexpensive to perform, since we only have to reconcile the additions introduced by Ann and Bob.

In line 9, we compute the mapping between `m1'` and `m2'` to drive the Merge in line 10. To compose `m1'_m2'`, we need to consider both “paths” between the two models. One of them includes the matches between the added elements, `m1x_m2x`, and the other goes over the original model `m`. Similarly, the mapping `m_m3` is obtained in line 13 by joining two paths, one going through `m1` and the other through `m2`, portions of which are computed in lines 11-12. In line 14, the results of the script execution are returned.

## 9. RELATED WORK

Many individual aspects of model-management have been studied extensively in the literature, which is too voluminous to cite here. We highlight only some key aspects. In previous work [2,5,11,12,13,19,20,23], schemas were typically represented as graphs whose nodes denote classes of entities that participate in various semantically rich relationships, such as is-a, has-a, functional dependencies, etc. In our approach, the graphs are syntactic structures, whose semantics is opaque to many operators. Morphisms have been used under varying names in many systems, e.g., as schema correspondences in Clío [22]. To our knowledge, selectors have been first introduced in this article.

Past papers on model management reified mappings as models [5,9,23]. One of the surprises of the present work is how much leverage one can get out of simple morphisms. However, morphisms clearly have their limits. Section 7 presents a scenario in which SQL views are used as reified mappings to describe instance transformations. Reified mappings add complexity to scripts and operator implementations. A general treatment of reified mappings is subject of our ongoing work.

The operators discussed in [5] include Diff, Enumerate, and Apply. As explained in Section 5.5, in Rondo we implemented the operator Diff using extraction of the unmatched portion of one of the input models. Operators Apply and Enumerate are invoked by passing a selector to native Java code. The change propagation script of Section 2 is an alternative realization of the round-trip engineering scenario presented in [5].

A substantial effort has been devoted recently to schema matching. To minimize the amount of manual post-processing, existing schema matching tools deploy various techniques surveyed in [24], such as machine learning [4], etc. In our prototype, we use the structural matcher of [17], which is available for download from the authors' website.

Our definition of the Merge operator was influenced by the schema join operation of [1]. Schema merging has been further addressed e.g. in [11,19,23]. The algorithms suggested there can exploit rich relationship types that are not available in the GraphMerge algorithm that we developed, and do not take the ordering of model elements into account. Our heuristic deployed in GraphMerge is only an initial step in the challenging research issue of semiautomatic conflict resolution. In [21], this issue has been addressed in the context of ontology merging.

Schema translation across different modeling languages has been explored e.g. in [2,13]. The techniques presented there could be used for implementing a generic operator for generating one model from another. Currently, we are using a less general approach, in which each converter is implemented as a custom, non-generic operator.

To our knowledge, the generic operators Extract and Delete have first been investigated and implemented in this article. Our algorithm for Extract was inspired by the discussion of schema merging in [11].

Algebraic and model-theoretic semantics of model-management structures and operators has been considered in [1,19], but is still a new and largely unexplored area. Currently we are working on a state-based characterization of morphism semantics, building on the approach of [16]. The next section highlights some more of our future work.

## 10. OUTLOOK: STRUCTURAL VS. STATE-BASED SEMANTICS

The operators presented in this article treat models and mappings to a large extent as syntactic structures. The semantics of the operators is defined in terms of structural transformations on graphs. We refer to this semantics as *structural semantics*. Structural semantics is aligned closely with how developers think of metadata manipulations. A precise specification of structural semantics is crucial for deployment of real applications. In fact, many applications rely on the fact that certain tables or complex types bear specific names or are arranged in a certain order. For example, although the order of attributes in a relational table is, in theory, unimportant, many APIs, such as JDBC, allow developers to refer to the attributes of a table by their ordinal numbers rather than by names. Furthermore, applications may rely on the fact that data is stored in a denormalized representation or grouped in certain ways for performance reasons. For example, grouping products by orders or the other way around in a database schema may impact storage and query efficiency, although in both cases the same information is represented.

Focusing on structural semantics may simplify operator implementation. For example, exploiting the graph representation of metadata artifacts allows the operators like Match or Merge to be implemented in a generic fashion for different kinds of models, as we illustrated in Sections 5.6 and 5.7.

And yet, the effect of applying “syntactic” operators to models ultimately needs to be expressed in terms of what the operators do to the instances of these models, such as whole database states. We call this other kind of semantics *state-based semantics*. For example, conditions (i)–(iv) for the Extract operator (Section 4.3), or (i)–(ii) for Merge (Section 4.5) reflect the state-based semantics of these operators to a limited degree. State-based semantics allows us to specify and verify that a merged schema can indeed accommodate all information that can be represented by its source schemas. It helps us to uncover and analyze surprises that we get when an extracted schema actually captures more information than the original schema.

On the one hand, state-based semantics can be viewed as a descriptive tool: we describe what the implemented structural operators do to the instances of models. In fact, currently we are developing precise formal definitions of the state-based semantics of the operators implemented in Rondo. The mappings are interpreted as binary relations between instances of models, i.e., between whole database states. Thus, a mapping can describe any conceivable database transformation.

On the other hand, the state-based semantics could play a prescriptive role: it could provide a formal specification of how the execution of a script affects the possible database states described by the manipulated schemas. This specification should be precise enough so that an engineer could implement it unambiguously and tools could be built to analyze it, for example to identify possible undesired side effects.

For example, the state-based specification of the schema evolution task (a special case of change propagation when a view defined on a schema breaks due to schema changes), could be formulated roughly as follows: the updated view must expose all extra information that has been added to the updated schema while preserving the capability of answering queries that impact only the unchanged portion of the original schema. Ideally, such specification would remain executable – just as the scripts in Rondo. And in the best

possible world, it would be declarative and amenable to automated rewriting and optimization.

While substantial competence has been gained by tool vendors with respect to structural transformations of metadata artifacts, expressing the state-based semantics has proven to be hard. At this point, it is not entirely clear how the precise state-based characteristics of the key operators could look like. The state-based semantics of morphisms and selectors, as defined in this article, are work in progress. Moreover, we lack a way of telling whether a given set of operators is “complete” according to some metric. As a consequence, we cannot say what kind of metadata manipulation tasks are or are not amenable to a model-management solution. These open questions stress the importance of future work on state-based semantics for the development of metadata-intensive applications.

## 11. CONCLUSIONS

In this article we presented a programming platform for model management that implements all generic operators suggested so far in the literature. We explored the use of morphisms and selectors and introduced several novel generic operators. We discussed the operator semantics and the algorithms that we developed for implementing them. We showed that introducing a new model type like SQL DDL schemas in our prototype requires a moderate programming effort, but brings a large new class of model-management tasks within reach.

The main conclusions that we draw are the following:

1. One can solve practical problems using the model management operators.
2. The solutions require a relatively small amount of code.
3. One can get quite far using a relatively weak representation for models and mappings.
4. A precise specification of both structural and state-based semantics of the operators is needed to provide a satisfactory programming platform.

Our implementation experience, backed by the in-depth investigation of the individual operations by other researchers, suggests that the question raised in [7] is likely to have a positive answer, i.e., generic metadata management is in fact feasible. Even if we cannot handle subtle and complex cases, if we can solve a large class of non-trivial problems then we are offering a useful programming platform. Still, resolving the debate of [7] to the full extent can be done only by writing scripts for a substantial number of real applications and demonstrating that they work.

Other hard challenges remain open. Examples are providing meaningful semantic constraints on operators and proving that certain syntactic transformations “play by the rules”, or supporting more powerful mapping languages, which can be deployed directly to transform data instances. A salient non-technical challenge is acceptance by the developer community. As with each new programming paradigm, the willingness of engineers to learn a new way of approaching old problems is critical for success of generic model management.

## 12. ACKNOWLEDGMENTS

We thank Gio Wiederhold for his insightful comments on the article. We are grateful to Serge Abiteboul, Paolo Atzeni, Stefano Ceri, Alon Halevy, Martin Kersten, Renée Miller, Rachel Pottinger, and Gerhard Weikum for helpful discussions. This work was supported in part by a grant from the Database Group at Microsoft Research.

## REFERENCES

1. S. Alagic, P. A. Bernstein: A Model Theory for Generic Schema Management. Proc. DBPL, pp. 228-246, 2001
2. P. Atzeni, R. Torlone: Management of Multiple Models in an Extensible Database Design Tool. pp. 79-95, EDBT 1996
3. S. Bergamaschi, S. Castano, M. Vincini: Semantic Integration of Semistructured and Structured Data Sources, SIGMOD Record 28(1), pp. 54-59, 1999
4. J. Berlin, A. Motro: Database Schema Matching Using Machine Learning with Feature Selection. pp. 452-466, CAiSE 2002
5. P. A. Bernstein: Applying Model Management to Classical Meta Data Problems. pp. 209-220, CIDR 2003
6. P. A. Bernstein, A. Halevy, R. A. Pottinger: A Vision for Management of Complex Models. SIGMOD Record 29(4), pp. 54-63, 2000
7. P. A. Bernstein (moderator), L. Hass, M. Jarke, E. Rahm, G. Wiederhold (panelists): Is Generic Metadata Management Feasible? Panel, pp. 660-662, VLDB 2000
8. P. A. Bernstein, T. Bergstraesser, J. Carlson, S. Pal, P. Sanders, D. Shutt: Microsoft Repository Version 2 and the Open Information Model. Inf. Systems 24(2), p. 71-98, 1999
9. P. A. Bernstein, E. Rahm: Data Warehousing Scenarios for Model Management. pp. 1-15, Proc. Intl. Conf. on Conceptual Modeling (ER) 2000
10. S. Bowers, L. Declambre: On Modeling Conformance for Flexible Transformation over Data Models, Workshop on Transformation for the Semantic Web, July 2002
11. P. Buneman, S. B. Davidson, A. Kosky: Theoretical Aspects of Schema Merging. pp. 152-167, EDBT 1992
12. K. T. Claypool, E. A. Rundensteiner: Sangam: A Framework for Modeling Heterogeneous Database Transformations, ICEIS 2003
13. S. Cluet, C. Delobel, J. Siméon, K. Smaga: Your Mediators Need Data Conversion! pp. 177-188, SIGMOD 1998
14. S. Davidson, P. Buneman, A. Kosky: Semantics of Database Transformations. In B. Thalheim, L. Libkin, Eds., Semantics in Databases, LNCS 1358, pp. 55-91, 1998
15. R. Hull: Relative Information Capacity of Simple Relational Database Schemata. SIAM J. Computing, 15(3), pp. 856-886, Aug 1986
16. J. Madhavan, P. A. Bernstein, P. Domingos, A. Y. Halevy: Representing and Reasoning about Mappings between Domain Models. pp. 80-86, AAAI/IAAI 2002
17. S. Melnik, H. Garcia-Molina, E. Rahm: Similarity Flooding: A Versatile Graph Matching Algorithm and its Application to Schema Matching. ICDE 2002
18. S. Melnik, E. Rahm, P. A. Bernstein. Rondo: A Programming Platform for Generic Model Management. Proc. ACM SIGMOD 2003
19. R. J. Miller, Y. E. Ioannidis, R. Ramakrishnan: Schema Equivalence in Heterogeneous Systems: Bridging Theory and Practice. Information Systems 19(1), pp. 3-31, 1994
20. P. Mitra, G. Wiederhold, M. L. Kersten: A Graph-Oriented Model for Articulation of Ontology Interdependencies. p. 86-100, EDBT 2000
21. N. F. Noy, M. A. Musen: PROMPT: Algorithm and Tool for Automated Ontology Merging and Alignment. Proc. AAAI/IAAI 2000.
22. L. Popa, Y. Velegrakis, R. J. Miller, M. A. Hernández, R. Fagin: Translating Web Data. VLDB 2002
23. R. A. Pottinger, P. A. Bernstein: Merging Models Based on Given Correspondences. Proc. VLDB 2003

24. E. Rahm, P. A. Bernstein: A Survey of Approaches to Automatic Schema Matching. VLDB Journal 10(4), 2001

## ABOUT THE AUTHORS



Sergey Melnik

Sergey Melnik is finishing a Ph.D. in Computer Science at Leipzig University, Germany. He serves as an invited expert in the RDF Core Working Group at the World-Wide Web Consortium and is a recipient of a best student paper award (ICDE'02). He spent three years (1999-2002) as a visiting researcher in the Stanford Database Group where he worked on a variety of topics including metadata management, database optimization, information retrieval, and Semantic Web.



Erhard Rahm

Prof. Dr. Erhard Rahm has been the Chair for Databases at the Institute of Computer Science at the University of Leipzig since 1994 (<http://dbs.uni-leipzig.de>). His current research areas are metadata management, data warehousing, XML databases and bio databases. He is responsible for industrial and other third-party funded research projects, and is author of several books and numerous publications. In 1988 he received his PhD in Computer Science from the University of Kaiserslautern, and in 1993 his Postdoctoral Lecture Qualification. He was a visiting researcher at both the IBM Research Center in Hawthorne, NY, as well as Microsoft Research in Redmond, WA.



Philip A. Bernstein

Dr. Phil Bernstein is a researcher at Microsoft Corporation. Over the past 25 years, he has been a product architect at Microsoft and at Digital Equipment Corp., a professor at Harvard University and Wang Institute of Graduate Studies, and a VP Software at Sequoia Systems. During that time, he has published over 100 articles on the theory and implementation of database systems, and coauthored three books, the latest of which is "Principles of Transaction Processing for the System Professional" (Morgan Kaufmann, 1997). He holds a B.S. from Cornell University and a Ph.D. from University of Toronto. A summary of his current research on meta data management can be found at <http://www.research.microsoft.com/~philbe>.