

# Mapping Polymorphism

Ryan Wisnesky  
Harvard University  
ryan@cs.harvard.edu

Mauricio A. Hernández, Lucian Popa  
IBM Research – Almaden  
{mauricio, lucian}@almaden.ibm.com

## ABSTRACT

We examine schema mappings from a type-theoretic perspective and aim to facilitate and formalize the reuse of mappings. Starting with the mapping language of Clio, we present a type-checking algorithm such that typable mappings are necessarily satisfiable. We add type variables to the schema language and present a theory of polymorphism, including a sound and complete type inference algorithm and a semantic notion of a principal type of a mapping. Principal types, which intuitively correspond to the minimum amount of schema structure required by the mappings, have an important application for mapping reuse. Concretely, we show that mappings can be reused, with the same semantics, on any schemas as long as these schemas are expansions (i.e., subtypes) of the principal types.

## 1. INTRODUCTION

Data exchange is the process of transforming data instances of one or more source schemas into instances of a target schema. Much research about data exchange has been done in the context where the process is described using a high-level and declarative formalism called *schema mappings* [25, 27, 11]. Schema mappings (or *mappings* in short) are logical assertions that express constraints between two or more data sources. In particular, schema mappings are used to capture how data conforming to a source schema corresponds to data conforming to a target schema.

The development of schema mappings and their application to data exchange was pioneered by the Clio project (see [10] for a recent retrospective on Clio). There, a great deal of the work concentrated on the problem of generating schema mappings from even higher-level constructs like *correspondences* (or matches) between schema elements, and then converting the generated mapping into queries or programs that capture the transformation semantics of the mappings. In particular, Clio developed a number of query generators to convert mappings into various XML or relational transformation languages (e.g., SQL, XQuery, XSLT, SQL/XML).

Figure 1 shows Clio in action. The source-to-target lines (also called *correspondences*) are entered by the user and indicate how atomic elements in the source relate to atomic elements in the target. The lines on the left- and right-hand-side of the schemas represent the foreign-key constraints that are given with the schemas. Given such schemas with constraints, given a set of correspondences, and also through a fair amount of user interaction, Clio generates a set of schema mapping expressions that best represent the “intention” behind the visual specification of the mapping. In this example, the five correspondences shown in the figure would be translated into one mapping expression that requires that each record in the source set `gradEnroll` is split over three target sets: `eval`, `Student` and `course`. (In the figure, the elements that have the `[0,*]` suffix denote sets.)

One important and desirable feature of the generated mapping is that it preserves data associations. For the above example, this means that the individual values of the input record (i.e., `sid`, `name`, `cid`, `grade`, `file`) will remain connected in the target (by exploiting the foreign key constraints and the schema structure). We shall give the exact mapping expression in Section 2, in a language that follows the source-to-target tuple-generating dependencies of [11].

IBM [19], Microsoft [6], Oracle [7], and others are building an ecosystem of tools around schema mappings, where mappings are building blocks for more complex data transformations. Models and semantics of schema mappings for data exchange [11], operations over mappings (e.g., composition [23, 12] and inversion [9, 2]), and the application of such operations to metadata management [3, 5] have been extensively studied within the database and information integration community over the past decade.

However, the problem of **reusing schema mappings** has been largely ignored. Schema mappings are, in general, neither parametric nor modular. That is, mappings are static expressions over a source and a target schema and can not be easily reused over other schemas. Large schemas often combine and reuse existing type definitions to define larger and more complex types. For example, often a *Person* type is defined and used whenever the attributes of *Person* are needed in the schema definition. Current mapping languages cannot take advantage of this modularity and cannot express a mapping between expansions of these types (without knowing the concrete expansions). Ideally, we would prefer to create generic mappings between source and target types

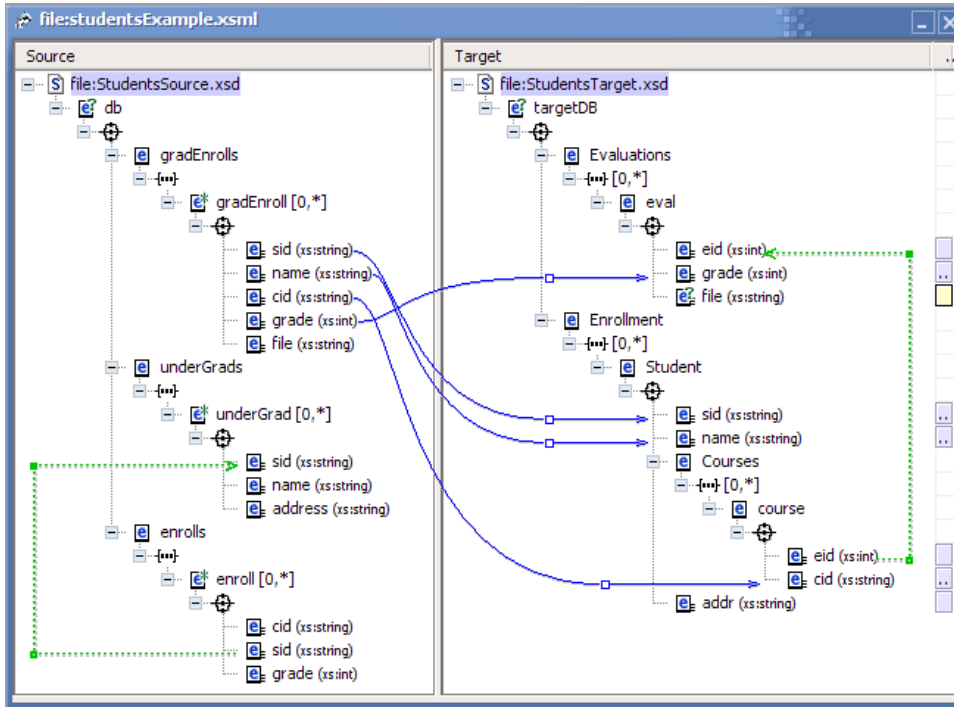


Figure 1: A schema mapping in Clío

(e.g., between the source *Person* type and a target *Student* type) and then reuse those mapping within larger mappings. Similarly to regular programming languages, users should be able to build libraries of reusable mappings between common data types. Mappings in these libraries can be reused on larger data transformation scenarios, as building blocks for more complex mappings.

This paper studies how to apply programming language techniques to mapping languages to formally defined when and how mapping are reusable. We propose *polymorphism* as a lightweight formal technique for reuse. We describe a mapping type and use standard type-checking techniques to determine if a mapping applies on a certain context. We further show that mappings have a natural semantic notion of *principal type*, corresponding intuitively to the minimum amount of schema structure required by the mappings. We then show that principal types can be soundly and completely inferred from the structure of the logical assertions alone (i.e. independently of the underlying schemas). In turn, this enables the following mapping reuse technique. Given a mapping expression  $M$ , which originally is defined from a concrete source schema  $S$  to a concrete target schema  $T$ , we first infer a principal source schema  $S'$  and a principal target schema  $T'$ . We can then reuse the same mapping  $M$  on any source schema that is an expansion (i.e., subtype) of  $S'$  and any target schema that is an expansion (i.e., subtype) of  $T'$ . Moreover, we show in a precise sense that the semantics of  $M$  remains invariant during reuse.

This paper is organized as follows. We begin by reviewing schema mappings in Section 2. We then formally define nested schemas and nested data in Section 3. The core mapping language that we focus in this paper is formally defined

in Section 4. We then formalize a type-checking algorithm in Section 5. We add type variables to the schema language and present a theory of polymorphism, including a sound and complete type inference algorithm in Section 6, and a semantic notion of principal type in Section 7. All the proofs will be available in the full version of this paper, which will appear online.

## 2. SCHEMA MAPPINGS: PRELIMINARIES

Several mapping languages have been proposed over the years. In the case of relational schemas, a language popular in the data exchange and data integration community is that of *source-to-target tuple generating dependencies (s-t tgds)* [11] or (*Global-and-Local-As-View*) *GLAV mappings* [14, 21]. Essentially, these are formulas of the form

$$\forall \mathbf{x}(\phi(\mathbf{x}) \rightarrow \exists \mathbf{y} \psi(\mathbf{x}, \mathbf{y}))$$

where  $\phi$  is a conjunction of relational atoms over the source schema, and  $\psi$  is a conjunction of relational atoms over the target schema.

The language of s-t tgds was extended to handle hierarchical (XML) schemas in [27]. There, the extended language must account for the fact that we can have relations nested inside relations. However, the  $\forall\exists$  shape of the dependencies remains the same. An orthogonal extension has been introduced in [15], where in addition to the nesting in the data, we may also have nesting in the mapping formulas themselves (i.e., a  $\forall\exists$  type of mapping as above can appear as one of the atoms of  $\psi$  in another  $\forall\exists$  mapping). It is this type of more flexible, nested mappings that are in use in the Clío system and that we target in this paper.

Although we are using Clío's mapping language [15] as a

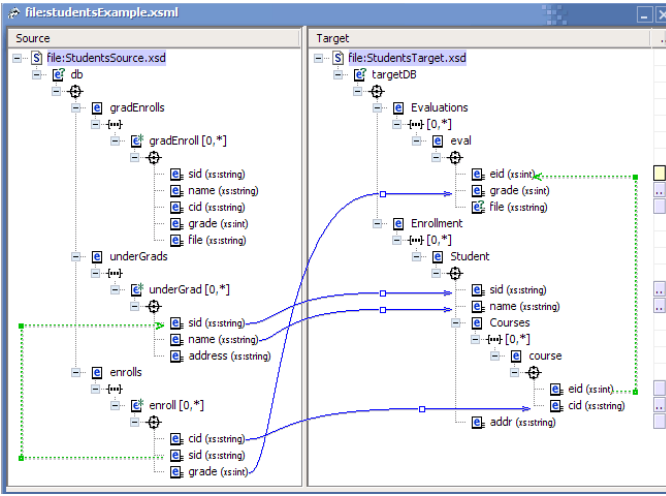


Figure 2: A different mapping

starting point<sup>1</sup>, for most of our results it will be convenient to use a simpler and more uniform mapping syntax, which we translate Clio mappings into. Furthermore, with an eye toward enabling interoperability between mapping systems and conventional programming languages, and unlike [27], we will explicitly represent nested relational data using simple algebraic datatypes.

In this section we will first informally describe the full language of Clio nested mappings after which we move on to the formal definitions of our data model and mapping expressions.

## 2.1 The Clio Mapping Language: Overview

The Clio mapping language uses a for - where  $\Rightarrow$  exists - where syntax. Intuitively, the for clause binds variables to tuples in the source, and the first (optional) where clause describes the source constraints to be satisfied by these source tuples (e.g., these constraints express filters or join conditions). The exists clause describes the tuples that are expected to exist in the target, and the second where clause describes the target constraints to be satisfied by the target tuples as well as the content of these target tuples in terms of the source tuples.

For our example in Figure 1, the mapping that was informally discussed in Section 1 can be written as:

$$\begin{aligned} & \text{for } g \text{ in } db.gradEnrolls.gradEnroll \\ & \Rightarrow \text{exists } s \text{ in } targetDB.Enrollment.Student, \\ & \quad c \text{ in } s.Courses.course, \\ & \quad e \text{ in } targetDB.Evaluations.eval \\ & \text{where } s.sid = g.sid \wedge s.name = g.name \wedge \\ & \quad c.cid = g.cid \wedge e.grade = g.grade \wedge \\ & \quad e.eid = c.eid \end{aligned}$$

Note that there may be dependencies between variables, re-

<sup>1</sup>We must actually omit one feature of [15]: grouping conditions, which use set-valued equalities. This is done for expediency, and it is certainly plausible that our results can be extended.

flecting the nesting in the data. For example, the binding  $c$  in  $s.Courses.course$  reflects the fact that  $c$  must be an element in the set  $course$  that is nested under  $s.Courses$ , while  $s$  itself is an element of the top-level set  $Student$ . In general, there are typing rules that dictate whether a mapping is well-typed with respect to a schema. We shall visit these typing rules in detail, after we formally define schemas.

The most direct semantics of such a mapping is that of a *constraint* between the source schema and the target schema. In this view, a mapping  $M$  defines a set of pairs of instances  $(I, J)$  with  $I$  over the source and  $J$  over the target such that  $(I, J)$  satisfies the constraint  $M$  (in the standard sense). For every  $I$  we call a  $J$  such that  $(I, J)$  satisfies  $M$  a *solution* for  $I$  with respect to  $M$  [11].

For our example, note that the where clause contains two types of equality conditions. The last equality condition is a target-target condition reflecting a constraint (present on the target schema), while the other four conditions are target-source conditions reflecting the correspondences between source elements and target elements.

Finally, note that not all fields in the schema are mentioned in the mapping (e.g., the target addr field). Such fields may be required to appear in the resulting instance. Synthetizing null values for such fields as well as for the fields that are mentioned by the mapping but are not constrained by the source (e.g., the two target fields eid) is the role of the data exchange process that implements the mapping. (See [11] for a canonical implementation via the chase that constructs *universal solutions*.)

Mapping expressions may be recursively nested inside the second where clause of another mapping expression. Such nesting, in general, is orthogonal to the nesting of the data. As an example of nested mappings, consider the Clio mapping depicted in Figure 2. The four correspondences in Figure 2 maps the source-side undergraduate information into two several target sets. Clio compiles those correspondence into the following nested mapping:

$$\begin{aligned} & \text{for } u \text{ in } db.underGrads.underGrad \\ & \Rightarrow \text{exists } s \text{ in } targetDB.Enrollment.Student \\ & \quad \text{where } s.sid = u.sid \wedge s.name = u.name \wedge \\ & \quad (\text{for } e \text{ in } db.enrolls.enroll \\ & \quad \quad \text{where } e.sid = u.sid \\ & \quad \quad \Rightarrow \text{exists } c \text{ in } s.Courses.course, \\ & \quad \quad \quad e' \text{ in } targetDB.Evaluations.eval \\ & \quad \quad \quad \text{where } c.cid = e.cid \wedge e'.grade = e.grade \wedge \\ & \quad \quad \quad e'.eid = c.eid ) \end{aligned}$$

Here, source tuples in the underGrad set are mapped into target tuples in the Student set by the outer mapping. The inner mapping requires that all the associated source enrollment tuples (determined by the join condition  $e.sid = u.sid$ ) are mapped into corresponding target tuples under course and eval. Note that the inner mapping is a constraint that must be satisfied for every binding of the variables in the outer mapping (i.e.,  $u$  and  $s$ ).

Even with the help of tools, the full-fledged development of

such mappings can be quite complex and can involve significant user effort. Being able to reuse mappings on similar schemas is a crucial feature for real-life metadata applications. In the subsequent part of the paper, we shall explore types and, in particular, polymorphism to show when and how can mappings be reused.

### 3. NESTED RELATIONAL MODEL

In this section we define the nested relational schema and the nested relational data over which our mappings operate. We define schema and data independently of each other, and then relate the two via a typing relation that can be used to type the data with a schema.

#### 3.1 Nested Relational Schema

The mappings that we consider operate over data instances whose shape can be described by *nested-relational (NR) schemas* [27]. NR schemas describe atomic types, (unordered) records, (unordered) sets of records, and (unordered) sets of choices (or variants).

DEFINITION 1 (NESTED RELATIONAL (NR) SCHEMA).

$$\begin{aligned} \text{Row} &::= \emptyset \mid \langle \mathcal{L} : \text{Schema}, \text{Row} \rangle \\ \text{Schema} &::= \text{ATOMIC } \mathcal{A} \mid \text{RCD } \text{Row} \mid \\ &\quad \text{SETRCD } \text{Row} \mid \text{SETCHC } \text{Row} \end{aligned}$$

Intuitively, a *Row* is the building block for either a record or a choice. A *Row* is a tuple of *label : type* pairs, where *label* is drawn from an infinite set  $\mathcal{L}$  of label names, while *type* is one of the four main types under *Schema*. We require that *Rows* contain only one instance of any label name, and we equate *Rows* that are equivalent up to permutation of record label name and schema pairs<sup>2</sup>. We shall often abbreviate  $\langle l_1 : t_1, \langle l_2 : t_2, \emptyset \rangle \rangle$  as  $\langle l_1 : t_1, l_2 : t_2 \rangle$ .

ATOMIC  $\mathcal{A}$  denotes an atomic type, where  $\mathcal{A}$  stands for any concrete base type (i.e., **Int**, **String**, etc.). For convenience, we will denote ATOMIC **Int** simply as **Int**.

We make several simplifying assumptions in the above definition of complex types. As a result, we depart somewhat from a completely general nested model, although, in practice, we can still capture most XML and relational schemas. Concretely, we do not allow choice values unless they are immediately under a set type. Hence, we explicitly “package” a set of choices into the SETCHC construct. Another restriction is that we disallow sets of sets or sets of atomic types. As a result, we only allow either sets of records or sets of choices, which are “packaged” as SETRCD or SETCHC. Note, however, that the resulting definition allows RCD, SETRCD, and SETCHC to be freely nested within each other.

We illustrate the above definition with the following two NR schemas, denoted as *src* and *dst* (these names also play the role of the *roots* of the schemas). Here, *src* describes a set of students records, each with a nested set of choices

<sup>2</sup>These restrictions are not captured in the syntax, and require special treatment during type inference.

reflecting the teaching/enrollment status of a student, while *dst* describes a set of employee records.

```
src, RCD ( ( students : SETRCD (
            fullname : String,
            status : SETCHC ( ( teaching : String,
                              taking : String ) ) ) ) )
dst, RCD ( ( employees : SETRCD ( ( name : String,
                                   job : String,
                                   id : Int ) ) ) )
```

Notice that each **status** record contains either a **teaching** or a **taking** element. Intuitively, if we assume **students** are represented in XML, the equivalent DTD representation is:

```
<!ELEMENT students (fullname , status)*>
<!ELEMENT status (teaching | taking)*> ...
<!ELEMENT employees (name, job, id)*> ...
```

#### 3.2 Nested Relational Data

Although there are many ways to represent nested relations we adopt here particular representation that is based on algebraic datatypes. The results of this paper are independent of the way that nested relational data are represented, but the primary advantage to using this representation is that it has both a simple set-theoretic definition and a simple definition using algebraic datatypes. As such, this representation can be used to exchange nested relational data between mapping systems and programming languages.

DEFINITION 2 (DATA INSTANCES). An instance is constructed inductively as one of the following:

- an atom (e.g. 1 or “IBM”),
- a pair  $(l : d)$  of a label  $l$  and an instance  $d$
- a set  $\{d_1, \dots, d_n\}$  of instances

To illustrate, the following is a valid instance, representing intuitively a set of person tuples, where each tuple is a set of *label : value* pairs.

```
{ { (name : John Doe), (age : 25) },
  { (firstname : Alice), (lastname : May), (age : 22) } }
```

An instance has no type a priori, and in general it could be typed with multiple types (schemas) or it may have no type at all. Concretely, a set of pairs could be typed as RCD, but it could also be typed as SETCHC. For this example, it is intuitive that the inner set  $\{(name : John Doe), (age : 25)\}$  should be typed as RCD, while the outer set should be typed as SETCHC.

We next define, inductively, the typing rules that describe how data can be associated with a schema. In effect, these rules define all the valid data instances for each schema construct. We use  $\mathcal{B}(A)$  to represent the domain of an atomic type  $A$ . Moreover, we use  $\llbracket X \rrbracket$  to denote the set of all data instances conforming to schema construct  $X$ .

DEFINITION 3 (TYPING DATA).

$$\frac{}{\emptyset \in \llbracket \text{SETCHC } \emptyset \rrbracket}$$

$$\frac{\forall d \in D, d \in \{(l : i) \mid i \in \llbracket t \rrbracket\} \cup \llbracket \text{SETCHC } r \rrbracket}{D \in \llbracket \text{SETCHC } (l : t, r) \rrbracket} \quad \frac{}{\emptyset \in \llbracket \text{RCD } \emptyset \rrbracket}$$

$$\frac{d \in \llbracket t \rrbracket \quad e \in \llbracket \text{RCD } r \rrbracket}{\{(l : d)\} \cup e \in \llbracket \text{RCD } (l : t, r) \rrbracket} \quad \frac{d \in \mathcal{B}(A)}{d \in \llbracket \text{ATOMIC } A \rrbracket}$$

$$\frac{\forall d \in D, d \in \llbracket \text{RCD } r \rrbracket}{D \in \llbracket \text{SETRCD } r \rrbracket}$$

For example, the RCD rule can be informally read as follows: the instance formed by adding a pair  $(l : d)$  to  $e$  is of type  $\text{RCD}(l : t, r)$ , if  $d$  is an instance of type  $t$  and  $e$  is an instance record defined by  $r$ . Similarly, the SETCHC rule says that a set of instances  $D$  are of type  $\text{SETCHC } (l : t, r)$  if for all instances  $d \in D$ ,  $d$  is either the pair  $(l : i)$  where  $i$  is of type  $t$ , or  $d$  is a pair in  $\text{SETCHC } r$ . An implicit assumption in the RCD and SETCHC rules is that the label  $l$  does not occur in  $r$  (otherwise  $(l : t, r)$  is not a well-formed Row). In effect, these two rules specify how to construct sets of “larger” choices (or “larger” records) from “smaller” data instances. The other interesting rule is the SETRCD rule which specifies how to construct sets of records from records.

For our two NR schemas defined earlier, the following are valid nested relational data instances (the first for  $src$ , and the last for  $dst$ ):

```
src, {(students : {
  {(fullname : John Doe),
   (status : {
     (teaching : CS100),
     (taking : CS200),
     (teaching : CS101)})
  },
  {(fullname : Mary Jane),
   (status : {
     (taking : CS100),
     (taking : CS200)})
  })
})}

dst, {(employees : {
  {(name : John Doe), (job : CS100), (id : 1)}
  {(name : John Doe), (job : CS101), (id : 2)}
})}
```

As a final example, suppose we have  $\{a_1, a_2\} = \llbracket \text{ATOMIC } A \rrbracket$  and  $\{b\} = \llbracket \text{ATOMIC } B \rrbracket$ . Then our typing relation gives the following:

- $\llbracket \text{RCD } (i : \text{ATOMIC } B, j : \text{ATOMIC } A) \rrbracket$  has two instances:
  1.  $\{(i : b), (j : a_1)\}$
  2.  $\{(i : b), (j : a_2)\}$
- $\llbracket \text{SETRCD } (i : \text{ATOMIC } B, j : \text{ATOMIC } A) \rrbracket$  has four instances:
  1.  $\{\{(i : b), (j : a_1)\}, \{(i : b), (j : a_2)\}\}$
  2.  $\{\{(i : b), (j : a_1)\}\}$
  3.  $\{\{(i : b), (j : a_2)\}\}$
  4.  $\{\}$

- $\llbracket \text{SETCHC } (i : \text{ATOMIC } B, j : \text{ATOMIC } A) \rrbracket$  has 8 instances:

1.  $\{(i : b), (j : a_1), (j : a_2)\}$
2.  $\{(i : b), (j : a_1)\}$
3.  $\{(i : b), (j : a_2)\}$
4.  $\{(i : b)\}$
5.  $\{(j : a_1), (j : a_2)\}$
6.  $\{(j : a_1)\}$
7.  $\{(j : a_2)\}$
8.  $\{\}$

### Connections to Algebraic Datatypes

The above definition of the nested relational data (with choices) differs from those traditionally found in database literature. However, this representation makes it easy to encode nested relational data using the standard algebraic datatypes  $0$ ,  $1$ ,  $+$ ,  $\times$  and a labeled pair construct. If nothing else, this encoding can be used to transfer data between nested relational systems and programming languages. We can break out a separate notion of powerset ( $\mathcal{P}$ ) and choice (CHC), which lets us write data in an equivalent way as:

$$\begin{aligned} \llbracket \text{SETCHC } r \rrbracket &= \mathcal{P}(\llbracket \text{CHC } r \rrbracket) \\ \llbracket \text{SETRCD } r \rrbracket &= \mathcal{P}(\llbracket \text{RCD } r \rrbracket) \\ \llbracket \text{CHC } \emptyset \rrbracket &= 0 \\ \llbracket \text{CHC } (l : t, r) \rrbracket &= \{(l : i) \mid i \in \llbracket t \rrbracket\} + \llbracket \text{CHC } r \rrbracket \\ \llbracket \text{RCD } \emptyset \rrbracket &= 1 \\ \llbracket \text{RCD } (l : t, r) \rrbracket &= \{(l : i) \mid i \in \llbracket t \rrbracket\} \times \llbracket \text{RCD } r \rrbracket \end{aligned}$$

This definition has translations into the functional programming languages ML and Haskell.

## 4. CORE MAPPINGS

For most of the results in this paper we do not need to work directly with Clio mapping expressions. Instead, we work with a simpler, more uniform mapping syntax that simplifies the presentation. We call this simplified language the *core* language. We note that this is not the surface syntax that a user will see, but rather an internal, more basic representation.

The main difference between the core and the Clio mapping language described in Section 2 is that, in the core language, we ignore the *sidedness* (and the quantification) of the variables. In the Clio language, each variable is marked with a side, either source or target, and correspondingly quantified (as universal or existential, respectively). This sidedness can be easily inferred by taking the convention that all variables descended from the source schema root(s) are universally quantified and all variables descended from the target schema root(s) are existentially quantified.

The syntax of core mapping expressions is given by the following definition.

DEFINITION 4 (CORE MAPPING EXPRESSIONS).

$$\begin{aligned} \text{Path} &::= v \mid \text{Path}.l \\ \text{NmE} &::= \text{TRUE} \mid \text{Path EQ Path} \mid \text{NmE AND NmE} \mid \\ &\quad v \text{ IN Path. NmE} \mid v \text{ OF } l \text{ FROM Path. NmE} \end{aligned}$$

In the above, a *Path* is an expression that navigates inside a record. A *Path* is always of the form  $v.l_1. \dots .l_n$ , where  $v$  is a variable (or one of the schema roots) and  $l_1, \dots, l_n$  are labels. Nested mapping expressions (*NmE*) are then formed via two forms of variables binders, one used to navigate inside sets of records and the other used to navigate inside sets of choices. Each bound variable can then be used in a nested mapping subexpression, which in turn can be either another binder, or an equality or a conjunction (possibly empty, or TRUE) of other mapping expressions.

For the  $v$  IN *Path* form, *Path* must resolve into a set of records (SETRCD) and  $v$  will bind to records of that set. For the  $v$  OF  $l$  FROM *Path* form, *Path* must resolve into a set of choice elements (SETCHC); this construct automatically selects only labeled pairs of the form  $l : d$ , and then  $v$  will bind to the data values  $d$ . (Recall that, in the surface syntax, each such binding will be either universal or existential.)

To illustrate, consider the following Clio mapping between the *src* and *dst* schemas defined in Section 3.1:

```
(m)  for s in src.students,
      t of teaching FROM s.status
      =>
      exists e in dst.employees
      where e.name = s.fullname & e.job = t
```

We write the above Clio mapping as the following core mapping expression:

```
s IN src.students .
t OF teaching FROM s.status .
e IN dst.employees .
e. name = s.fullname AND e.job = t
```

We note that the above definition ignores, for simplicity, any distinction between the types of equality conditions that may appear in a Clio mapping. In particular, in a Clio mapping, source-to-source equalities always appear on the left-hand side of an implication ( $\Rightarrow$ ), since they play the role of selection or join conditions on the source side. In order to keep the correspondence to Clio mappings, we simply take the convention that all source-to-source equalities are implicitly assumed, in the core mapping syntax, to be on the left-hand side of such source-to-target implications. To illustrate, consider the source-to-source equality  $e.sid = u.sid$  in the nested mapping of Section 2.1. In the core mapping syntax, this equality will be put together (via AND) with the subsequent encoding of the exists clause. However, when translating back to the Clio mapping, we would know to replace AND with  $\Rightarrow$ , since the equality is source-to-source. In general, it is relatively simple to translate each core mapping expression, unilaterally, back into a Clio mapping.

A mapping expression can range over multiple source schema roots and target schema roots, which will appear as free variables. The association between a root and its schema (type) is captured by a *context*  $\Gamma$ , which in general is a finite map of bindings from variables to Schema:

DEFINITION 5 (CONTEXT).

$$\Gamma ::= - \mid (v, \text{Schema}); \Gamma$$

A *mapping* is a core mapping expression together with a context. (In the full Clio language, a mapping is actually a set of mapping expressions that share a context.)

DEFINITION 6 (MAPPING). A (not necessarily well-typed) *mapping* is an association  $(\Gamma, m)$  between a core mapping expression  $m$  and a context  $\Gamma$  that contains exactly the free variables of  $m$ .

**Satisfaction** A mapping  $M$  can be given a meaning as a constraint between a set of source and target data instances, where one instance is associated with each schema root (whether source or target) of  $M$ . Formally, we define an *environment* to be an association from schema roots  $v$  to data instances  $I$ , as follows

DEFINITION 7 (ENVIRONMENT).

$$\Delta ::= - \mid (v, I); \Delta$$

There is a natural correspondence between contexts and environments. We will write  $\Delta \in \llbracket \Gamma \rrbracket$  to indicate that each binding  $(v, t) \in \Gamma$  has a corresponding binding  $(v, I) \in \Delta$  such that  $I \in \llbracket t \rrbracket$ .

Satisfaction, written  $\Delta \models M$ , is a relation between environments  $\Delta$  and mappings  $M$ , and means that the constraints expressed by  $M$  are true when interpreted in the structure  $\Delta$ . Satisfaction is in general “untyped”: it is independent of any notion of schema and it may apply to ill-typed mappings and to data that is not an instance of any schema. Much of the utility of type-checking, which we address next, comes from carving out a subset of mapping expressions (the well-typed ones) that are well-behaved (i.e., are satisfiable, see Theorem 1) over data instances that conform to the nested relational model.

## 5. TYPE CHECKING

One of the goals of our type system is to ensure that well-typed mappings are satisfiable. Our typing relation  $\vdash$  is between a core mapping and a context. We give the typing relation by means of inductive inference rules. First, we need rules for typing a path, which we indicate with  $::$  and define below.

DEFINITION 8 (TYPE-CHECKING PATHS).

$$\frac{\text{VAR}}{(v, t) \in \Gamma} \Gamma \vdash v :: t \qquad \frac{\text{RCD-ELIM}}{\Gamma \vdash p :: \text{RCD}(l : t, r)} \Gamma \vdash p.l :: t$$

With this in hand, the main type-checking relation is:

DEFINITION 9 (TYPE-CHECKING).

$$\begin{array}{c}
\text{WF-EQ} \\
\frac{\Gamma \vdash p :: \text{ATOMIC } a \quad \Gamma \vdash p' :: \text{ATOMIC } a}{\Gamma \vdash p \text{ EQ } p'} \quad \text{WF-TRUE} \\
\frac{}{\Gamma \vdash \text{TRUE}} \\
\\
\text{WF-AND} \\
\frac{\Gamma \vdash m \quad \Gamma \vdash m'}{\Gamma \vdash m \text{ AND } m'} \\
\\
\text{SETRCD-ELIM} \\
\frac{\Gamma \vdash p :: \text{SETRCD } r \quad (v, \text{RCD } r); \Gamma \vdash m}{\Gamma \vdash v \text{ IN } p. m} \\
\\
\text{SETCHC-ELIM} \\
\frac{\Gamma \vdash p :: \text{SETCHC } (l : t, r) \quad (v, t); \Gamma \vdash m}{\Gamma \vdash v \text{ OF } l \text{ FROM } p. m}
\end{array}$$

The typing rules are syntax directed and are easily read bottom up. For instance, the rule WF-EQ says that to check if an equality constraint is well-formed, we must check if each of the paths is well formed and, moreover, that the two paths have the same atomic type. Checking well-formedness of the paths, in turn, requires repeated uses of RCD-ELIM to check if the required projections exist.

The two more complex rules are SETRCD-ELIM and SETCHC-ELIM. For SETRCD-ELIM, to check that  $v \text{ IN } p. m$  is well-formed with respect to a context  $\Gamma$ , we must perform two things. First, we must verify that  $p$  types to  $\text{SETRCD } r$  for some Row  $r$ . Then we must check, recursively, that  $m$  is well-formed in a new context where  $\Gamma$  is extended with the pair  $(v, \text{RCD}r)$ . The SETCHC-ELIM rule is somewhat similar, and involves the additional check that the label  $l$  must be one of the valid choices in the  $\text{SETCHC}$  type for  $p$ .

The following theorem shows that the typing relation achieves our goal of satisfiability.

THEOREM 1 (SATISFIABILITY). *Suppose  $\Gamma \vdash M$ . Then we can compute an environment  $\Delta \in \llbracket \Gamma \rrbracket$  such that  $\Delta \models M$ .*

**Proof Sketch.** We follow the general technique of constructing canonical databases (or canonical models). While there is some analogy with satisfiability of conjunctive queries [1], our proof is more involved, since the mapping language is more complex, involves nested relational data, and uses both universal and existential quantification.

We construct a canonical model by initializing first an environment where all sets are empty. This environment is then enriched based on the typing derivation tree of the mapping, by adding an element for each universally quantified variable. Each such element (either a record or a choice) is minimally constructed in the sense that we populate its required structure while leaving empty any of its nested sets. Atomic type fields (whenever they appear) are initialized with fresh new nulls or symbols.

This environment is then enriched with elements for the existentially quantified variables, via a process similar to the chase with source-to-target tgds that constructs a canonical

universal solution [11]. Equalities are used to identify the symbolic values of atomic fields. Since there are no equalities between constants in the language, and the instances are entirely symbolic, this process never fails. In the end, we replace all variables with distinct constants and obtain a satisfying environment.  $\square$

We note that, for this general form of nested mappings, the above theorem cannot be strengthened to say that we can always find target solutions with respect to  $M$  when given an *arbitrary* set of instances over the source schema roots. This is in contrast to the simpler language of s-t tgds in [11] which always admit solutions. It is also in contrast to the more restricted class of nested mappings in [15], which include a complex syntactic check to always guarantee the existence of solutions. However, such syntactic check can always be added as an orthogonal ingredient on top of our typing system.

The following example shows why our nested mappings may not always have solutions (although they are always satisfiable). (For ease of presentation, we revert here to the surface syntax of Clio mappings.)

```

for s in src.students
⇒
  exists e in dst.employees
  where e.name = s.fullname AND
    ( for t of teaching FROM s.status
    ⇒
      e.job = t)

```

If we look at our earlier *src* instance in Section 3.2, it is easy to see that we cannot construct a target solution, since we would have to construct an employee whose name is *John Doe* and whose job is equal to both *CS100* and *CS101*. Nevertheless, the mapping is satisfiable, since we can always pick some other *src* instance (e.g., with one *teaching* element in the *status* set) for which there is a solution.

Finally, we note that we can permit more expressive mappings at the cost of weakening the semantic guarantees provided by the type system. For instance, naively adding atomic valued constants results in typeable mappings that contain unsatisfiable constraints (e.g.  $1 = 2$ ). Similarly, we can add  $\text{CHC}$  types that are not guarded by a set type (i.e., not packaged as  $\text{SETCHC}$ ) to give rise to typeable mappings that are unsatisfiable.

## 6. POLYMORPHISM

Our typing relation allows for a typable  $M$  to have distinct  $\Gamma$  such that  $\Gamma \vdash M$ . In other words, there can be different schemas for which  $M$  is valid. In this section we extend the schema language with type variables to obtain a formalism for expressing the *principal typings* [30] of mapping expressions, where principal typings completely capture the contexts for which a mapping type-checks. A principal typing corresponds intuitively to the minimum amount of structure that is needed by a mapping to type-check.

## 6.1 Polymorphic Schema

To begin, we extend the schema language to include row variables  $\rho$ , schema variables  $\sigma$ , and atomic variables  $\alpha$  in place of atomic type names. We collectively call these variables *type variables*.

DEFINITION 10 (POLYMORPHIC SCHEMA).

$$\begin{aligned} \text{Row} &::= \langle \rangle \mid \langle \text{Row}, \mathcal{L} : \text{Schema} \rangle \mid \rho \\ \text{Schema} &::= \text{ATOMIC } \alpha \mid \text{RCD Row} \mid \\ &\quad \text{SETRCD Row} \mid \text{CHC Row} \mid \sigma \end{aligned}$$

We shall often apply substitutions to polymorphic schemas; a substitution  $\phi$  maps type variables to schema constructs of the same sort (i.e., a row variable can be mapped to a *Row*, a schema variable can be mapped to a *Schema*, while an atomic variable can be mapped to a concrete atomic type name or another atomic variable). Polymorphic schemas are not closed under arbitrary substitutions in the sense that row variables cannot be substituted arbitrarily. For instance, in  $\langle \rho, l : t \rangle$ ,  $\rho$  can only range over rows that do not include  $l$ . When we write substitutions we assume that the result is well-formed.

## 6.2 Principal Typings

First, note that our earlier typing rules (in Section 5) may be used without modification with polymorphic schemas. Our notion of principal typing applies *only* to polymorphic schema, however.

DEFINITION 11 (PRINCIPAL TYPING).  $\Gamma$  is a principal typing for  $M$  if (1) for every substitution  $\phi$ , we have that  $\phi\Gamma \vdash M$ , and (2) for every  $\Gamma'$  such that  $\Gamma' \vdash M$ , there is some  $\phi$  such that  $\phi\Gamma = \Gamma'$ .

Thus, the contexts for which a mapping  $M$  type-checks are exactly those that can be obtained (via substitution) from the principal typings.

The following is a principal typing for the earlier schema mapping  $m$  in Section 4. Here,  $\rho_1, \dots, \rho_5$  are distinct row variables and  $\alpha_1, \alpha_2$  are distinct atomic variables.

$$\begin{aligned} \text{src}, \text{RCD} \langle \rho_1, \text{students} : \text{SETRCD} \langle \rho_2, \\ \quad \text{fullname} : \text{ATOMIC } \alpha_1, \text{status} : \text{SETCHC} \langle \rho_3, \\ \quad \quad \text{teaching} : \text{ATOMIC } \alpha_2 \rangle \rangle \rangle \\ \text{dst}, \text{RCD} \langle \rho_3, \text{employees} : \text{SETRCD} \langle \rho_5, \text{name} : \text{ATOMIC } \alpha_1, \\ \quad \quad \text{job} : \text{ATOMIC } \alpha_2 \rangle \rangle \end{aligned}$$

Note that the `taking` and `id` fields are absent, intuitively because they are not mentioned in the mapping expression. However, by applying a substitution that sends  $\rho_3$  to a row containing a `taking` label, and sends  $\rho_5$  to a row containing an `id` label, we obtain polymorphic schemas that correspond to the original schemas *src* and *dst* of Section 3.1.

PROPOSITION 1. *Principal types are unique up to  $\alpha$ -equivalence (that is, one-to-one renaming of type variables).*

## 6.3 Type Inference

In this section we give a sound and complete type inference algorithm that computes the principal typing of a core mapping expression, or fails if one does not exist. Note that the input to our algorithm is a set of mapping expressions without a context (i.e., a set of “schema-less” mappings).

The general approach of the inference algorithm is to use iterated unification, extended to account for permutation of rows. Intuitively, during inference we need to unify row expressions like  $\langle l_1 : t_1, l_2 : t_2 \rangle$  and  $\langle l_2 : t_2, l_1 : t_1 \rangle$  but traditional unification distinguishes these permutations and cannot unify them. A detailed discussion of such extended unification algorithm can be found in [17]. We give next our unification algorithm as the set of rules in Definition 12.

These rules define a reflexive and symmetric relationship  $x \stackrel{\phi}{\sim} x'$  between type expressions  $x$  and  $x'$ , where  $x \stackrel{\phi}{\sim} x'$  means that  $x$  and  $x'$  are equivalent under the substitution  $\phi$ . (The reflexivity and symmetry rules are not shown explicitly in the definition, but nevertheless assumed.) The substitution  $\phi$  is synthesized (inductively) by the rules.

DEFINITION 12 (SCHEMA UNIFICATION).

$$\begin{array}{c} \text{BIND} \quad \frac{v \notin fv(x)}{v \stackrel{\phi}{\sim} x} \\ \text{APPLY} \quad \frac{x \stackrel{\phi}{\sim} x'}{Cx \stackrel{\phi}{\sim} Cx'} \\ \text{ROW} \quad \frac{(l : t) \stackrel{\phi}{\in} r' \quad \phi(r) \stackrel{\psi}{\sim} \phi(r') - l}{(l : t, r) \stackrel{\psi\phi}{\sim} r'} \quad \text{INVAR} \quad \frac{r' \text{ fresh} \quad r \notin fv(t)}{(l : t) \stackrel{r \mapsto (l:t, r')}{\in} r} \\ \text{INHEAD} \quad \frac{t \stackrel{\phi}{\sim} t'}{(l : t) \stackrel{\phi}{\in} (l : t', r)} \quad \text{INTAIL} \quad \frac{(l : t) \stackrel{\phi}{\in} r \quad l' \neq l}{(l : t) \stackrel{\phi}{\in} (l' : t', r)} \end{array}$$

In the above, BIND and APPLY are typical unification rules. In BIND, for instance,  $x$  is a type or row expression,  $v$  is a type variable, and the  $v \notin fv(x)$  represents the “occurs check” that  $v$  must not occur among the free variables of  $x$ . If the premise is satisfied, then  $v$  is equivalent to  $x$  under a substitution that maps  $v$  to  $x$ . In APPLY, the notation  $C$  is used to mean one of ATOMIC, SETRCD, RCD, and SETCHC.

The rules ROW, INVAR, INTAIL, and INHEAD are needed for row unification. They define and use an additional *inserter* substitution  $\phi$ , of  $(l : t)$  into  $r$ , written  $(l : t) \stackrel{\phi}{\in} r$ , if  $(l : \phi(t)) \in \phi(r)$ . The  $-$  operator removes a label from a row. The notation  $\psi\phi$  represents the composition of substitutions (i.e., apply  $\phi$  first and then apply  $\psi$ ).

Unification may generate row expressions with duplicate labels, and the inference algorithm must explicitly check for this. For brevity we have omitted these checks in the rules.

PROPOSITION 2. *Schema unification produces most general unifiers.*

Based on schema unification, we are now ready to define the type inference algorithm. As with typechecking, the rules for type inference are syntax directed. Substitutions are synthesized (returned), and contexts are inherited (passed as arguments when we go up the rules). We use  $\Vdash$  to indicate inference. We begin by performing type inference on paths:

DEFINITION 13 (TYPE INFERENCE FOR PATHS).

$$\begin{array}{c} \text{VAR-INF} \\ \frac{(v, t) \in \Gamma}{\Gamma \Vdash v :: t} \end{array} \quad \begin{array}{c} \text{RCD-ELIM-INF} \\ \frac{\phi\Gamma \Vdash p :: t \quad \text{RCD } (\langle l : \sigma, \rho \rangle \overset{\psi}{\sim} t \quad \sigma, \rho \text{ fresh}}{\psi\phi\Gamma \Vdash p.l :: \psi\sigma} \end{array}$$

We explain the second, more complex rule. The input is an initial context  $\Gamma$  and an expression  $p.l$ . We first infer that  $p$  has type  $t$  (under some substitution  $\phi$ ). We then verify that  $t$  can be written equivalently (via some other substitution  $\psi$ ) as  $\text{RCD } (\langle l : \sigma, \rho \rangle)$ , for some type variable  $\sigma$  and row variable  $\rho$ . Here we use the earlier unification algorithm. If this verification succeeds, then  $p.l$  has type  $\psi\sigma$ , under a new context obtained from  $\Gamma$  by applying both substitutions  $\psi$  and  $\phi$ . Note that in this rule all we need to infer about the type of  $p$  is that it is a record that contains the label  $l$ .

The complete inference algorithm for core mapping expressions is given by the following set of rules.

DEFINITION 14 (TYPE INFERENCE FOR MAPPING EXPS).

$$\begin{array}{c} \text{WF-EQ-INF} \\ \frac{\phi_1\Gamma \Vdash p :: t \quad \phi_2\phi_1\Gamma \Vdash p' :: t' \quad \phi_2 t \overset{\phi_3}{\sim} t' \quad \alpha \text{ fresh} \quad \phi_3 t' \overset{\phi_4}{\sim} \text{ATOMIC } \alpha}{\phi_{4..1}\Gamma \Vdash p \text{ EQ } p'} \end{array} \quad \begin{array}{c} \text{WF-TRUE-INF} \\ \frac{}{\Gamma \Vdash \text{TRUE}} \end{array}$$

$$\begin{array}{c} \text{WF-AND-INF} \\ \frac{\phi_1\Gamma \Vdash m \quad \phi_2\phi_1\Gamma \Vdash m'}{\phi_2\phi_1\Gamma \Vdash m \text{ AND } m'} \end{array}$$

SETRCD-ELIM-INF

$$\frac{\phi_1\Gamma \Vdash p :: t \quad \text{SETRCD } \rho \overset{\phi_2}{\sim} t \quad \rho \text{ fresh} \quad \phi_3\phi_2(\langle v, \text{RCD } \rho \rangle; \phi_1\Gamma) \Vdash m}{\phi_3\phi_2\phi_1\Gamma \Vdash v \text{ IN } p. m}$$

SETCHC-ELIM-INF

$$\frac{\phi_1\Gamma \Vdash p :: t \quad \text{SETCHC } (\langle l : \sigma, \rho \rangle \overset{\phi_2}{\sim} t \quad \sigma, \rho \text{ fresh} \quad \phi_3\phi_2(\langle v, \sigma \rangle; \phi_1\Gamma) \Vdash m}{\phi_3\phi_2\phi_1\Gamma \Vdash v \text{ OF } l \text{ FROM } p. m}$$

To give an idea of how the rules work, consider the SETRCD-ELIM-INF rule. We are given an initial (partially inferred) context  $\Gamma$  and a core mapping expression  $v \text{ IN } p. m$ . First, we infer the type  $t$  for  $p$  (under some substitution  $\phi_1$ ). We then check that this type  $t$  unifies with a  $\text{SETRCD } \rho$  type, for some row variable  $\rho$  (and under another substitution  $\phi_2$ ). We then extend the context  $\phi_1\Gamma$  with a new pair that binds  $v$  to  $\text{RCD } \rho$  (under the substitution  $\phi_2$ ). We then pass the new context and the mapping expression  $m$  to a recursive call to the type inference algorithm. In return, we will obtain a new context and substitution  $\phi_3$ .

In practice, having this algorithm means that given  $M$ , we can compute a principal typing  $\Gamma$ , or fail exactly when  $M$  is not typable. The following two main theorems of this section capture this precisely.

THEOREM 2 (SOUNDNESS). *For all  $\phi\Gamma \Vdash M, \varphi\Gamma \vdash M$ .*

THEOREM 3 (COMPLETENESS). *For all  $\varphi\Gamma \vdash M$ , there exists  $S$  and  $s$  such that  $S\Gamma \Vdash M$  and  $\varphi = s \circ S$ .*

These properties and the inference algorithm extend straightforwardly to additional operations that have types describable using the system of qualified types in [17]. For instance, atomic constants and function symbols are easy to add, and so is an “erase present field  $l$ ” operation. (Given a record with fields  $(\langle l : t, r \rangle)$ , the “erase present field  $l$ ” operation would return a record with fields  $r$ ). However, an “erase field  $l$  if it is present, otherwise do nothing” operation cannot be added as it cannot be typed in the language of [17].

One motivation for using this particular type inference algorithm (and this particular choice for the row unification algorithm, which is in the spirit of the algorithm in [17]) is that the resulting system is compatible with modern functional programming languages like Haskell. This compatibility means that mapping expressions embedded in languages like Haskell may have their principal typings inferred “for free” – surely a win for mapping reusability.

## 7. POLYMORPHISM AND SEMANTICS

In this section we investigate the meaning of polymorphism and give a semantic notion of principal type. We begin by making precise the notion that schema structure can be unnecessary for a mapping, and show how instances can have corresponding unnecessary data removed in a satisfiability preserving way. We conclude by giving a condition that guarantees a mapping semantics is indifferent to unnecessary structure and data.

### 7.1 Subtyping

To structurally compare schemas we define a subtyping relation ( $\leq, \preceq$ ) as the reflexive transitive closure of the following rules.

DEFINITION 15 (NR SCHEMA SUBTYPING).

$$\begin{array}{c} \text{WIDTH} \\ \frac{}{\langle l : t, r \rangle \prec r} \end{array} \quad \begin{array}{c} \text{DEPTH} \\ \frac{t' < t}{\langle l : t', r \rangle \prec \langle l : t, r \rangle} \end{array}$$

$$\begin{array}{c} \text{SUB-SETCHC} \\ \frac{r' \prec r}{\text{SETCHC } r' < \text{SETCHC } r} \end{array} \quad \begin{array}{c} \text{SUB-SETREC} \\ \frac{r' \prec r}{\text{SETRCD } r' < \text{SETRCD } r} \end{array}$$

$$\begin{array}{c} \text{SUB-REC} \\ \frac{r' \prec r}{\text{RCD } r' < \text{RCD } r} \end{array}$$

Since NR schema definition is mutually inductive (in defines both *Row* and *Schema*), our definition of subtyping contains

both rules for rows ( $\preceq$ ) and for schema ( $\leq$ ). This definition can be used with concrete NR schema, or polymorphic NR schema. Subtyping lifts to contexts and mappings pointwise and respects typability:

**THEOREM 4 (MAPPING SUBTYPING).**  $\Gamma \vdash M$  and  $\Gamma' \leq \Gamma$  implies  $\Gamma' \vdash M$ .

## 7.2 Erasure

We might hope that  $X \leq Y$  implies  $\llbracket X \rrbracket \subseteq \llbracket Y \rrbracket$ . However, this fails for our schema semantics: for example,  $\text{RCD}(\ell : \text{int}) \leq \text{RCD}(\emptyset)$ , but a data instance for the first record type is never a data instance for the second record type. In the extreme, we could achieve this property by changing the schema semantics to be more inclusive (so that, for instance,  $\{\ell : 0\} \in \llbracket \text{RCD}(\emptyset) \rrbracket$ ): define  $\llbracket X \rrbracket_{\leq} = \bigcup_{X' \leq X} \llbracket X' \rrbracket$ . From a mapping perspective this is a radical and unintuitive departure from the nested relational model, so we will instead relate subtyping to a different, more intuitive, semantic operation of erasure. Principal typings then correspond to spaces of instances that are “maximally” erased.

If  $\Gamma$  and  $\Gamma'$  are two contexts such that  $\Gamma' \leq \Gamma$ , we can define an operation  $\text{erase}(\Gamma' \leq \Gamma) : \llbracket \Gamma' \rrbracket \rightarrow \llbracket \Gamma \rrbracket$  over a derivation of  $\Gamma' \leq \Gamma$  that removes data from the instances in  $\llbracket \Gamma' \rrbracket$  so that they become instances in  $\llbracket \Gamma \rrbracket$ . The definition of  $\text{erase}$  in Figure 3 is meant to apply pointwise to the instances.

The  $\text{erase}$  operation (and auxiliary operation  $\text{erase}'$ ) is similar to a projection operator applied from a subtype to a supertype. As a very simple example,  $\text{erase}(\text{SETRCD}(A : t_1, B : t_2) \leq \text{SETRCD}(A : t_1))$  has the same effect as the standard relational algebra projection of a relation with the  $A, B$  attributes to a relation with just the  $A$  attribute. But  $\text{erase}$  is more general – in programming language parlance, it is a *subtyping coercion*.

Erasure is directed by the derivation of  $X' \leq X$ , and in general there may be distinct derivations. For instance,  $\text{erase}'$  is non-deterministic because the first two  $\text{erase}'$  rules may both apply to a given  $r' \preceq r$ . However, the following proposition holds.

**PROPOSITION 3.** For any two derivations  $a, b$  of  $X' \leq X$ , for any  $(v, x) \in X'$ ,  $\text{erase}(a)(x) = \text{erase}(b)(x)$ .

For this reason, we will treat  $\text{erase}$  as being a function parameterized by the subtyping relations itself and not by subtyping derivations.

The following theorem, which is of importance for mapping reuse, states that  $\text{erase}$  removes unnecessary data in a way that preserve the semantics of a mapping.

**THEOREM 5.** Suppose  $\Gamma \vdash M$  and  $\Gamma' \leq \Gamma$  and  $\Delta' \in \llbracket \Gamma' \rrbracket$ . Then  $\Delta' \models M$  if and only if  $\text{erase}(\Gamma' \leq \Gamma)(\Delta') \models M$ .

As an example, consider the context  $\Gamma'$  given by the  $\text{src}$  and  $\text{dst}$  schemas in Section 3.1. This context is a subtype of the following context  $\Gamma$ :

```
src, RCD(⟦students : SETRCD(⟦fullname : String,
                             status : SETCHC(⟦teaching : String⟦)⟦)⟦)
dst, RCD(⟦employees : SETRCD(⟦name : String, job : String⟦)⟦)
```

Furthermore, consider now the  $\text{src}$  and  $\text{dst}$  instances given in Section 3.2. These instances form an environment  $\Delta' \in \llbracket \Gamma' \rrbracket$ , which, moreover, satisfies the mapping  $m$  given in Section 4 (i.e.,  $\Delta' \models m$ ). If we apply  $\text{erase}(\Gamma' \leq \Gamma)(\Delta')$ , we obtain the following pair of instances:

```
src, {(students : {(fullname : John Doe,
                    (status : {(teaching : CS100,
                                teaching : CS101)})),
                    {(fullname : Mary Jane), (status : {}}))}}
dst, {(employees : {(name : John Doe), (job : CS100)},
                    {(name : John Doe), (job : CS101)}})}
```

It is then immediate to see that  $\text{erase}(\Gamma' \leq \Gamma)(\Delta') \models m$ . In other words the satisfaction of  $m$  is preserved when we move between instances of  $\Gamma'$  and instances of  $\Gamma$  via  $\text{erase}$ . Of course, the intuition behind this preservation is that  $m$  does not use any of the “extra” fields in  $\Gamma'$  (i.e., the fields taking and id).

Principal typings are defined using polymorphic schemas, but we do not have a semantics for polymorphic schemas. However, we can *concretize* principal typings by using a canonical substitution to remove row, schema, and atomic variables. This canonical substitution takes row variables to the empty row, atomic variables to arbitrary (distinct) atomic types, and schema variables to the empty record. Principal typings that have been concretized in such a way denote spaces of instances that cannot be further erased (while still preserving the semantics of mappings). For instance,  $\Gamma$  above is a concretized principal typing of  $m$ . A mapping will not type-check with respect to any schema that has less structure than a concretized principal type, and erasure may not be satisfiability preserving for supertypes of the concretized principal type.

## 7.3 Parametricity

We now relate the above notion of preservation of satisfaction that is based on subtyping and erasure with the more general notion of parametricity. In general, for an arbitrary semantics, we have no guarantees that whenever  $\Gamma \vdash M$  and  $\Gamma' \vdash M$  the meaning of the mapping  $(M, \Gamma)$  is related to the meaning of the mapping  $(M, \Gamma')$ . In contrast, with a parametric semantics (defined below), the meaning of a mapping depends only on the mapping expression, and not on the context with respect to which it type-checks. We show in this section that such a parametric semantics for mappings does exist.

**DEFINITION 17 (PARAMETRICITY).** A mapping meaning function  $\llbracket \cdot \rrbracket$  is parametric if  $\Gamma \vdash M$  and  $\Gamma' \vdash M$  imply  $\llbracket (M, \Gamma') \rrbracket = \llbracket (M, \Gamma) \rrbracket$

In general, not all semantics are parametric. Taking the meaning of a mapping to be the query that a system like Clio generates to implement the mapping (see [19]) usually

DEFINITION 16 (ERASURE).

$$\begin{aligned}
\text{erase}'(\llbracket l : t', r \rrbracket \preceq r)(x) &= \{ (l' : v) \mid (l' : v) \in x \wedge l' \neq l \} \\
\text{erase}'(\llbracket l : t', r \rrbracket \preceq \llbracket l : t, r \rrbracket)(x) &= \{ (l' : \text{erase}(t' \leq t)(y)) \mid (l' : y) \in x \wedge l' = l \} \cup \{ (l' : y) \mid (l' : y) \in x \wedge l' \neq l \} \\
\text{erase}'(r' \preceq z \preceq r)(x) &= \text{erase}(z \preceq r)(\text{erase}(r' \preceq z)(x)) \\
\text{erase}'(r = r)(x) &= x \\
\text{erase}(\text{RCD } r' \leq \text{RCD } r)(x) &= \text{erase}'(r' \preceq r)(x) \\
\text{erase}(\text{SETCHC } r' \leq \text{SETCHC } r)(x) &= \text{erase}'(r' \preceq r)(x) \\
\text{erase}(\text{SETRCD } r' \leq \text{SETRCD } r)(x) &= \{ \text{erase}(\text{RCD } r' \prec \text{RCD } r)(y) \mid y \in x \} \\
\text{erase}(t' \leq z \leq t)(x) &= \text{erase}(z \leq t)(\text{erase}(t' \leq z)(x)) \\
\text{erase}(t = t)(x) &= x
\end{aligned}$$

Figure 3: Erasure

results in a non-parametric semantics, because fields that do not appear in the mapping expressions may still need to be mentioned in the query (typically, the query must explicitly set these fields to NULL). Likewise, a standard satisfaction-based semantics,

$$\llbracket (M, \Gamma) \rrbracket = \{ \Delta \mid \Delta \in \llbracket \Gamma \rrbracket \wedge \Delta \models M \}$$

is not parametric because as the schemas in  $\Gamma$  vary, so do the spaces of instances. However, concretized principal types are unique, so we can give a parametric semantics by taking:

$$\llbracket M \rrbracket = \{ \Delta \mid \Delta \in \llbracket \hat{\Gamma} \rrbracket \wedge \Delta \models M \}$$

which, by our earlier Theorem 5, is equivalent to a semantics of erased solutions:

$$\llbracket M \rrbracket = \bigcup_{\Gamma} \{ \text{erase}(\Gamma \leq \hat{\Gamma})(\Delta) \mid \Delta \in \llbracket \Gamma \rrbracket \wedge \Gamma \leq \hat{\Gamma} \wedge \Delta \models M \}$$

(It is immediate that the above semantics is parametric, because no matter the choice of  $\Gamma$  for which  $\Gamma \vdash M$ , the meaning of  $M$  is given in terms of  $\hat{\Gamma}$  and therefore invariant.)

A parametric semantics such as above allows a mapping to be applied at different schemas (as long as they are subtypes of a concretized principal typing) with the same meaning. A possible scenario for mapping reuse is one in which a developer creates a mapping  $M$  from  $\Gamma_{\text{SRC}}$  to  $\Gamma_{\text{DST}}$ . Then the system infers a tighter schema, namely,  $\hat{\Gamma}$ , which gets rid of all the unused parts. Then  $M$  can be automatically applied to other contexts  $\Gamma \leq \hat{\Gamma}$ . Furthermore, the meaning is the same:  $M$  will be satisfied by a set of instances that is the same modulo erasure (i.e., the set of erased instances with respect to  $\hat{\Gamma}$  is the same).

## 8. RELATED WORK

Nested relational data can be represented using trees, and the programming languages community has a wealth of knowledge about transformations on tree-like data [13], including bi-directional tree transformations [18]. The XML processing languages and systems XDuce [20] and Xtatic [16] aim to create general XML processing languages where XML values are first-class. The languages are functional in nature and have an intuitive semantics for XML processing. They introduce a rich language of types to describe XML values (including regular expressions). The specificity of types for XML, however, leads to restrictions on polymorphism, function types, and type inference that have only recently been addressed [29]. These systems are, in a certain sense, the

XML counterparts of LINQ [24]. There is also a fair amount of work on schema inference for SQL and relational algebra (e.g. [28, 8]).

The idea of reusing parts of mappings to construct other mappings has appeared sparingly in the data exchange literature. The work in [22] describes how to use previously computed correspondences between schema elements (a.k.a. *schema matchings*) to enhance the discovery of new schema matchings. In our work, we try to match and reuse *entire* schema mappings that encode source and target schema constraints and a (potentially large) number of correspondences. More recently, [26] explores mapping reuse as part of their *schema exchange* framework. In that work, mappings are expressed between meta-schema models called *schema templates*. Given a mapping between two schema templates and a source schema that is an “instance” of the source schema template, their framework computes a new target schema that is an “instance” of the target template and derives a schema mapping between them. Our work does not depend on creating or maintaining these higher-level mappings between templates and we assume the target schema is given, not created as part of the data exchange process. Further, [26] only reports on relational schemas while our work consider nested-relational schemas.

Finally, one could argue that a different mapping reuse strategy is possible if mappings can be composed as described in [12, 4]. Given a mapping from schema A into schema B and another mapping from schema B to schema C, mapping composition creates (under the correct conditions) a mapping from A to C. In effect, we are “reusing” two existing mappings to create a new mapping between A and C. Our work can be used in this scenario to help find the existing intermediate mappings. For instance, if we are trying to create a mapping between schemas A and C and we already have a mapping from A to B, we can use the techniques in this paper to find and coerce an existing mapping into a mapping between B and C.

## 9. CONCLUSION AND FUTURE WORK

We have implemented our ideas as an extension to Clio. The extension is able to infer types of mappings, reuse mappings at different schema, and can automatically populate mapping graphs through schema-analysis. The extension also has experimental support for features not described in this paper, including reuse in the presence of target-side foreign key

constraints, the ability to rewrite a mapping from a schema to apply it to a larger (in a schema containment sense) schema, and support for a mapping language extension that is able to express mappings that depend on other mappings. The latter feature has applications, for instance, in dataflow graphs of mappings when a mapping downstream depends on a mapping upstream.

As a future direction, we are studying how to support recursive schema languages. Some schema languages used in data exchange allow recursion (e.g., XML schemas) but the mapping and schema languages defined in this paper do not have syntax for recursion. Clio deals with recursive XML schema by unfolding them a set number of times while translating them into NR schema. We are currently working on extending the mapping and NR schema languages to handle recursion and investigating reuse with the resulting languages.

## 10. REFERENCES

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison Wesley Publishing Co, 1995.
- [2] M. Arenas, J. Pérez, and C. Riveros. The Recovery of a Schema Mapping: Bringing Exchanged Data Back. In *PODS*, pages 13–22, 2008.
- [3] P. A. Bernstein. Applying Model Management to Classical Meta Data Problems. In *CIDR*, pages 209–220, 2003.
- [4] P. A. Bernstein, T. J. Green, S. Melnik, and A. Nash. Implementing Mapping Composition. In *VLDB*, pages 55–66, 2006.
- [5] P. A. Bernstein and S. Melnik. Model Management 2.0: Manipulating Richer Mappings. In *SIGMOD*, pages 1–12, 2007.
- [6] P. A. Bernstein, S. Melnik, and J. E. Churchill. Incremental Schema Matching. In *VLDB (demo)*, pages 1167–1170, 2006.
- [7] V. Borkar, M. Carey, D. Engovatov, D. Lychagin, T. Westmann, and W. Wong. XQSE: An XQuery Scripting Extension for the AquaLogic Data Services Platform. In *ICDE*, pages 1307–1316, 2008.
- [8] J. V. den Bussche, D. V. Gucht, and S. Vansummeren. A crash course on database queries. In *PODS*, pages 143–154, 2007.
- [9] R. Fagin. Inverting schema mappings. *ACM TODS*, 32(4), 2007.
- [10] R. Fagin, L. M. Haas, M. A. Hernández, R. J. Miller, L. Popa, and Y. Velegrakis. Clio: Schema Mapping Creation and Data Exchange. In *Conceptual Modeling: Foundations and Applications, Essays in Honor of John Mylopoulos*, pages 198–236. Springer, 2009.
- [11] R. Fagin, P. G. Kolaitis, R. J. Miller, and L. Popa. Data exchange: semantics and query answering. *Theor. Comput. Sci.*, 336(1):89–124, 2005.
- [12] R. Fagin, P. G. Kolaitis, L. Popa, and W. Tan. Composing Schema Mappings: Second-Order Dependencies to the Rescue. *TODS*, 30(4):994–1055, 2005.
- [13] J. N. Foster, B. C. Pierce, and A. Schmitt. A logic your typechecker can count on: Unordered tree types in practice. In *PLAN-X, informal proceedings*, Jan. 2007.
- [14] M. Friedman, A. Y. Levy, and T. D. Millstein. Navigational Plans For Data Integration. In *AAAI/IAAI*, pages 67–73, 1999.
- [15] A. Fuxman, M. A. Hernández, H. Ho, R. J. Miller, P. Papotti, and L. Popa. Nested Mappings: Schema Mapping Reloaded. In *VLDB*, pages 67–78, 2006.
- [16] V. Gapeyev, M. Y. Levin, B. C. Pierce, and A. Schmitt. The Xtatic experience. In *PLAN-X*, Jan. 2005. University of Pennsylvania Technical Report MS-CIS-04-24, Oct 2004.
- [17] B. R. Gaster and M. P. Jones. A polymorphic type system for extensible records and variants. Technical Report NOTTCS-TR-96-3, Department of Computer Science, University of Nottingham, November 1996.
- [18] M. Greenwald, J. Moore, B. Pierce, and A. Schmitt. A language for bi-directional tree transformations. Manuscript. [www.cis.upenn.edu/~bcpierce/papers/lenses.pdf](http://www.cis.upenn.edu/~bcpierce/papers/lenses.pdf), 2003.
- [19] L. M. Haas, M. A. Hernández, H. Ho, L. Popa, and M. Roth. Clio Grows Up: From Research Prototype to Industrial Tool. In *SIGMOD*, pages 805–810, 2005.
- [20] H. Hosoya and B. C. Pierce. Xduce: A statically typed xml processing language. *ACM Trans. Inter. Tech.*, 3(2):117–148, 2003.
- [21] M. Lenzerini. Data Integration: A Theoretical Perspective. In *PODS*, pages 233–246, 2002.
- [22] J. Madhavan, P. A. Bernstein, A. Doan, and A. Y. Halevy. Corpus-based Schema Matching. In *ICDE*, pages 57–68, 2005.
- [23] J. Madhavan and A. Y. Halevy. Composing Mappings Among Data Sources. In *VLDB*, pages 572–583, 2003.
- [24] E. Meijer, B. Beckman, and G. M. Bierman. LINQ: reconciling object, relations and XML in the .NET framework. In *SIGMOD*, page 706, 2006.
- [25] R. J. Miller, L. M. Haas, and M. A. Hernández. Schema Mapping as Query Discovery. In *VLDB*, pages 77–88, 2000.
- [26] P. Papotti and R. Torlone. Schema exchange: Generic mappings for transforming data and metadata. *Data Knowl. Eng.*, 68(7):665–682, 2009.
- [27] L. Popa, Y. Velegrakis, R. J. Miller, M. A. Hernández, and R. Fagin. Translating Web Data. In *VLDB*, pages 598–609, 2002.
- [28] J. van den Bussche and E. Waller. Type inference in the polymorphic relational algebra. In *PODS*, pages 80–90, 1999.
- [29] J. Vouillon. Polymorphic regular tree types and patterns. In *POPL 06*, pages 103–114, New York, NY, USA, 2006. ACM.
- [30] J. B. Wells. The essence of principal typings. In *ICALP '02*, pages 913–925, London, UK, 2002. Springer-Verlag.