

# Data Integration by Bi-Directional Schema Transformation Rules

Peter M<sup>c</sup>Brien

Dept. of Computing,  
Imperial College,  
London SW7 2BZ  
pjm@doc.ic.ac.uk

Alexandra Poulouvasilis

School of Computer Science  
and Information Systems,  
Birkbeck College,  
Univ. of London,  
ap@dcs.bbk.ac.uk

## Abstract

*In this paper we describe a new approach to data integration which subsumes the previous approaches of **local as view (LAV)** and **global as view (GAV)**. Our method, which we term **both as view (BAV)**, is based on the use of reversible schema transformation sequences. We show how LAV and GAV view definitions can be fully derived from BAV schema transformation sequences, and how BAV transformation sequences may be partially derived from LAV or GAV view definitions. We also show how BAV supports the evolution of both global and local schemas, and we discuss ongoing implementation of the BAV approach within the AutoMed project.*

## 1 Introduction

Data integration is a process by which several databases, which associated **local schemas**, are integrated to form a single virtual database with an associated **global schema**. Up to now, data integration approaches have been either **global as view (GAV)** or **local as view (LAV)** [12]. Examples of the GAV approach are TSIMMIS [4], InterViso [26] and Garlic [25] while examples of the LAV approach are IM [14] and Agora [15].

In this paper, we present a unifying framework for GAV and LAV, based on the use of reversible schema transformations. An important feature of our approach is that it is possible to extract a definition of the global schema as a view over the local schemas, and it is also possible to extract definitions of the local schemas as views over the global schema. Hence we term our approach **both as view (BAV)**. As we will see in this paper, one advantage of our BAV approach over GAV

and LAV is that it readily supports the evolution of both global and local schemas.

Before introducing BAV, we briefly review GAV and LAV via an example which we will be referring to again throughout the paper. Figure 1 shows two local schemas  $S_1$  and  $S_2$ , and an integration of these schemas into a global schema  $S_g$ .

Schema  $S_1$  is a database of past and present undergraduate students at a university, and also records the member of staff who was their undergraduate tutor. The `ug` relation holds student records identified by the key attribute `id` (key attributes are shown underlined). The optional attribute `left` is the date the student left the university (optional attributes are suffixed with a #). The attribute `sno` is a foreign key indicating which member of staff in the `tutor` relation is the student's tutor (foreign keys are shown in italics).

Schema  $S_2$  is a database of past and present PhD students at the university, and also records which member(s) of staff are the graduate supervisor(s) of the student. A PhD student keeps the same `id` if they previously studied an undergraduate course at the university.

Schema  $S_g$  is designed for use by the office handling enquiries for references about students, and needs to have within the `student` relation all current and past students, with the last degree that they studied, and when they left the university. The office also needs a record of the members of staff who last had contact with each student, as they will be asked to complete any requested references.

In GAV, the constructs of a global schema are described as views over the local schemas. These view definitions are used to rewrite queries over a global schema into distributed queries over the local databases. In our example, using a calculus-like nota-

|       |  |
|-------|--|
| $S_g$ | student(id,name,left#,degree)<br>monitors( <u>sno</u> ,id)<br>staff( <u>sno</u> ,sname,dept#)  |
| $S_1$ | ug(id,name,left#,degree,sno)<br>tutor( <u>sno</u> ,sname)                                      |
| $S_2$ | phd(id,name,left#,title)<br>supervises( <u>sno</u> ,id)<br>supervisor( <u>sno</u> ,sname,dept) |

**Figure 1.** Example global and local schemas

tion, the constructs of  $S_g$  can be defined as follows<sup>1</sup>:

$$\begin{aligned}
G_1 \quad & \text{student}(id, name, left, degree) = \{x, y, z, w \mid \\
& \quad \langle x, y, z, w, - \rangle \in \text{ug} \wedge \langle x, -, -, - \rangle \notin \text{phd} \\
& \quad \vee \langle x, y, z, - \rangle \in \text{phd} \wedge w = \text{'phd'}\} \\
G_2 \quad & \text{monitors}(sno, id) = \{x, y \mid \\
& \quad \langle x, -, -, y \rangle \in \text{ug} \wedge \langle x, -, -, - \rangle \notin \text{phd} \\
& \quad \vee \langle x, y \rangle \in \text{supervises}\} \\
G_3 \quad & \text{staff}(sno, sname, dept) = \{x, y, z \mid \\
& \quad \langle x, y \rangle \in \text{tutor} \wedge \langle x, -, - \rangle \notin \text{supervisor} \\
& \quad \vee \langle x, y, z \rangle \in \text{supervisor}\}
\end{aligned}$$

In LAV, the constructs of the local schemas are defined as views over the global schema, and processing queries over the global schema involves rewriting queries using views [13]. In our example, schema  $S_1$  can be defined by derivation rules  $L_1$  and  $L_2$  below, and schema  $S_2$  by derivation rules  $L_3$ ,  $L_4$  and  $L_5$ :

$$\begin{aligned}
L_1 \quad & \text{tutor}(sno, sname) = \{x, y \mid \\
& \quad \langle x, y, - \rangle \in \text{staff} \wedge \langle x, z \rangle \in \text{monitors} \\
& \quad \wedge \langle z, -, -, w \rangle \in \text{student} \wedge w \neq \text{'phd'}\} \\
L_2 \quad & \text{ug}(id, name, left, degree, sno) = \{x, y, z, w, v \mid \\
& \quad \langle x, y, z, w \rangle \in \text{student} \wedge \langle v, x \rangle \in \text{monitors} \\
& \quad \wedge w \neq \text{'phd'}\} \\
L_3 \quad & \text{phd}(id, name, left, title) = \{x, y, z, w \mid \\
& \quad \langle x, y, z, v \rangle \in \text{student} \wedge v = \text{'phd'} \\
& \quad \wedge w = \text{null}\} \\
L_4 \quad & \text{supervises}(sno, id) = \{x, y \mid \\
& \quad \langle x, y \rangle \in \text{monitors} \wedge \langle x, -, -, z \rangle \in \text{student} \\
& \quad \wedge z = \text{'phd'}\} \\
L_5 \quad & \text{supervisor}(sno, sname, dept) = \{x, y, z \mid \\
& \quad \langle x, y, z \rangle \in \text{staff} \wedge \langle x, w \rangle \in \text{monitors} \\
& \quad \wedge \langle w, -, -, v \rangle \in \text{student} \wedge v = \text{'phd'}\}
\end{aligned}$$

<sup>1</sup>We use a calculus-like notation for our examples in this paper because this is syntactically close to our actual intermediate query language, which is a comprehensions-based functional query language [21]. Such languages subsume high-level query languages such as SQL and OQL in expressiveness [3, 7].

In rule  $L_3$ , note the use of *null* for the value of the title, since no information can be derived from  $S_g$  regarding the topic title of a phd student's thesis. We will later give a more precise method of reasoning about such information, using the distinguished constant *void* to differentiate information that is not derivable from information that is actually recorded as *null*.

Note that rule  $L_1$  does *not* define the complete extent of the relation *tutor* in  $S_1$ , but just those tuples of *tutor* which can be derived from  $S_g$ . In particular, any staff member who supervises PhD students and tutors undergraduates will not appear in this view. This is a subtle aspect of the LAV approach, since in general it is not possible to take any  $S_g$  and provide views that can identify the source of all of the data. In the terminology of [11], this means  $L_1$  is a **complete** LAV rule. A similar argument holds for  $L_2$ . Both  $L_1$  and  $L_2$  could become **sound** LAV rules by the removal of the  $w \neq \text{'phd'}$  term, since they would then derive a superset of the source relations. Rules  $L_3, L_4$  and  $L_5$  are **exact** LAV rules, since they exactly derive the source relations from the global schema. In this particular example there is no possible exact LAV rule for  $L_1$ , as there is no way of distinguishing those members of staff who supervise PhD students and undergraduates from those staff who supervise PhD students only. A similar argument holds for why no exact definition of  $L_2$  is possible. A similar categorisation into sound, complete and exact is possible for GAV rules, but in practice it is usually assumed that the rules are exact in GAV.

A variation of LAV called GLAV [8] allows for the head of the view definition rules in GLAV to contain conjunctions of source relations as opposed to the single relations allowed by LAV, and to contain variables that do not appear in the body of the rule. This would be required if  $S_2$  had defined *supervises*(sno,title), where a GLAV rule could provide a definition *phd*(id, name, left, title)  $\wedge$  *supervises*(sno, title) without providing values to the title attribute. By contrast, a LAV approach would require that title appear in the global schema merely to allow the definition of this join.<sup>2</sup>

The principal disadvantage of GAV is that it does not readily support the evolution of the local schemas. For example, adding an age attribute to the *phd* relation in  $S_2$  would invalidate view definitions  $G_1$  and  $G_2$  above. Since this maintenance problem worsens as the number of source databases increases, GAV has been

<sup>2</sup>If we do not allow the use of *null* values, then  $L_3$  would not be valid. In that case a LAV rule defining *phd* would force the title attribute of  $S_2$  to appear in  $S_g$ , regardless as to whether it is otherwise required to be in the global schema. However, removing the term  $w = \text{null}$  from  $L_3$  would result in a valid GLAV rule.

criticised as not scaling well.

LAV isolates changes to local schemas to impact only on the derivation rules defined for that schema. Hence, adding the extra `age` attribute to `phd` would only effect rule  $L_3$  above. However, LAV has problems if one needs to change the global schema, since all the rules for defining local schemas as views of the global schema will need to be reviewed.

The remainder of this paper discusses how our BAV method may be used as a unifying framework for LAV and GAV, and how several practical advantages result. In Section 2, we present our BAV integration method. In Section 3, we show how BAV can capture all the semantic information that is present in LAV or GAV derivation rules. In particular, we show how to derive exact or complete LAV or GAV rules from BAV rules. Section 4 discusses how BAV supports the evolution of both global and local schemas. Section 5 discusses implementation of the BAV approach within the ongoing AutoMed project. Section 6 gives our concluding remarks and directions of further work.

## 2 The BAV Integration Method

In previous work [22, 17] we have developed a general framework to support schema transformation and integration in heterogeneous database architectures. The framework consists of a low-level **hypergraph-based data model (HDM)** and a set of primitive schema transformations defined for this model. Higher-level data models and primitive schema transformations for them are defined in terms of this lower-level common data model.

In our framework, schemas are incrementally transformed by applying to them a sequence of primitive transformation steps  $t_1, \dots, t_n$ . Each primitive transformation  $t_i$  makes a ‘delta’ change to the schema, adding, deleting or renaming just one schema construct. Each add or delete step is accompanied by a query specifying the extent of the new or deleted construct in terms of the rest of the constructs in the schema.

In [16, 22] we showed how our primitive transformations are rich enough to express all of the transformations commonly used in the integration of ER schemas, and we formally derived precisely what, if any, constraints on the actual data sources these transformations are conditional upon.

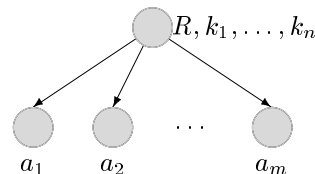
One advantage of using a low-level common data model such as the HDM is that semantic mismatches between modelling constructs are avoided. Another advantage is that it provides a unifying semantics for higher-level modelling constructs. In particular, in [18]

we described how our framework can be applied to different high-level modelling languages such as relational, ER and UML. In [19] we showed how XML data sources can also be handled within our framework.

### 2.1 A simple relational data model

For our purposes in the present paper, we assume that all schemas are specified in the very simple relational data model we define below. But our discussion here is equally applicable to any data modelling language supported by our framework.

In our simple relational model, schemas are constructed from primary key attributes, non-primary key attributes, and the relationships between them. The underlying graph representation of a relation  $R$  with primary key attributes  $k_1, \dots, k_n$  and other attributes  $a_1, \dots, a_m$  is:



Thus the relation `staff` of  $S_g$  is modelled as by the Rel scheme  $\langle\langle\text{staff}, \text{sno}\rangle\rangle$  with unary-tuple instances (since in this case  $n = 1$ ), together with two Att schemes  $\langle\langle\text{staff}, \text{sname}\rangle\rangle$  and  $\langle\langle\text{staff}, \text{dept}\rangle\rangle$  with binary-tuple instances. Thus we write  $x \in \langle\langle\text{staff}, \text{sno}\rangle\rangle$  to find instances of the key, and  $\langle x, y \rangle \in \langle\langle\text{staff}, \text{sname}\rangle\rangle$  and  $\langle x, y \rangle \in \langle\langle\text{staff}, \text{dept}\rangle\rangle$  to find values of the non-key attributes for this key.

The set of primitive transformations for schemas expressed in this data model is as follows:

- `addRel( $\langle\langle R, k_1, \dots, k_n \rangle\rangle, q$ )` adds to the schema a new relation  $R$  with primary key attribute(s)  $k_1, \dots, k_n$ ,  $n \geq 1$ .

The query  $q$  specifies the set of primary key values in the extent of  $R$  in terms of the already existing schema constructs (this will be a set of  $n$ -tuples).

- `addAtt( $\langle\langle R, a \rangle\rangle, c, q$ )` adds to the schema a non-primary key attribute  $a$  for relation  $R$ . The parameter  $c$  can be either `null` or `nonnull`.

The query  $q$  specifies the extent of the binary relationship between the primary key attribute(s) of  $R$  and this new attribute  $a$  in terms of the already existing schema constructs (this extent will be a set of pairs).

- `delRel( $\langle\langle R, k_1, \dots, k_n \rangle\rangle, q$ )` deletes from the schema the relation  $R$  with primary key attribute(s)

$k_1, \dots, k_n$  (provided all its non-primary key attributes have first been deleted).

The query  $q$  specifies how the set of primary key values in the extent of  $R$  can be restored from the remaining schema constructs.

- $\text{delAtt}(\langle\langle R, a \rangle\rangle, c, q)$  deletes from the schema the non-primary key attribute  $a$  of relation  $R$ .

The query  $q$  specifies how the extent of the binary relationship between the primary key attribute(s) of  $R$  and  $a$  can be restored from the remaining schema constructs.

All primitive transformation rules have a optional additional argument which specifies a constraint on the data that must hold if the transformation is to apply. In the case of our relational model, this can be used to enforce any foreign key constraints that are required to hold. Each of the above primitive transformations,  $t$ , has an automatically derivable reverse transformation,  $\bar{t}$ , defined as follows:

| $t : S_x \rightarrow S_y$                                 | $\bar{t} : S_y \rightarrow S_x$                           |
|---|---|
| $\text{addRel}(\langle\langle R, k_n \rangle\rangle, q)$  | $\text{delRel}(\langle\langle R, k_n \rangle\rangle, q)$  |
| $\text{addAtt}(\langle\langle R, a \rangle\rangle, c, q)$ | $\text{delAtt}(\langle\langle R, a \rangle\rangle, c, q)$ |
| $\text{delRel}(\langle\langle R, k_n \rangle\rangle, q)$  | $\text{addRel}(\langle\langle R, k_n \rangle\rangle, q)$  |
| $\text{delAtt}(\langle\langle R, a \rangle\rangle, c, q)$ | $\text{addAtt}(\langle\langle R, a \rangle\rangle, c, q)$ |

Here, the notation  $\overset{\rightarrow}{k_n}$  is an abbreviation for the sequence of attributes  $k_1, \dots, k_n$ , and we will use similar abbreviations for sequences of attributes in the remainder of the paper.

In [17] we show how this reversibility of schema transformations allows automatic query translation between schemas. In particular, for the simple relational data model above, if a schema  $S$  is transformed to a schema  $S'$  by a single primitive transformation step, the only cases that need to be considered in order to translate a query  $Q$  posed on  $S$  to a query  $Q'$  posed on  $S'$  are if the transformation step is a  $\text{delRel}$  or  $\text{delAtt}$ , in which case occurrences in  $Q$  of the deleted construct need to be substituted by the query  $q$  specified in the transformation. For sequences of primitive transformations, these substitutions are successively applied in order to obtain the final translated query  $Q'$ .

This translation scheme can be applied to each of the constructs of a global schema in order to obtain the derivation of each construct from a set of local schemas. These derivations can then be substituted into any query over the global schema in order to obtain an equivalent query distributed over the local schemas [10], as in the GAV approach.

The above translation scheme can equivalently be applied to each of the constructs of a local schema in order to obtain its derivation from a global schema, as in the LAV approach.

## 2.2 Integrating $S_1$ and $S_2$ into $S_g$

Using the above set of primitive transformations on our simple relational data model, Table 1 lists the transformation steps necessary to integrate the local schemas  $S_1$  and  $S_2$  into the global schema  $S_g$ .

An implicit first step in this integration process is to combine the local schemas into a single schema  $S_1 \cup S_2$  whose constructs are precisely those of  $S_1$  and  $S_2$ . In our particular example, there are no commonly-named relations in  $S_1$ ,  $S_2$  and  $S_g$ , and so for simplicity we just use the relation names directly in Table 1. In general, each construct from a local schema would be tagged with a unique schema identifier, as would the constructs in the target global schema. Note, if required, this implicit step can be formalised in our approach using a series of **extend** steps (see below) to add to  $S_2$  the constructs of  $S_1$ , and to add to  $S_1$  the constructs of  $S_2$ .

Table 1 uses two more primitive transformations to the ones we defined in Section 2.1: **conRel** and **conAtt**. These behave in the same way as **delRel** and **delAtt** except that they indicate that their accompanying query  $q$  may only *partially* restore the extent of the deleted construct; hence we may say that they are **complete** in the same sense as LAV/GAV rules are complete. These two **contract** (abbreviated as **con**) transformations have corresponding reverse transformations **extRel** and **extAtt** respectively, which behave in the same way as **addRel** and **addAtt** except that they indicate their accompanying query  $q$  may only *partially* construct the extent of the new construct. For these **contract** and **extend** (abbreviated as **ext**) transformations, the query  $q$  may be just the distinguished constant *void*, which indicates that there is no information about how to derive the extent of the deleted/new construct from the rest of the schema constructs, even partially.

This distinction between partial and complete derivations allows a precise statement to be made of the relationship between the source and target schemas, since it provides a means of stating how much information is known about a particular schema construct. To illustrate, suppose that we wanted to add a relation **cs(id)** containing all students that study in the computer science department, and that we knew that all students with **degree='G500'** study computer science, as do some students with **degree='GG15'** (and that no other students study computer science):

1.  $\text{addRel}(\langle\langle \text{cs, id} \rangle\rangle, \{x \mid \langle x, y \rangle \in \langle\langle \text{ug, degree} \rangle\rangle \wedge y = \text{'G500'}\})$  would be incorrect, since it states that *cs* contains *just* those students on G500, and no other students.
2.  $\text{extRel}(\langle\langle \text{cs, id} \rangle\rangle, \{x \mid \langle x, y \rangle \in \langle\langle \text{ug, degree} \rangle\rangle \wedge y = \text{'G500'}\})$  would be correct, since it states that *cs* contains those students on G500, plus some others which we have no way of determining.
3.  $\text{extRel}(\langle\langle \text{cs, id} \rangle\rangle, \text{void})$  would be incorrect, since it states that we have no way of determining any of the students that belong to *cs*, when in fact we know that all G500 students belong to *cs*.

If instead it was the case that only students with G500 study computer science, then (1) would be correct. Alternatively, if instead some G500 and some GG15 students study computer science, then (3) would be correct. To make this incompleteness of knowledge explicit in query processing, an annotation *partial* can be attached to a query or sub-query to indicate that it returns a partial result. Thus  $\{x \mid \langle x, y \rangle \in \langle\langle \text{ug, degree} \rangle\rangle \wedge y = \text{'G500'}\}^{\text{partial}}$  means that  $\{x \mid \langle x, y \rangle \in \langle\langle \text{ug, degree} \rangle\rangle \wedge y = \text{'G500'}\}$  returns instances which may be only a subset of the complete answer.

Returning now to Table 1, we see that the transformation steps ① to ⑧ incrementally define the constructs of  $S_g$  from the constructs of  $S_1$  and  $S_2$ . This can be regarded as the GAV aspect of our framework.

Steps ⑨ to ⑳ then incrementally remove the constructs of  $S_1$  and  $S_2$  from this intermediate schema, finally leaving only the constructs of  $S_g$ , as desired. The queries accompanying steps ⑨ to ⑮ show how the constructs of  $S_1$  can be restored from those of  $S_g$ , and the queries accompanying steps ⑯ to ⑳ show how the constructs of  $S_2$  can be restored from  $S_g$ . This can be regarded as the LAV aspect of our framework.

Generally, this is the form that schema transformation/integration takes in our framework, namely a growing phase in which new schema constructs are added, and a shrinking phase in which redundant constructs are deleted.

Note that ⑧ has the optional constraint allowed in our transformation rules, which enforces the foreign keys of the monitors relation in  $S_g$ . Similarly ⑯ has a constraint to enforce the foreign key constraints of the supervises table.

### 3 Correspondence between BAV and GAV/LAV

In this section we examine in more detail the correspondence between BAV and GAV, and between BAV and LAV. We first show how a GAV or LAV definition can be converted into a partial BAV definition. We then show how a complete GAV or LAV definition can be derived from a BAV definition.

BAV thus combines the benefits of both GAV and LAV in the sense that any reasoning or processing which is possible with the view definitions of GAV or LAV will also be possible with the BAV definition.

Given its expressiveness, it could be argued that data integration using the BAV approach is more complex than with GAV or LAV. However, by specifying well-known schema equivalences as BAV transformation macros the production of BAV definitions is greatly simplified, and we conclude this section with a discussion of this issue.

In what follows, we will need to use the fact that a relation  $R$  with primary key attributes  $k_1, \dots, k_n$  and other attributes  $a_1, \dots, a_m$ , has a lossless-join decomposition into  $m$  relations, each containing key attributes  $k_1, \dots, k_n$  and one of the other attributes  $a_1, \dots, a_m$ . Hence we can rewrite any derivation rule of the form

$$R(\vec{k}_n, \vec{a}_m) = \{\vec{x}_n, \vec{y}_m \mid E\}$$

as  $m + 1$  equivalent derivation rules:

$$R(\vec{k}_n) = \{\vec{x}_n \mid E\}$$

$$R(\vec{k}_n, a_1) = \{\vec{x}_n, y_1 \mid E\}$$

⋮

$$R(\vec{k}_n, a_m) = \{\vec{x}_n, y_m \mid E\}$$

(Note that some sub-expressions of the expression  $E$  may be logically redundant for some of the above derivation rules, and can be automatically removed by standard query simplification techniques.) The original derivation rule for the relation  $R$  can thus be represented by the following sequence of primitive transformation steps, where the function  $D(E)$  rewrites the expression  $E$  so that any relations it references are decomposed in the same way as  $R$ :

$$\text{addRel}(\langle\langle R, \vec{k}_n \rangle\rangle, \{\vec{x}_n \mid D(E)\})$$

$$\text{addAtt}(\langle\langle R, a_1 \rangle\rangle, \{\vec{x}_n, y_1 \mid D(E)\})$$

⋮

$$\text{addAtt}(\langle\langle R, a_m \rangle\rangle, \{\vec{x}_n, y_m \mid D(E)\})$$

In what follows, we will refer to this equivalence between derivation rules for relations and sequences of

- ① addRel( $\langle\langle\text{student,id}\rangle\rangle, \{x \mid x \in \langle\langle\text{ug,id}\rangle\rangle \vee x \in \langle\langle\text{phd,id}\rangle\rangle\}$ )
- ② addAtt( $\langle\langle\text{student,name}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{ug,name}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y \rangle \in \langle\langle\text{phd,name}\rangle\rangle\}$ )
- ③ addAtt( $\langle\langle\text{student,left}\rangle\rangle, \text{null}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{ug,left}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y \rangle \in \langle\langle\text{phd,left}\rangle\rangle\}$ )
- ④ addAtt( $\langle\langle\text{student,degree}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{ug,degree}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee x \in \langle\langle\text{phd,id}\rangle\rangle \wedge y = \text{'phd'}\}$ )
- ⑤ addRel( $\langle\langle\text{staff,sno}\rangle\rangle, \{x \mid x \in \langle\langle\text{tutor,sno}\rangle\rangle \vee x \in \langle\langle\text{supervisor,sno}\rangle\rangle\}$ )
- ⑥ addAtt( $\langle\langle\text{staff,sname}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{tutor,sname}\rangle\rangle \wedge x \notin \langle\langle\text{supervisor,sno}\rangle\rangle \vee \langle x, y \rangle \in \langle\langle\text{supervisor,sname}\rangle\rangle\}$ )
- ⑦ addAtt( $\langle\langle\text{staff,dept}\rangle\rangle, \text{null}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{supervisor,dept}\rangle\rangle \vee x \in \langle\langle\text{tutor}\rangle\rangle \wedge x \notin \langle\langle\text{supervisor,sno}\rangle\rangle \wedge y = \text{null}\}$ )
- ⑧ addRel( $\langle\langle\text{monitors,sno,id}\rangle\rangle, \{x, y \mid \langle y, x \rangle \in \langle\langle\text{ug,sno}\rangle\rangle \wedge y \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y \rangle \in \langle\langle\text{supervises,sno,id}\rangle\rangle, \langle x, y \rangle \in \langle\langle\text{monitors,sno,id}\rangle\rangle \rightarrow x \in \langle\langle\text{staff,sno}\rangle\rangle \wedge y \in \langle\langle\text{student,id}\rangle\rangle\}$ )
- ⑨ conAtt( $\langle\langle\text{tutor,sname}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{staff,sname}\rangle\rangle \wedge x \in \langle\langle\text{tutor,sno}\rangle\rangle\}$ )
- ⑩ conRel( $\langle\langle\text{tutor,sno}\rangle\rangle, \{x \mid x \in \langle\langle\text{staff,sno}\rangle\rangle \wedge \langle y, x \rangle \in \langle\langle\text{ug,sno}\rangle\rangle\}$ )
- ⑪ conAtt( $\langle\langle\text{ug,name}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{student,name}\rangle\rangle \wedge x \in \langle\langle\text{ug,id}\rangle\rangle\}$ )
- ⑫ conAtt( $\langle\langle\text{ug,left}\rangle\rangle, \text{null}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{student,left}\rangle\rangle \wedge x \in \langle\langle\text{ug,id}\rangle\rangle\}$ )
- ⑬ conAtt( $\langle\langle\text{ug,degree}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{student,degree}\rangle\rangle \wedge x \in \langle\langle\text{ug,id}\rangle\rangle\}$ )
- ⑭ conAtt( $\langle\langle\text{ug,sno}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle y, x \rangle \in \langle\langle\text{monitors,sno,id}\rangle\rangle \wedge x \in \langle\langle\text{ug,id}\rangle\rangle\}$ )
- ⑮ conRel( $\langle\langle\text{ug,id}\rangle\rangle, \{x \mid x \in \langle\langle\text{student,id}\rangle\rangle \wedge \langle x, \text{'phd'} \rangle \notin \langle\langle\text{student,degree}\rangle\rangle\}$ )
- ⑯ delRel( $\langle\langle\text{supervises,sno,id}\rangle\rangle, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{monitors,sno,id}\rangle\rangle \wedge y \in \langle\langle\text{phd,id}\rangle\rangle, \langle x, y \rangle \in \langle\langle\text{supervises,sno,id}\rangle\rangle \rightarrow x \in \langle\langle\text{supervisor,sno}\rangle\rangle \wedge y \in \langle\langle\text{phd,id}\rangle\rangle\}$ )
- ⑰ delAtt( $\langle\langle\text{supervisor,sname}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{staff,sname}\rangle\rangle \wedge x \in \langle\langle\text{supervisor,sno}\rangle\rangle\}$ )
- ⑱ delAtt( $\langle\langle\text{supervisor,dept}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{staff,dept}\rangle\rangle \wedge x \in \langle\langle\text{supervisor,sno}\rangle\rangle\}$ )
- ⑲ delRel( $\langle\langle\text{supervisor,sno}\rangle\rangle, \{x \mid x \in \langle\langle\text{staff,sno}\rangle\rangle \wedge \langle x, y \rangle \in \langle\langle\text{staff,dept}\rangle\rangle \wedge y \neq \text{null}\}$ )
- ⑳ delAtt( $\langle\langle\text{phd,name}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{student,name}\rangle\rangle \wedge x \in \langle\langle\text{phd,id}\rangle\rangle\}$ )
- ㉑ delAtt( $\langle\langle\text{phd,left}\rangle\rangle, \text{null}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{student,left}\rangle\rangle \wedge x \in \langle\langle\text{phd,id}\rangle\rangle\}$ )
- ㉒ conAtt( $\langle\langle\text{phd,title}\rangle\rangle, \text{nonnull}, \text{void}$ )
- ㉓ delRel( $\langle\langle\text{phd,id}\rangle\rangle, \{x \mid x \in \langle\langle\text{student,id}\rangle\rangle \wedge \langle x, \text{'phd'} \rangle \in \langle\langle\text{student,degree}\rangle\rangle\}$ )

**Table 1. A complete BAV specification of the integration of  $S_1$  and  $S_2$  into  $S_g$**

primitive transformation steps as the **decomposition rule**.

### 3.1 Deriving BAV from GAV

A GAV definition can be used to derive *some* of the information present in the BAV definition, as follows:

First, each GAV derivation rule has the decomposition rule applied to it. Thus, in our running example, rule  $G_1$  generates steps ①–④,  $G_2$  generates ⑧, and  $G_3$  generates ⑤–⑦.

Next, each construct  $c$  of type  $T$  in the source schema is removed using a transformation step of the form  $\text{con}T(c, \text{void})$ . This indicates that the construct is being removed, but (as a consequence of GAV) there is an absence of information about how to restore its extent. In our running example, this means that steps ⑨–㉓ are partially derived as:

⑨ conAtt( $\langle\langle\text{tutor,sname}\rangle\rangle, \text{nonnull}, \text{void}$ )

⋮  
 ㉒ conAtt( $\langle\langle\text{phd,title}\rangle\rangle, \text{nonnull}, \text{void}$ )  
 ㉓ conRel( $\langle\langle\text{phd,id}\rangle\rangle, \text{void}$ )

The box placed around the circle of a transformation step indicates that the same construct is derived as in Table 1, but the query only returns a partial result compared to that in Table 1. For example, ⑨ has a void query, as so clearly returns a partial result compared ⑨. By contrast, ㉒ is completely defined by virtue of the fact that it was defined with a *void* query in the complete BAV integration, and hence no box is placed around the circle.

### 3.2 Deriving BAV from LAV

A LAV definition can also be used to derive *some* of the information present in the BAV definition:

Starting with the set of LAV rules, the decomposition rule is applied to each of them. In our running

example, LAV rules  $L_1$  to  $L_5$  generate  $\textcircled{23}$ – $\textcircled{9}$ , the reverse transformation steps of  $\textcircled{9}$ – $\textcircled{23}$ . However, in this case, *all* the BAV transformation steps generated must be **extend** rather than **add** ones, since with LAV it may be the case that not all the information in a local construct is derivable from the global schema:

$\textcircled{23}$  extRel( $\langle\langle\text{phd,id}\rangle\rangle, \{x \mid x \in \langle\langle\text{student,id}\rangle\rangle \vee \langle x, \text{'phd'} \rangle \in \langle\langle\text{student,degree}\rangle\rangle\}$ )  
 $\vdots$   
 $\textcircled{10}$  extRel( $\langle\langle\text{tutor,sno}\rangle\rangle, \{x \mid x \in \langle\langle\text{staff,sno}\rangle\rangle \wedge \langle y, x \rangle \in \langle\langle\text{ug,sno}\rangle\rangle\}$ )  
 $\textcircled{9}$  extAtt( $\langle\langle\text{tutor,sname}\rangle\rangle, \text{notnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{staff,sname}\rangle\rangle \wedge x \in \langle\langle\text{tutor,sno}\rangle\rangle\}$ )

Note that in this case  $\textcircled{10}$  and  $\textcircled{9}$  remain unboxed, by virtue of the fact that the inversion of contract transformations  $\textcircled{10}$  and  $\textcircled{9}$  are the **extend** transformations shown above with the same queries.

Next, a sequence of BAV steps analogous to  $\textcircled{8}$ – $\textcircled{1}$  is generated as a sequence of **contract** steps, the difference from the BAV specification being that each step is accompanied by a *void* query:

$\textcircled{8}$  conRel( $\langle\langle\text{monitors,sno,id}\rangle\rangle, \text{void}$ )  
 $\vdots$   
 $\textcircled{2}$  conAtt( $\langle\langle\text{student,name}\rangle\rangle, \text{notnull}, \text{void}$ )  
 $\textcircled{1}$  conRel( $\langle\langle\text{student,id}\rangle\rangle, \text{void}$ )

### 3.3 Deriving GAV from BAV

In any BAV transformation which defines a global schema  $S_g$  from local schemas  $S_1, \dots, S_n$ , it is possible to identify the transformation steps that define the constructs of  $S_g$  by taking the subset,  $\mathcal{G}$ , of the **add** and **ext** transformation steps in the total transformation sequence from  $S_1 \cup \dots \cup S_n$  to  $S_g$ .

The GAV derivation rules can then be constructed by taking each **addRel**/**extRel** step in  $\mathcal{G}$ , together with all **addAtt**/**extAtt** steps for the same relation, and applying the decomposition rule in reverse (*i.e.* joining the relations  $\langle\langle R, a_1 \rangle\rangle, \dots, \langle\langle R, a_m \rangle\rangle$  to restore  $R$ ). If the GAV rule is composed from solely **addRel**/**addAtt** BAV rules it will be exact, otherwise it will be complete.

For the example in Table 1, the steps that form  $\mathcal{G}$  are  $\textcircled{1}$ – $\textcircled{8}$ . Taking  $\textcircled{1}$  as one example of an **addRel** step, there are three steps  $\textcircled{2}$ – $\textcircled{4}$  which add attributes to the **student** relation. This automatically gives the following definitions of the constituent parts of **student**:

$\langle\langle\text{student,id}\rangle\rangle = \{x \mid x \in \langle\langle\text{ug,id}\rangle\rangle \vee x \in \langle\langle\text{phd,id}\rangle\rangle\}$   
 $\langle\langle\text{student,name}\rangle\rangle = \{x, y \mid \langle x, y \rangle \in \langle\langle\text{ug,name}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y \rangle \in \langle\langle\text{phd,name}\rangle\rangle\}$   
 $\langle\langle\text{student,left}\rangle\rangle = \{x, y \mid \langle x, y \rangle \in \langle\langle\text{ug,left}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y \rangle \in \langle\langle\text{phd,left}\rangle\rangle\}$   
 $\langle\langle\text{student,degree}\rangle\rangle = \{x, y \mid \langle x, y \rangle \in \langle\langle\text{ug,degree}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee x \in \langle\langle\text{phd,id}\rangle\rangle \wedge y = \text{'phd'}\}$

The **student** relation can now be recomposed using the decomposition rule in reverse, to get a definition of the **student** relation in BAV constructs:

$\text{student}(\text{id,name,left,degree}) = \{x, y_1, y_2, y_3 \mid (x \in \langle\langle\text{ug,id}\rangle\rangle \vee x \in \langle\langle\text{phd,id}\rangle\rangle) \wedge (\langle x, y_1 \rangle \in \langle\langle\text{ug,id}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y_1 \rangle \in \langle\langle\text{phd,name}\rangle\rangle) \wedge (\langle x, y_2 \rangle \in \langle\langle\text{ug,left}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee \langle x, y_2 \rangle \in \langle\langle\text{phd,left}\rangle\rangle) \wedge (\langle x, y_3 \rangle \in \langle\langle\text{ug,degree}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \vee x \in \langle\langle\text{phd,id}\rangle\rangle \wedge y_3 = \text{'phd'})\}$

The **ug** and **phd** relations are then recomposed from the BAV constructs that appear on the RHS of the rule. After some logical simplification, this gives a rule that is equivalent to  $G_1$ :

$\text{student}(\text{id,name,left,degree}) = \{x, y_1, y_2, y_3 \mid \langle x, y_1, y_2, y_3 \rangle \in \text{ug} \wedge \langle x, y_1, y_2, y_3 \rangle \notin \text{phd} \vee \langle x, y_1, y_2, - \rangle \in \text{phd} \wedge y_3 = \text{'phd'}\}$

By a similar process,  $G_2$  can be derived from  $\textcircled{8}$  and  $G_3$  from  $\textcircled{5}$ – $\textcircled{7}$ .

### 3.4 Deriving LAV from BAV

In any BAV transformation which defines a global schema  $S_g$  from local schemas  $S_1, \dots, S_n$ , it is possible to identify the transformation steps that define the constructs of some local schema  $S_i$  by taking the subset,  $\mathcal{L}$ , of the **del** and **con** transformations on constructs of  $S_i$ . Construction of the LAV view definitions from  $\mathcal{L}$  proceeds in a similar fashion to how the GAV view definitions were constructed from  $\mathcal{G}$  in the previous subsection. If the LAV rule is composed from solely **delRel**/**delAtt** BAV rules it will be exact, otherwise it will be complete.

For the example, in Table 1, the steps forming  $\mathcal{L}$  for schema  $S_1$  are  $\textcircled{9}$ – $\textcircled{15}$ . Rule  $L_1$  can then be derived from  $\textcircled{9}$ – $\textcircled{10}$  and rule  $L_2$  from  $\textcircled{11}$ – $\textcircled{15}$ .

### 3.5 Using BAV transformation macros

Although the BAV transformations of Table 1 form a complete definition of the integration of  $S_1$  and  $S_2$  into  $S_g$ , it could be argued that this has been achieved at a greater complexity compared to using GAV or LAV alone. Indeed we have shown how some transformations of the BAV pathway can derive GAV rules, and other distinct transformations derive LAV rules, and thus it could be argued that we are using the two approaches in parallel. There are two arguments against this statement. Firstly, BAV defines the integration with *more* precision than just GAV and BAV in parallel, since BAV differentiates between completely defined and partially defined constructs, rather than exact and complete rules (which in general define many constructs together). Secondly, the complexity may be greatly reduced by specifying well-known schema equivalences as **BAV transformation macros** [24]. We argue that it is possible to specify BAV transformations with as much ease as LAV or GAV derivation rules. The approach involves two phases:

1. Applying well-known schema equivalences from the schema integration literature [1] to make relations from different local databases union-compatible; this includes renaming attributes to be the same in any relations to be unioned.
2. Applying a union operation to integrate the relations into a global schema construct, possibly with a bias to prefer values from a particular local database where there is conflict between the two.

To illustrate, Table 2 shows the integration of  $S_1$  and  $S_2$  into  $S_g$  using this approach, where steps 24–27 are from phase 1 and 28–30 are from phase 2. The schema transformations undertaken by each step are as follows:

- 24 Perform a lossless decomposition of `ug` in  $S_1$  so that the `sno` attribute is held in a separate relation `tutors`. Note that the new `tutors` relation is union-compatible with `supervises` from  $S_2$ .
- 25 Extend `tutor` in  $S_1$  with an attribute `dept`, always with value null (since there is no record of departments in  $S_1$ ). This makes the `tutor` relation union-compatible with `supervisor` from  $S_2$ .
- 26 Contract the `title` attribute from `phd`, since this is not required in the global schema.
- 27 Add an attribute `degree` with value ‘`phd`’ to `phd` in  $S_2$ , making the assumption that all PhD students

are registered for a PhD degree. This transformation combined with 24 and 26 makes the `ug` and `phd` relations union-compatible.

- 28 Define the `student` relation as a union of `ug` and `phd`, but if the same key appears in both relations keep only the attribute values that appear in `phd`. We term this a **right union**.
- 29 Define the `staff` relation as a right union of `tutor` and `supervisor`.
- 30 Define the `monitors` relation as a right union of `tutors` and `supervises`.

The macro used in 24, which implements the decomposition of a table by moving one of its attributes `Att` out into a new table `Rel`, is defined as follows:

```
decompose(Rel,Att)=
  addRel(Rel, {x,y | <y,x> ∈ Att})
  delAtt(Att, {x,y | <y,x> ∈ NewRel})
```

The body of this macro forms a ‘template’ such that `Rel` can be substituted by the scheme of any relation, and `Att` by the scheme of any attribute. In particular, step 24 above expands to:

- 31 `addRel(⟨⟨tutors,sno,id⟩⟩, {x,y | <y,x> ∈ ⟨⟨ug,sno⟩⟩})`
- 32 `delAtt(⟨⟨ug,sno⟩⟩,notnull {x,y | <y,x> ∈ ⟨⟨tutors,sno,id⟩⟩})`

The macro used in 28–30 is defined as follows:

```
rightUnion(RU,R1,C1,R2,C2)=
  addRel(Rel,{x_n^→ | <x_n^→> ∈ R1 ∨ <x_n^→> ∈ R2}
  for i=1 to AttCount(RU)
    addAtt(Att(RU,i),{x_n^→,a_i |
      <x_n^→,a_i> ∈ Att(R1,i) ∧ <x_n^→> ∉ R2
      ∨ <x_n^→,a_i> ∈ Att(R2,i)})
  conRel(R1,{x_n^→ | <x_n^→> ∈ RU ∧ C1}
  for i=1 to AttCount(RU)
    conAtt(Att(R1,i),{x_n^→,a_i |
      <x_n^→,a_i> ∈ Att(RU,i) ∧ <x_n^→> ∈ R1})
  delRel(R2,{x_n^→ | <x_n^→> ∈ RU ∧ C2}
  for i=1 to AttCount(RU)
    delAtt(Att(R2,i),{x_n^→,a_i |
      <x_n^→,a_i> ∈ Att(RU,i) ∧ <x_n^→> ∈ R2})
```

This forms what we term a right union relation, `RU`, between two union-compatible relations `R1` and `R2`, where tuples from `R1` appear only if the key values of that tuple do not appear as a tuple in `R2`. Defined



- ⑭ decompose( $\langle\langle\text{ug,sno}\rangle\rangle, \langle\langle\text{tutors,sno,id}\rangle\rangle$ )
- ⑮ extAtt( $\langle\langle\text{tutor,dept}\rangle\rangle, \text{nonnull}, \{x, y \mid x \in \langle\langle\text{tutor,sno}\rangle\rangle \wedge y = \text{null}\}$ )
- ⑯ conAtt( $\langle\langle\text{phd,title}\rangle\rangle, \text{nonnull}, \text{void}$ )
- ⑰ addAtt( $\langle\langle\text{phd,degree}\rangle\rangle, \text{nonnull}, \{x, y \mid x \in \langle\langle\text{phd,id}\rangle\rangle \wedge y = \text{'phd'}\}$ )
- ⑱ rightUnion( $\langle\langle\text{student,id}\rangle\rangle, \langle\langle\text{ug,id}\rangle\rangle, \langle x, y \in \langle\langle\text{student,degree}\rangle\rangle \wedge y \neq \text{'phd'} \rangle\rangle, \langle\langle\text{phd,id}\rangle\rangle, \langle x, y \in \langle\langle\text{student,degree}\rangle\rangle \wedge y = \text{'phd'} \rangle\rangle$ )
- ⑲ rightUnion( $\langle\langle\text{staff,sno}\rangle\rangle, \langle\langle\text{tutor,sno}\rangle\rangle, \langle x, y \in \langle\langle\text{staff,dept}\rangle\rangle \wedge y = \text{null} \rangle\rangle, \langle\langle\text{supervisor,sno}\rangle\rangle, \langle x, y \in \langle\langle\text{staff,dept}\rangle\rangle \wedge y \neq \text{null} \rangle\rangle$ )
- ⑳ rightUnion( $\langle\langle\text{monitors,sno,id}\rangle\rangle, \langle\langle\text{ug,sno}\rangle\rangle, \langle y, x \in \langle\langle\text{monitors,sno,id}\rangle\rangle \wedge x \notin \langle\langle\text{phd,id}\rangle\rangle \rangle\rangle, \langle\langle\text{supervises,sno,id}\rangle\rangle, \langle x, y \in \langle\langle\text{monitors,sno,id}\rangle\rangle \wedge y \in \langle\langle\text{phd,id}\rangle\rangle \rangle\rangle$ )

**Table 2. A complete BAV macro-based specification of the integration of  $S_1$  and  $S_2$  into  $S_g$**

in the relational algebra,  $\text{RU} = (\text{R1} - (\text{R1} \times \text{R2})) \cup \text{R2}$ , where  $\times$  is a semi-join on the key attributes. The function  $\text{AttCount}(\text{Rel})$  returns the number of non-primary key attributes of  $\text{Rel}$ . The function  $\text{Att}(\text{Rel}, n)$  returns the  $n^{\text{th}}$  attribute of  $\text{Rel}$  (where for the sake of argument alphabetical ordering on their names can be used to sort the attributes). The constraint C1 holds for each tuple of  $\text{RU}$  that can be determined to have originated from  $\text{R1}$ , and the constraint C2 holds for each tuple that can be determined to have originated from  $\text{R2}$ .

The transformation steps of Table 1 can, to a large extent, be generated from these macro definitions. As one might expect, the transformations are not precisely the same since the macro-based approach generates a few extra steps. None-the-less, the two integrations in Table 1 and Table 2 are logically equivalent.

In particular, ⑲ will expand to generate several of the transformation steps in Table 1. The  $\text{addRel}$  statement in the right-union macro will generate ①. Since the attributes of  $\text{student}$  in alphabetical order are  $\text{degree}$ ,  $\text{left}$  and  $\text{name}$ , the first for loop of the macro will generate ④, ③ and ②. The  $\text{conRel}$  statement generates ⑮, and the second for loop generates ⑬, ⑫ and ⑪. The  $\text{delRel}$  statement generates ⑳, and the final for loop generates a step ㉓ not in Table 1:

- ㉓ delAtt( $\langle\langle\text{phd,degree}\rangle\rangle, \text{nonnull}, \{x, y \mid \langle x, y \rangle \in \langle\langle\text{student,degree}\rangle\rangle \wedge x \in \langle\langle\text{phd,id}\rangle\rangle\}$ )

plus two steps ⑳ and ㉑ from Table 1.

Note that step ㉓ ‘cancels out’ step ⑰, in the sense that ⑰ introduced an attribute into the  $\text{phd}$  relation to make it union-compatible with  $\text{ug}$ , and therefore there was an extra attribute to remove from  $\text{phd}$  than was the case in Table 1. Using the  $\text{rightUnion}$  macro on ⑲ and ㉓ will produce similar ‘extra’ transformations to balance ⑮ and ㉑.

We conclude by noting that a  $\text{leftUnion}$  macro would be defined by exchanging  $\text{R1}$  and  $\text{R2}$  in the body of the

$\text{rightUnion}$  macro. An ‘ordinary’ union macro would be defined in a similar manner to  $\text{rightUnion}$ , except for omitting the  $\langle \vec{x}_n \rangle \notin \text{R2}$  constraint, and changing the  $\text{delRel}$  and  $\text{delAtt}$  to  $\text{conRel}$  and  $\text{conAtt}$ .

## 4 Schema Evolution

We next discuss how BAV supports the evolution of both global and local schemas, allowing the transformation pathways and schemas to be incrementally modified as opposed to having to be regenerated.

### 4.1 Evolution of a global schema

Suppose that  $n$  local schemas  $S_1, \dots, S_n$  have been integrated into a global schema  $S_g$ . Let us consider the evolution of the global schema  $S_g$  to  $S'_g$ . The first step is to define this evolution as a schema transformation in our framework. We consider first how to handle evolutions consisting of a *single* primitive transformation.

Suppose that  $S_g$  has evolved to  $S'_g$  by a primitive transformation  $t$ . If  $T^{\text{old}}$  was the transformation from  $S_1 \cup \dots \cup S_n$  to  $S_g$ , the new transformation from  $S_1 \cup \dots \cup S_n$  to  $S'_g$  can automatically be generated by suffixing  $t$  to  $T^{\text{old}}$ :

$$T^{\text{new}} = T^{\text{old}}; t$$

There are three cases to consider for  $t$ :

1. If  $t$  is an  $\text{add}$  or  $\text{del}$  transformation, then  $S'_g$  will be semantically equivalent to  $S_g$ . Any information available from  $S_g$  is also available from  $S'_g$ , and there is nothing further that needs to be done to  $T^{\text{new}}$ .
2. If  $t$  is a  $\text{contract}$  transformation, then there will be some information that used to be present in  $S_g$  that will no longer be available from  $S'_g$ . Thus, it

may be the case that there is now no representation of some local schema construct(s) within  $S'_g$ . Nothing further needs to be done to  $T^{new}$ .

3. If  $t$  is an extend transformation with a *void* accompanying query, then the relationship of the new construct to  $S_1 \dots S_n$  needs to be examined as it may actually be partially or completely derivable from  $S_1 \dots S_n$  by some transformation. This requires domain knowledge e.g. from a human expert or an ontology.

If the new construct is not derivable, then nothing further needs to be done to  $T^{new}$ . If it is derivable, then  $t$ , the last transformation step in  $T^{new}$ , needs to be replaced by a more informative extend or add step (as per our discussion in Section 2.2).

Thus, we see that the first of these two cases is handled totally automatically. The third requires domain knowledge, but is still handled as a simple evolution of the transformation pathway between  $S_1 \cup \dots \cup S_n$  and  $S_g$ .

Evolutions of  $S_g$  consisting of a sequence of primitive transformation steps can be handled by applying the above treatment to each step.

## 4.2 Evolution of a local schema

This is a little more complex since it may require the global schema to be changed. Suppose that some  $S_i$  evolves, to  $S'_i$  say, by a primitive transformation  $t$  (again, we need only consider evolutions consisting of a single primitive transformation, and composite transformations can be handled as a series of primitive transformation steps).

If  $T^{old}$  was the transformation from  $S_1 \cup \dots \cup S_i \cup \dots \cup S_n$  to  $S_g$ , a new transformation from  $S_1 \cup \dots \cup S'_i \cup \dots \cup S_n$  to  $S_g$  can be automatically generated by prefixing the reverse of  $t$  to  $T^{old}$ :

$$T^{new} = \bar{t}; T^{old}$$

Again there are three cases to consider for  $t$ :

1. If  $t$  is an add or del transformation, then the new schema  $S'_i$  will be semantically equivalent to  $S_i$ . Thus any information in  $S_g$  derived from  $S_i$  can now be derived from  $S'_i$  and no changes are required to  $S_g$  or  $T^{new}$ .
2. If  $t$  is a contract transformation, then some information that used to be present in  $S_i$  will no longer be available from  $S'_i$ . In particular, it may now be the case that  $S_g$  contains constructs about which no information is derivable from any local schema

— this can be determined automatically by inspection of  $T^{new}$ . Any such constructs can be removed from  $S_g$ , and the corresponding extend steps from each of the local schemas removed from  $T^{new}$ .

3. If  $t$  is an extend transformation, then the relationship of the new construct to  $S_g$  needs to be examined, as the new construct may actually be derivable from  $S_g$  by some transformation. This requires domain knowledge.

If the new construct is indeed derivable, then the transformation that does so,  $T$  say, is appended to  $T^{new}$ , in order to add the new construct to  $S_g$ . The resulting transformation  $T^{new}; T$  is equal to  $\bar{t}; T^{old}; T$  which just simplifies to  $T^{old}$ .

If the new construct is not derivable, then an extend step is appended to  $T^{new}$ , in order to add the new construct to  $S_g$ . This step matches up with the contract transformation in  $\bar{t}$ . This pair of steps are now removed in order for the new extent of the new construct in  $S'_i$  to be usable by queries posed on the global schema, leaving the original transformation  $T^{old}$ .

Again, the first two of the above cases can be handled totally automatically and the third semi-automatically — [20] gives further details of the second and third cases.

The above treatment also applies in the case that a new schema  $S_{n+1}$  is added to the set of local schemas from which  $S_g$  was derived, since  $S_{n+1}$  can be treated as an empty schema which is successively expanded with new constructs. Similarly, the treatment covers the removal of a local schema  $S_i$  since its constructs can be successively contracted to leave an empty schema.

## 5 Implementation

We are currently implementing the BAV integration approach described here within the ongoing AutoMed project at Birkbeck and Imperial Colleges (see <http://www.doc.ic.ac.uk/automed/>). Figure 2 illustrates the main components of the AutoMed system.

The Model Definitions Tool allows the modelling constructs and primitive transformations of high-level modelling languages to be specified in terms of those of the lower-level hypergraph data model (HDM). These definitions are stored in the Model Definitions Repository (MDR) [2].

The Schemas & Transformations Repository (STR) stores the schemas of data sources, intermediate schemas, global schemas, and the transformation pathways between them [2]. The Schema Transformation &

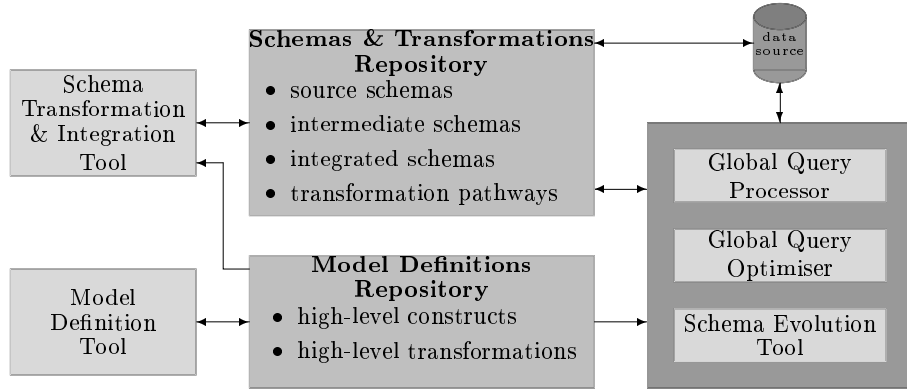


Figure 2. The AutoMed Architecture

Integration Tool allows the creation of new intermediate or global schemas together with the corresponding transformation pathways.

The Global Query Processor undertakes execution of global queries using the schemas and transformation pathways in the STR. Currently we are adopting a GAV approach to query processing [10], but from our earlier discussion in Section 3.4 a LAV approach would also be possible. We are using a functional query language as our intermediate query language (IQL) [21]. Global queries are translated into this IQL and are then reformulated to operate on the data sources, using the GAV views derived from the schema transformation pathways. After optimisation of the global query, subqueries of it are translated into the query languages supported by the data sources, and are submitted to the sources for evaluation. Any further query post-processing is undertaken by the IQL evaluator.

The Schema Evolution Tool supports the evolution of schemas and transformation pathways in the MDR, as described in Section 4 above. These evolved pathways may in some cases be amenable to simplification. For example, a composite transformation  $T; t; T'; \bar{t}; T''$ , where  $T, T', T''$  are arbitrary sequences of primitive transformations,  $t$  is primitive transformation and  $\bar{t}$  is its inverse, can be simplified to  $T; T'; T''$  provided that there are no references within  $T'$  to the construct being renamed, added or deleted by  $t$ . Renaming transformations may also be redundant, and hence removable. For example,  $ren\ c\ c'; del\ c' \equiv del\ c, add\ c'; ren\ c'\ c \equiv add\ c, ren\ c'\ c''; ren\ c''\ c \equiv ren\ c'\ c$ . This kind of simplification is automatically carried out by the Schema Evolution Tool [27]. More generally, these same simplifications can be applied to any schema transformation pathway generated (either automatically or manually) by the Schema Transformation & Integration Tool.

As well as implementing the above components, we are currently investigating techniques for automatic or semi-automatic generation of model definitions in the MDR and transformation pathways in the STR. For the former, we are developing a framework for defining correspondences between different high-level data models. For the latter we are investigating data mining [5] and schema matching approaches [23, 9]. We are also extending the scope of AutoMed to semi-structured data (XML, formatted files) and unstructured text. Finally, we are investigating materialised (as opposed to just virtual) data integration using AutoMed schema transformations, and in particular issues of incremental view maintenance and data lineage tracing [6].

## 6 Conclusions

In this paper we have presented a new approach to data integration which we term **both as view** (BAV). BAV is based on the use of reversible schema transformation sequences. We have shown that LAV and GAV view definitions can be fully derived from BAV transformation sequences, and that BAV transformation sequences can be partially derived from LAV or GAV view definitions.

BAV thus combines the benefits of GAV and LAV in the sense that any reasoning or processing which is possible with GAV or LAV view definitions will also be possible with the BAV specification.

A key advantage of BAV is that it readily supports the evolution of both local and global schemas, allowing transformation pathways and schemas to be incrementally modified as opposed to having to be regenerated.

Our original motivation for developing the BAV approach was as a general framework for schema transformation and integration, which was sufficiently expressive to capture all of the schema transformations

proposed in the literature and used in practical schema integration, as well as having a formal basis. Thus, it may be argued that the transformation pathways resulting from BAV are likely to be more fine-grained, and hence more costly to reason with and process (e.g. in query optimisation and query processing) than the corresponding LAV or GAV view definitions. However, as we discussed in the previous section, BAV transformation pathways are amenable to considerable simplification. Furthermore, standard query optimisation techniques can be applied to the view definitions derived from BAV transformation pathways. Thus, we believe that query processing in AutoMed will not be significantly slower than in other LAV or GAV systems.

We are currently implementing the BAV integration approach within the ongoing AutoMed project, and this is one of the questions to be verified by practical experimentation. Our major application testbed in AutoMed is bioinformatics, and in particular the semantic integration of genomic data sources.

We are currently working on a number of other directions: development of graphical end-user tools for model definition and schema integration; techniques for automatic or semi-automatic generation of model definitions and schema transformation pathways; extension of the approach to support integration of semi-structured and unstructured information; and extension of the approach to materialised data integration.

## References

- [1] C. Batini, M. Lenzerini, and S. Navathe. A comparative analysis of methodologies for database schema integration. *ACM Comput. Surv.*, 18(4):323–364, 1986.
- [2] M. Boyd, P. McBrien, and N. Tong. The automed schema integration repository. In *Proc. BNCOD02, LNCS 2405*, pages 42–45, 2002.
- [3] P. Buneman, S. Naqvi, V. Tannen, and L. Wong. Principles of programming with complex objects and collection types. *Theoretical Computer Science*, 149(1):3–48, 1995.
- [4] S. Chawathe *et al.* The TSIMMIS project: Integration of heterogeneous information sources. In *Proc. 10th Meeting of the Information Processing Society of Japan*, pages 7–18, October 1994.
- [5] A. Doan, P. Domingos, and A. Y. Halevy. Reconciling schemas of disparate data sources: A machine-learning approach. In *ACM SIGMOD*, 2001.
- [6] H. Fan and A. Poulouvasilis. Tracing data lineage using schema transformation pathways. In *Workshop on Knowledge Transformation for the Semantic Web (with ECAI’02)*, 2002.
- [7] L. Fegaras and D. Maier. Towards an effective calculus for object query languages. In *ACM SIGMOD*, pages 47–58, 1995.
- [8] M. Friedman, A. Levy, and T. Millstein. Navigational plans for data integration. In *Proc. 16th National Conf. on AI*, pages 67–73. AAAI Press, 1999.
- [9] M. A. Hernandez, R. J. Miller, and L. M. Haas. Clio: A semi-automatic tool for schema mapping. In *ACM SIGMOD*, 2001.
- [10] E. Jasper. Global query processing in the AutoMed heterogeneous database environment. In *Proc. BNCOD02, LNCS 2405*, pages 46–49, 2002.
- [11] M. Lenzerini. Data integration: A theoretical perspective. In *Proc. PODS02*, pages 247–258, 2002.
- [12] A. Levy. Logic-based techniques in data integration. In J. Minker, editor, *Logic Based Artificial Intelligence*. Kluwer Academic Publishers, 2000.
- [13] A. Levy, A. Mendelzon, Y. Sagiv, and D. Srivastava. Answering queries using views. In *Proc. PODS’95*, pages 95–104. ACM Press, May 1995.
- [14] A. Levy, A. Rajaraman, and J. Ordille. Querying heterogeneous information sources using source description. In *Proc. VLDB’96*, pages 252–262, 1996.
- [15] I. Manolescu, D. Florescu, and D. Kossmann. Answering XML queries on heterogeneous data sources. In *Proc. VLDB’01*, pages 241–250, 2001.
- [16] P. McBrien and A. Poulouvasilis. A formalisation of semantic schema integration. *Information Systems*, 23(5):307–334, 1998.
- [17] P. McBrien and A. Poulouvasilis. Automatic migration and wrapping of database applications — a schema transformation approach. In *Proc. ER’99, LNCS 1728*, pages 96–113, 1999.
- [18] P. McBrien and A. Poulouvasilis. A uniform approach to inter-model transformations. In *Proc. CAiSE’99, LNCS 1626*, pages 333–348, 1999.
- [19] P. McBrien and A. Poulouvasilis. A semantic approach to integrating XML and structured data sources. In *Proc. CAiSE’01, LNCS 2068*, pages 330–345, 2001.
- [20] P. McBrien and A. Poulouvasilis. Schema evolution in heterogeneous database architectures, a schema transformation approach. In *Proc. CAiSE’02, LNCS 2348*, pages 484–499, 2002.
- [21] A. Poulouvasilis. The AutoMed Intermediate Query Language. Technical report, AutoMed Project, 2001.
- [22] A. Poulouvasilis and P. McBrien. A general formal framework for schema transformation. *Data and Knowledge Engineering*, 28(1):47–71, 1998.
- [23] E. Rahm and P. Bernstein. A survey of approaches to automatic schema matching. *VLDB Journal*, 10:334–350, 2001.
- [24] N. Rizopoulos. A database integration tool. Technical report, Imperial College, 2001.
- [25] M. Roth and P. Schwarz. Don’t scrap it, wrap it! A wrapper architecture for data sources. In *Proc. VLDB’97*, pages 266–275, Athens, Greece, 1997.
- [26] M. Templeton, H. Henley, E. Maros, and D. V. Buer. InterViso: Dealing with the complexity of federated database access. *The VLDB Journal*, 4(2):287–317, 1995.
- [27] N. Tong. Database schema transformation optimisation techniques for the AutoMed system. Technical report, AutoMed Project, 2002.