

# Conjunctive Constraint Mapping for Data Translation

Chen-Chuan K. Chang  
Department of Electrical Engineering  
Stanford University  
Stanford, CA 94305, USA  
E-mail: changcc@cs.stanford.edu

Héctor García-Molina  
Department of Computer Science  
Stanford University  
Stanford, CA 94305, USA  
E-mail: hector@cs.stanford.edu

## ABSTRACT

In this paper we present a mechanism for translation of information in heterogeneous digital library environments. We model information as a set of conjunctive constraints that are satisfied by real-world objects (*e.g.*, documents). Through application of semantic rules and value transformation functions, constraints are mapped into ones that are understood and supported in another context. Our machinery can also deal with hierarchically structured information.

**KEYWORDS:** constraint mapping, data translation, semantic interoperability.

## 1 INTRODUCTION

Digital libraries often need to interact with autonomous systems that structure and represent their information differently. To ensure semantic interoperability [1], information must be appropriately mapped from its source context to its target context where it will be used. For example, the condition [author = "John Smith"] in a query, may need to be translated to [name = "Smith, J. "]. Similarly, the component [date = "Jan. 1998"] in the resulting answer (either metadata about a document or the document itself) may have to be translated to the pair of components [year = 1998] and [month = "January"].

For translation purposes we view “information” as a collection of *constraints* of the form [attrib = val] or simply [attrib val]. If a constraint refers to a document, it indicates a “fact” about it, *i.e.*, that the attribute or property attrib has the value val. If it refers to a query, then the constraint specifies a desired fact about matching documents. Our attributes can be hierarchical, *e.g.*, publication.date.year (*i.e.*, paths). In general, constraints could have non-equality operators, *e.g.*, [date > "Jan. 1998"]. Due to space limitations, we only consider the equality operator, and therefore do not explicitly write it. Our work can be generalized to other operators.

In this paper, we interpret a set of constraints as a *conjunction*, *i.e.*, the desired or described document satisfies all of the constraints. This is adequate for describing documents and vector-space queries (where a query is simply a document similar to the desired ones). However, Boolean queries may include disjunctions, *e.g.*, [author "Smith"] OR [date "Jan. 1998"]. In [2] we describe how our work can be extended to deal with constraint disjunction. It is important to note that translating a set of conjunctive constraints into another is not simply a matter of mapping each constraint separately. In general, the mappings can be many-to-many. For example, the constraints [car-type "ford-taurus"] AND [year 1994] may yield [make "ford"] AND [model "taurus-94"] at the target.

In this paper we present a framework and algorithms for information translation. The translation is performed by *rules*. Intuitively, each rule specifies how one or more source constraints are to be transformed into one or more target constraints. For translating the values themselves, *e.g.*, "John Smith" into the first initial and last name strings "J.", "Smith," we rely on *functions* that can be written in any programming language. The rules and function are written by human experts such as librarians, either for pairs of sources or pairs of attribute sets (*e.g.*, Dublin Core to US-MARC). Eventually, we expect that libraries of common translation functions will be available, so the librarian will mainly have to select which functions (with what parameters) have to be invoked to perform the translation. At run time, queries and documents can then be translated automatically using the rules provided. Fig. 1 gives an overview of the translation process: The librarians analyze the source and target contexts (*i.e.*, their schemas), and define the mapping rules, which then drive the constraint mapping algorithm to perform data translation.

There are many ways one can go about translating digital library information [3, 4, 5, 6]. We believe that our approach offers the following advantages over other proposals.

- Because we view information as constraints, our machinery can provide “minimal” translations without explicit user guidance. For example, if several rules are available for translating say first, middle and last names of authors,

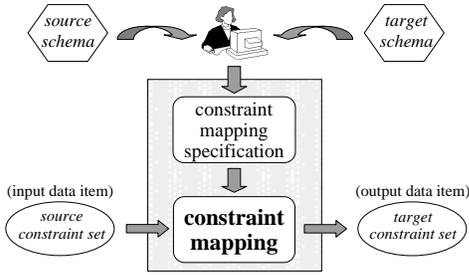


Figure 1: Overview of the data translation process.

the one that covers the most constraints will be applied.

- Our framework provides a natural way for dealing with hierarchically structured documents (including HTML, SGML). We encode the location of each value into a hierarchical attribute name, and then use the same simple machinery for translating the value into its appropriate target location.
- We include several features that have turned out to be very useful in a digital library environment. For instance, we provide facilities for dealing with lists of values (*e.g.*, for converting a list of authors to a single string containing all names), and for adding default constraints *e.g.*, if no author is given, generate [author "anonymous"].
- Our last “advantage” is subjective, but we believe that our framework provides the correct balance between simplicity and power that is needed in a digital library environment. At one end of the spectrum, we could say that a single function translates full documents or queries. This is a very simple framework, but leaves all the work for the function. At the other end, we could eliminate functions and do all processing within the translation framework. For example, to translate a PowerPoint image to a GIF, we would not call a function, but we would write rules that specify how each component of the image is translated. This would require the framework to be as powerful as a general purpose programming language (including recursion). As we will see, our framework focuses on the logical mappings between constraints, leaving the pure procedural transformation of *values* to the functions. This makes it easy to both implement and use our framework.

We start by defining the subsuming mapping and other fundamental notions for our framework. In Section 3 we develop the rule system framework and the associated algorithm for data translation, initially for “flat” data (*i.e.*, those naturally represented as attribute-value pairs). Section 4 then generalizes the translation framework for hierarchical data (*i.e.*, those with nested structures). Finally, in Section 5 we briefly review related efforts on data translation.

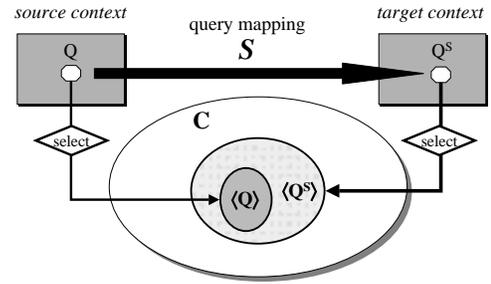


Figure 2: Conceptual illustration of query mapping.

## 2 PRELIMINARIES

Our universe is a collection  $C$  of information *objects*, each representing a real world entity of interest (*e.g.*, a book, a document, a car, a person). We define a *query* (or a discriminator) to be a set of conjunctive constraints that select or identify one or more objects. For example, the set  $Q = \{c_1, c_2\}$ , where  $c_1 = [\text{lname "Smith"}]$ , and  $c_2 = [\text{fname "John"}]$  is a query that selects the qualified objects (*i.e.*, persons named "Smith, John") from our universe of objects. To stress the fact that  $Q$  is a conjunction (selected object must satisfy all constraints), we write it as  $c_1 \text{ AND } c_2$ , or as  $\text{AND}(\{c_1, c_2\})$ . An individual object (*e.g.*, a document) can be represented by the query that uniquely identifies it.

Fig. 2 conceptually illustrates the problem we face. We are given a query  $Q$  in a *source context*. This context may limit how constraints are expressed, semantically or syntactically. For instance, a context may not support some attributes, or support them with different attribute names, value formats, unit, currency, *etc.* Within the source context,  $Q$  can be interpreted and identifies a set of objects  $\langle Q \rangle$  that satisfy all the  $Q$  constraints.

Our problem is to transform the information encoded by  $Q$  into another set of constraints  $Q^S$  that can be interpreted in a different *target context* (Fig. 2). In particular, we would like to have a *query mapping function*  $S(\cdot)$  that transforms any query like  $Q$  into a “good” query  $Q^S$  for the target context. What does “good” mean here? First,  $Q^S$  should be interpretable in the target context, *e.g.*, it should use appropriate attributes and values. Second, the objects selected by  $Q^S$ ,  $\langle Q^S \rangle$ , should be a superset of  $\langle Q \rangle$  (we do not want to miss any objects). For example, consider again the "Smith, John" query  $c_1 \text{ AND } c_2$ . Suppose that the target context does not support the *fname* attribute, but does support last names. Then we may have to translate this query into  $c_1$ . This query will identify more objects than the original one, and is acceptable if there is no other way to identify at the target a smaller set of objects that includes “Smith, John.” The latter minimality requirement represents our third condition for a good mapping.

The next definition formalizes these three properties. Note that the notion of query subsumption is directly related to *query containment* in deductive databases [7] and has been applied extensively in information integration [8].

**Definition 1 (Minimal Subsuming Mapping):** Let  $\mathcal{S}(\cdot)$  be a mapping on queries of a source system  $T_s$ . (We use the notation  $Q^S$  for  $\mathcal{S}(Q)$ ).  $\mathcal{S}(\cdot)$  is the *minimal subsuming mapping* w.r.t. some target system  $T$ , if for any query  $Q$  in  $T_s$ ,

1.  $Q^S$  is expressible in the context of  $T$ ,
2.  $Q^S$  subsumes  $Q$ , i.e.,  $\langle Q \rangle$  is a subset of  $\langle Q^S \rangle$ ;  $\langle Q \rangle \subseteq \langle Q^S \rangle$ ,
3.  $Q^S$  is *minimal*, i.e., there is no query  $Q'$  such that (i)  $Q'$  satisfies 1 and 2, and (ii)  $Q^S$  subsumes  $Q'$ . ■

A good mapping  $\mathcal{S}(\cdot)$  cannot translate the individual constraints independently [9], as the following example illustrates.

**Example 1:** Consider again the mapping of  $Q = c_1 \text{ AND } c_2$ , where  $c_1$  is [fname "John"], and  $c_2$  is [lname "Smith"]. Suppose that, corresponding to lname and fname, the target system supports the name attribute, but requires that at least the last name component be specified. That is, a name can be stated as "Smith, John", or simply "Smith", if the first name is not known.

If we assume that the constraints are “independent” and process them separately, we obtain:

$$\begin{aligned} \mathcal{S}(Q) &= \mathcal{S}(c_1) \text{ AND } \mathcal{S}(c_2) \\ &= \text{True AND [name "Smith"]} \\ &= [\text{name "Smith"}], \end{aligned}$$

where  $\mathcal{S}(c_1) = \text{True}$  because there is no smaller mapping when only the first name is known. Note that *True* actually means no constraint.

Clearly, the above is not the minimal translation, because the fname constraint is effectively discarded. In fact, constraints  $c_1$  and  $c_2$  are “interrelated” because they collectively decide the target constraint on name. In other words, by considering  $c_1$  and  $c_2$  together, we obtain the minimal translation  $\mathcal{S}(Q) = [\text{name "Smith, John"}]$ . ■

As Example 1 shows, to obtain the optimal mappings, the translation framework cannot blindly assume that constraints are independent. Unfortunately, there is no automatic way for computers to figure out the interdependence of constraints; that knowledge must be supplied by human experts in a computable way. In Section 3 we introduce a rule system framework for data translation that encodes the “dependence” of constraints, and that defines how constraints should be mapped. In Section 4 we extend this rule system to hierarchical information models (such as SGML).

### 3 CONJUNCTIVE QUERY MAPPING

Given a constraint set  $I$  representing a conjunctive query  $\text{AND}(I)$ , our goal is to find its minimal subsuming mapping,  $\mathcal{S}(\text{AND}(I))$ , which is a conjunctive query in the target context. Fig. 3 shows a running example that we use in this section to illustrate the mapping process. A site Car maintains information on used cars; a second site Auto would like

Car instance $I_{car}$		Auto instance $I_{auto}$	
[lname "Smith"]	( $c_1$ )	[dealer "Ford 101"]	( $a_1$ )
[fname "John"]	( $c_2$ )	[owner "John Smith"]	( $a_2$ )
[license "4HQD973"]	( $c_3$ )	[id "4HQD973"]	( $a_3$ )
[price 9000]	( $c_4$ )	[price 9000]	( $a_4$ )
[car-type "ford-taurus"]	( $c_5$ )	[model "taurus-94"]	( $a_5$ )
[year "1994"]	( $c_6$ )	[make "Ford"]	( $a_6$ )
[mileage 80000]	( $c_7$ )	[mileage 50000]	( $a_7$ )
[unit "km"]	( $c_8$ )	[color-code "b1"]	( $a_8$ )
[color "blue"]	( $c_9$ )		

Figure 3: Mapping instance from Car to Auto.

$K_{car} \equiv \{\mathcal{R}_{c.1}, \mathcal{R}_{c.2}, \mathcal{R}_{c.3}, \mathcal{R}_{c.4}, \mathcal{R}_{c.5}, \mathcal{R}_{c.6}, \mathcal{R}_{c.7}\}$	
$\mathcal{R}_{c.1}$	[A1 X]; SimpleMapping(A1) $\mapsto$ A2 = AttrNameMapping(A1); emit: [A2 X]
$\mathcal{R}_{c.2}$	[lname L]; [fname F] $\mapsto$ A = LnFnToName(L,F); emit: [owner A]
$\mathcal{R}_{c.3}$	[lname L] $\mapsto$ emit: [owner L]
$\mathcal{R}_{c.4}$	[car-type T]; [year Y] $\mapsto$ M = Model(T,Y); A = Make(T); emit: [model M] AND [make A]
$\mathcal{R}_{c.5}$	[mileage M1]; [unit U1] $\mapsto$ M2 = CvtLengUnit(M1,U1,"mile") emit: [mileage M2]
$\mathcal{R}_{c.6}$	[color C] $\mapsto$ CC = ColorCode(C) emit: [color-code CC]
$\mathcal{R}_{c.7}$	$\mapsto$ emit: [dealer "Ford 101"]

Figure 4: Mapping specification  $K_{car}$  for translating instances from Car to Auto.

to collect data from the first site (and possibly other sites). Thus, information from the source context Car must be translated into the Auto target context. Fig. 3 shows on the left a particular object as represented by Car. That is, if  $I_{car}$  represents the constraints on the left of the figure ( $I_{car} \equiv \{c_1, c_2, \dots, c_9\}$ ), then the query  $\text{AND}(I_{car})$  represents that object. On the right of Fig. 3 we show the same object as represented in the target context. If  $I_{auto}$  are the constraints on the right ( $I_{auto} \equiv \{a_1, a_2, \dots, a_8\}$ ), then we would like our mapping to yield the target query  $\text{AND}(I_{auto})$ . That is,  $\mathcal{S}(\text{AND}(I_{car})) = \text{AND}(I_{auto})$ .

Given a source constraint set, our algorithm first maps individual constraints as directed by a rule system, and then formulates the mapping of the source query as a whole. The rule system, called a *mapping specification*, consists of a set of *mapping rules*, each specifying the mapping of some constraint patterns. For instance, Fig. 4 shows a mapping specification  $K_{car}$  for translating instances of Car to Auto. In the remainder of this section, we first discuss the rule system framework, and then the algorithm that uses the rule system to handle conjunctive queries.

#### 3.1 Mapping Rules

A constraint mapping rule specifies how a source constraint, or the conjunction of some related source constraints, can be mapped to the target side. Intuitively, we may think of the operation of a rule for constraint mapping as *pattern matching*, similar to that of the rules used in Lex and Yacc. In particular, a rule specifies some *pre-conditions* in its left hand side,

the *head*, and some *actions* in its right hand side, the *tail*, as separated by the  $\mapsto$  symbol. The pre-conditions state the source constraints to be matched, and the actions direct how the corresponding target constraints can be generated.

Rule  $\mathcal{R}_{c.6}$  (Fig. 4) states that a source constraint on the color attribute is mapped to a target constraint on the color-code attribute, with the color value replaced with some standard code. The rule calls upon the function `ColorCode()` to look up the color code. Similarly, rule  $\mathcal{R}_{c.2}$  specifies that the source constraints on `fname` and `lname` are translated jointly by first merging them into a full name, and then formulating the target constraint on `owner`.

The head of a rule consists of constraint patterns and conditions to be matched against the source constraints. A *constraint pattern* is a 2-tuple [attribute value], where each component can be either a constant (of the corresponding domain) or *variable*. A source constraint *matches* a pattern if both agree on the non-variable components. As the result of a matching, the variables of the pattern is then *bound* to the corresponding constants in the matching constraint. For instance, after matching the pattern [color C] with the constraint [color "blue"], we obtain the binding of C to the constant "blue".

A *condition* is a predicate function that takes bound variables as arguments, and return either *True* or *False*. Conditions can be used in the head of a rule, in addition to constraint patterns, to restrict matchings to only those variable bindings that hold the conditions. For instance, in  $\mathcal{R}_{c.1}$ , the condition `SimpleMapping(·)` tests if the attribute A1 is one that requires only straightforward name mapping. If so, the matching constraint will be “copied” with only the attribute name mapped appropriately. For example, referring to Fig. 3, license and price are such “simple” attributes.

The tail (right hand side) of a rule directs, for each set of matching source constraints, how the corresponding target constraints can be generated. It consists of two parts: a list of *function statements* and an “*emit:*” clause, with the former responsible for converting value formats, attribute mapping, and so on, and the latter specifying the corresponding target constraints to be generated.

A function statement is of the form:

$$Y_1, Y_2, \dots, Y_m = \text{FunctionName}(X_1, X_2, \dots, X_n),$$

where  $X_i$ 's are input constants or variables that are already bound in the rule head or the preceding function statements, and  $Y_j$ 's are variables to be bound by the function output. For instance, function `LnFnToName(·)` (in  $\mathcal{R}_{c.2}$ ) takes the bound variable L and F as its argument, and binds A as the full name. Similarly, function `CvtLengUnit(M1,U1,"mile")` converts length M1 from unit U1 to "mile" (*i.e.*, the target mileage is always specified in miles).

Note that the functions in the tail as well as the conditions in the head are supplied externally, and in principle can be

written in any programming languages.

As the last component of a rule, the *emit:* clause specifies the corresponding target constraints. Note that it can generate more than one constraint connected with AND. (Recall that constraint mapping is in general many-to-many.) For instance, rule  $\mathcal{R}_{c.4}$  produces the conjunction of the constraints on `model` and `make`.

Having described the components of mapping rules, we next discuss their evaluation and semantics. The evaluation of a rule, given a set of source constraints as inputs, involves finding matching constraints, from which the target constraints will be emitted as directed by the rule actions. Moreover, the semantics of the rule requires that the emitted constraints be the minimal subsuming mapping of the matching constraints (except for a special case, see later). In Definition 2 we formalize the notion of matchings and emissions. Example 2 then illustrates these notions and uses them to show the rule semantics, which we will formalize in Definition 3.

**Definition 2 (Matchings and Emissions):** Let  $r$  be a mapping rule and  $I$  be a set of source constraints.

1. A *matching set* (or *matching* for short) of  $r$  w.r.t.  $I$  is a subset of  $I$ , of which the constraints together satisfy the head of  $r$ . We denote the set of *all the matching sets* of  $r$  w.r.t.  $I$  by  $\mathcal{M}(r, I)$ .
2. The *emission* of  $r$  w.r.t. a matching  $m$ , *i.e.*,  $m \in \mathcal{M}(r, I)$ , is the query generated by the *emit:* clause of  $r$  w.r.t.  $m$ , and is denoted  $\mathcal{E}(r, m)$ . ■

**Example 2:** Let's consider the evaluation of the rules in  $K_{car}$  (Fig. 4) with respect to the source constraint set  $I_{car}$  (Fig. 3).

For  $\mathcal{R}_{c.1}$ , the pattern [A1 X] can match any constraint in  $I$ . Assuming that the function `SimpleMapping(·)` returns *True* for only license and price, we obtain two matchings:  $\{c_3\}$  and  $\{c_4\}$ , *i.e.*,  $\mathcal{M}(\mathcal{R}_{c.1}, I) = \{\{c_3\}, \{c_4\}\}$ .

With the matching  $\{c_3\}$ , assume  $\mathcal{R}_{c.1}$  emits the target constraint  $a_3$ , *i.e.*,  $\mathcal{E}(\mathcal{R}_{c.1}, \{c_3\}) = a_3$ . By definition, this means  $\mathcal{S}(c_3) = a_3$ .

For  $\mathcal{R}_{c.2}$ ,  $\mathcal{M}(\mathcal{R}_{c.2}, I) = \{\{c_1, c_2\}\}$ , *i.e.*, constraints  $c_1$  and  $c_2$  together form a matching. Since  $\mathcal{E}(\mathcal{R}_{c.2}, \{c_1, c_2\}) = a_2$ , we have

$$\mathcal{S}(\text{AND}(\{c_1, c_2\})) = \mathcal{S}(c_1 \text{ AND } c_2) = a_2.$$

Finally, observe that rule  $\mathcal{R}_{c.7}$  does not require any constraints to form a matching, because the rule head is empty. In this case,  $\mathcal{M}(\mathcal{R}_{c.7}, I) = \{\emptyset\}$  (the empty set is a matching). Note that the rule will be fired no matter what the source constraints are. The special case of “empty” matching is useful in making implicit constraints explicit at the target side. For  $\mathcal{R}_{c.7}$ , even based on no source constraints, the rule still emits [dealer "Ford 101"], an implicit constraint that all Car objects assume.

The next definition tells us that each rule emits “good” constraints for the target, in the sense discussed in Section 2. However, this notion of goodness does not apply for empty matchings. In our example,  $\mathcal{S}(\text{AND}(\emptyset)) \neq [\text{dealer "Ford 101"}]$ , because  $\text{AND}(\emptyset) \equiv \text{True}$  (i.e., no constraints). ■

**Definition 3 (Rule Semantics):** Let  $r$  be a mapping rule and  $I$  be a set of source constraints. For any matching  $m$ , s.t.  $m \in \mathcal{M}(r, I)$

1. if  $m \neq \emptyset$ , then  $\mathcal{S}(\text{AND}(m)) \equiv \mathcal{E}(r, m)$ . That is, the emission of the rule *defines* the minimal subsuming mapping of the matching constraints.
2. otherwise, if  $m = \emptyset$ , then the emission of the rule represents the mapping of some *implicit constraints* assumed at the source system. ■

### 3.2 Algorithm for Conjunctive Query Mapping

The conjunctive query mapping problem can be stated as follows. Given a set of source constraints  $I$  representing a conjunctive query  $\text{AND}(I)$ , and a set of rules  $K$  as the constraint mapping specification w.r.t. some target system  $T$ , the problem is to find the minimal subsuming mapping of the source query  $\text{AND}(I)$  w.r.t.  $T$ , i.e.,  $\mathcal{S}(\text{AND}(I))$ .

The algorithm is relatively straightforward; essentially, the mapping of the whole query (as a constraint set) is the conjunction of the mappings of some subsets of the constraints. First, we evaluate the mapping rules with respect to the source constraints to find the matchings. Note that each matching is a subset of constraints that are deemed “interrelated” and thus must be processed together (e.g., constraints  $c_1$  and  $c_2$  in Example 2). In other words, the mapping rules effectively partition the source constraint set into subsets of interrelated constraints, and define their mapping (as the emissions). The mappings of those matching subsets are then assembled in a conjunctive form as the mapping of the whole source query.

However, some matching subsets may be redundant and should be removed. For instance, in  $K_{car}$  (Fig. 4)  $\mathcal{R}_{c.2}$  defines the mapping to owner from both the source lname and fname, while  $\mathcal{R}_{c.3}$  from only the lname constraint. Note that  $\mathcal{R}_{c.3}$  can be useful to generate a partial name of owner if fname is optional and can be omitted in some Car instances. However, for a source instance with both lname and fname, such as  $I_{car}$ ,  $\mathcal{R}_{c.3}$  gives a redundant matching  $\{c_1\}$ , because in this case  $\mathcal{R}_{c.2}$  can generate better mappings (i.e., full names) with the “larger” matching  $\{c_1, c_2\}$ .

In general, if a matching is a subset of some other matchings, we can eliminate the former because, with a larger set of source constraints, the latter will generate a “smaller” mapping. However, when the matching is the empty set, it cannot be eliminated because it actually represents some implicit constraints, as discussed in the preceding section. Due to space limitation, we are not able to discuss the formalism

supporting this *sub-matching suppression*; please refer to [2] for more details.

In the following, we present the mapping algorithm, and then illustrate it in Example 3. Readers can refer to [2] for the proof that the algorithm does generate minimal subsuming mappings.

#### Algorithm 1 (Conjunctive Query Mapping):

• **Input:**

- $I$ : a set of constraints.
- $K$ : the constraint mapping specification w.r.t. the target system  $T$ .

• **Output:**  $\mathcal{S}(\text{AND}(I))$ , the minimal subsuming mapping of  $\text{AND}(I)$  w.r.t.  $T$ .

• **Procedure:**

1. Find all the matchings for all the rules:

$$A := \cup(\mathcal{M}(r, I), \text{ for all } r \in K).$$

2. Remove any non-empty matching that is a subset of other matchings (sub-matching suppression):

for all  $m_i \in A$ :

for all  $m_j \in A$  ( $j \neq i$ ),

if  $m_j \subseteq m_i$  and  $m_j \neq \emptyset$ , remove  $m_j$  from  $A$ .

3. Output the conjunction of all the emissions:

$$\mathcal{S}(\text{AND}(I)) = \text{AND}(\mathcal{S}(\text{AND}(m_i))) = \text{AND}(\mathcal{E}(r, m_i)),$$

for all  $m_i \in A$ , such that  $m_i$  matches  $r$ , i.e.,  $m_i \in \mathcal{M}(r, I)$ . ■

**Example 3:** To illustrate Algorithm 1 for data mapping, we use it to translate the Car instance  $I_{car}$  (Fig. 3). That is, we run the algorithm with the constraint set  $I_{car}$  and mapping specification  $K_{car}$  (Fig. 4) as inputs, and show that it does output  $I_{auto}$  as the mapping. In other words, we show that  $\mathcal{S}(\text{AND}(I_{car})) = \text{AND}(I_{auto})$ . We illustrate the process step by step.

1.  $A = \cup(\mathcal{M}(\mathcal{R}_{c.1}, I), \mathcal{M}(\mathcal{R}_{c.2}, I), \dots, \mathcal{M}(\mathcal{R}_{c.7}, I), )$   
 $= \cup( \{ \{c_3\}, \{c_4\} \}, \{ \{c_1, c_2\} \}, \{ \{c_1\} \}, \{ \{c_5, c_6\} \},$   
 $\{ \{c_7, c_8\} \}, \{ \{c_9\} \}, \{ \emptyset \} )$   
 $= \{ \{c_3\}, \{c_4\}, \{c_1, c_2\}, \{c_1\},$   
 $\{c_5, c_6\}, \{c_7, c_8\}, \{c_9\}, \emptyset \}$
2.  $\{c_1\}$  is removed, because it is a subset of  $\{c_1, c_2\}$ .
3.  $\mathcal{S}(\text{AND}(I_{car}))$   
 $= \text{AND}(\mathcal{E}(\mathcal{R}_{c.1}, \{c_3\}), \mathcal{E}(\mathcal{R}_{c.1}, \{c_4\}),$   
 $\mathcal{E}(\mathcal{R}_{c.2}, \{c_1, c_2\}), \mathcal{E}(\mathcal{R}_{c.4}, \{c_5, c_6\}), \dots,$   
 $\mathcal{E}(\mathcal{R}_{c.7}, \{ \emptyset \} )$   
 $= \text{AND}(a_3, a_4, a_2, a_5 \text{ AND } a_6, a_7, a_8, a_1)$   
 $= \text{AND}(I_{auto})$  ■

## 4 HIERARCHICAL DATA TRANSLATION

We have studied the translation machinery for “flat” data (as sets of attribute-value pairs); in this section we generalize the framework for hierarchical (*i.e.*, nested) data models. While the flat representation of data has been widely used (*e.g.*, BibTex, Dublin Core [10], and relational databases), hierarchical structures have also been developed to provide richer information abstraction. In particular, they have been used for structured documents (*e.g.*, SGML [11], XML [12]), metadata (*e.g.*, USMARC [13], Warwick Framework [14]), semi-structured data (*e.g.*, the OEM model in [5]), and scientific data [6].

This section is organized as follows. In Section 4.1, we discuss what we view as hierarchical data by introducing the underlying “conceptual” model as well as our restrictions (*e.g.*, we do not handle recursive schemas, see later). We then generalize the framework for translating data in this model. First, to base translation on the constraint mapping framework, we need to extend the notion of constraints for tree-like hierarchical data, so that they can be represented as conjunctive constraints. Second, the evaluation of mapping rules (in particular, the matching of constraint patterns) is slightly complicated because patterns can match “macro” constraints that represent, say, subtrees. We discuss each of them in Section 4.2 and 4.3 respectively.

### 4.1 Conceptual Grammatical Model

As a basis for hierarchical data translation, we conceptually abstract data items (*e.g.*, documents, metadata records, *etc.*) as EBNF (Extended Backus-Naur Form [15]) grammar-generated trees. Recognizing the fundamental characteristics of data to be its hierarchical and grammatical structure, many have proposed using the well-understood concept of formal language grammars to model data [16, 11]. For instance, in SGML [11] (or its simplified version, XML [12], an evolving standard), documents are defined with a DTD (document type definition), which is essentially a grammar. In particular, we adopt and extend the formalism of *grammatical model* [16] to model hierarchical data.

**Definition 4 (Conceptual Grammatical Model):** A *conceptual grammatical schema* is a five-tuple  $G = (V, T, P, S, D)$ :

- $(V, T, P, S)$  is a context free grammar in EBNF, of which  $V$  and  $T$  are finite sets of *structural attributes* and *terminal attributes* respectively,  $S$  a *start symbol*,  $S \in V$ , and  $P$  a finite set of EBNF production of the form  $A \rightarrow \alpha$  where  $A \in V$ ,  $\alpha$  is a non-empty regular expression over the alphabet  $V \cup T$ .
- $D$  is a mapping on  $T$ , *s.t.*  $D(A)$  is a set called the domain of  $A$ , for  $A \in T$ .

An *instance* (or *instance tree*) of  $G$  is a *derivation tree* for  $G$ , of which the leaves (labeled with terminal attributes) are assigned *values* from the corresponding domains as defined by  $D$ . ■

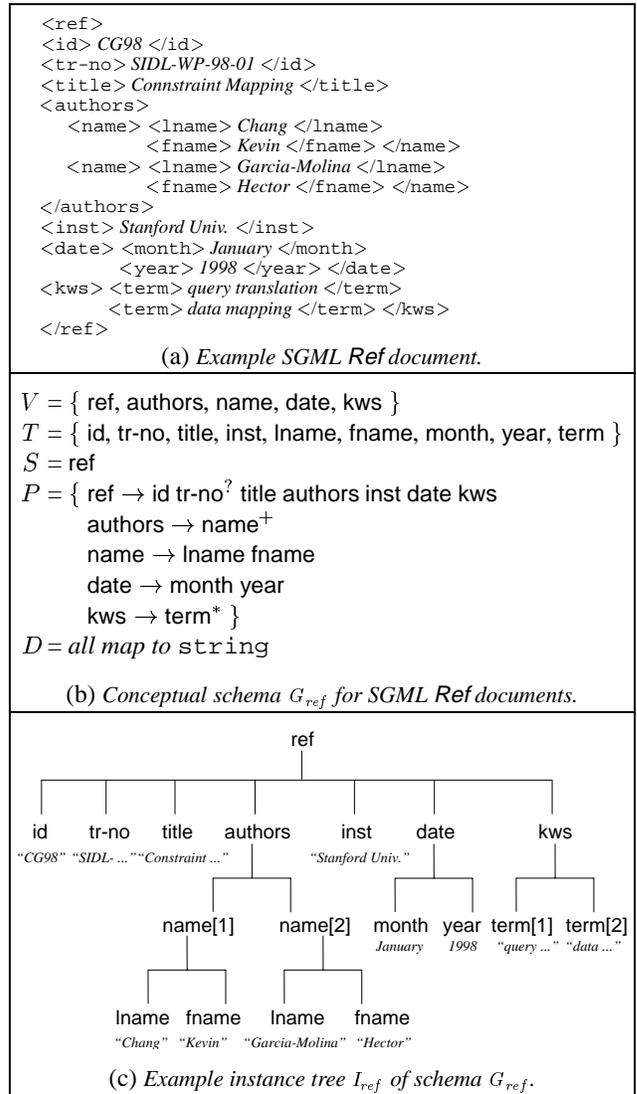


Figure 5: Schema and instance for SGML Ref.

We illustrate the notions with the running examples in SGML (Fig. 5) and BibTex (Fig. 6). Specifically, Fig. 5(a) shows an (assumed) SGML document of type Ref for recording bibliographic reference information. In (b) we illustrate the conceptual schema that defines the logical structure of such documents. The example document is then represented as an instance tree in (c). Furthermore, Fig. 6 shows the same information in the BibTex TechReport format; in particular, it also shows the corresponding BibTex entry in (a), the conceptual schema in (b), and the instance tree in . We will use these examples to illustrate the translation process, assuming we want to translate information from Ref to TechReport. Note that, in the regular expressions of grammars, we use  $*$  for Kleene closure (0 or more times),  $+$  for positive closure (1 or more times),  $?$  for optionality (0 or 1 time), and  $|$  for separating alternatives. For instance,  $\text{tr-no}^?$  in  $G_{\text{ref}}$  indicates that the attribute is optional.

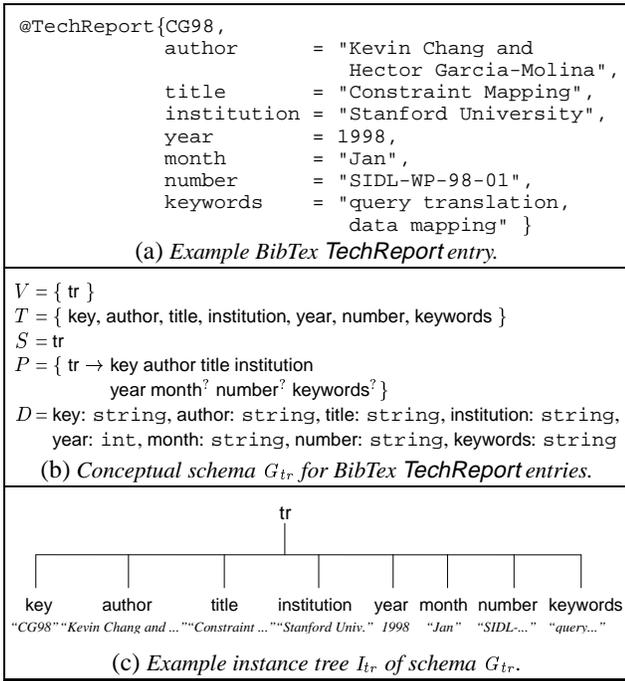


Figure 6: Schema and instance for Bibtex TechReport.

The iteration constructs in EBNF ("+" and "\*") provide a way for specifying *lists* of repeating attributes, which we use to model *collections*, for instance, a collection of names as the value of authors (authors  $\rightarrow$  name<sup>+</sup>). Note that, to unambiguously reference list elements, the repeated occurrences of an attribute (as generated by the iteration constructs) in instance trees (e.g., Fig. 5(c)) are ordered and labeled with the usual list operators, e.g., name[1] for the first element in the name list. In addition, the notation  $A[k_1:k_2]$  refers to the sub-list from  $k_1$ -th through the  $k_2$ -th elements. Either  $k_1$  or  $k_2$  can be omitted and defaulted to 1 and the length of the list, e.g.,  $A[:]$  refers to the whole list.

We believe the conceptual model is sufficient, for translation purposes, in describing data items observing some implicit or explicit grammatical structures. In particular, because a list is by definition an order and duplicate-sensitive collection, one may argue that we need other constructs such as *bags* (order-insensitive) and *sets* (order and duplicate-insensitive). However, it suffices to conceptually model them as lists, *i.e.*, by assuming some ordering in a set or bag, because the additional orderings should have no implication on the translation correctness. For instance, although the schema  $G_{ref}$  conceptually represents keywords as a list of terms (by keywords  $\rightarrow$  term\*), the underlying data source may actually model it as a set.

Without loss of generality, we assume *normalized* schemas, which ensure that the generated instance trees have distinct paths (*i.e.*, no two different paths have the same sequence of labels). This path uniqueness is critical because (for hierarchical data) constraints will be specified with paths (Sec-

instance $I_{ref} = \{c_1, c_2, \dots, c_{12}\}$	
[id "CG98"]	(c <sub>1</sub> )
[tr-no "SIDL-WP-98-01"]	(c <sub>2</sub> )
[title "Constraint ..."]	(c <sub>3</sub> )
[authors.name[1].lname "Chang"]	(c <sub>4</sub> )
[authors.name[1].fname "Kevin"]	(c <sub>5</sub> )
[authors.name[2].lname "Garcia-Molina"]	(c <sub>6</sub> )
[authors.name[2].fname "Hector"]	(c <sub>7</sub> )
[inst "Stanford Univ."]	(c <sub>8</sub> )
[date.month "January"]	(c <sub>9</sub> )
[date.year 1998]	(c <sub>10</sub> )
[kws.term[1] "query ..."]	(c <sub>11</sub> )
[kws.term[2] "data ..."]	(c <sub>12</sub> )

Figure 7:  $I_{ref}$  as a set of terminal constraints.

tion 4.2). Formally, a schema is normalized, if every grammar production has the following characteristics:

1. The iteration operators (\* and +) associate only with individual attributes. For instance,  $A \rightarrow (X|Y)^+$  is not allowed; instead, it should be normalized to  $A \rightarrow B^+$  and  $B \rightarrow X|Y$ .
2. Attributes are unique within a concatenation. For instance,  $A \rightarrow XYX$  is not allowed; instead it should be normalized (by annotation) to, say,  $A \rightarrow X_1YX_2$ .

We believe the restrictions are reasonable for grammars whose purpose is to describe the logical structures of information; most “natural” schemas are already normalized, for instance, the schemas of relational databases, BibTex, Dublin Core, USMARC, *etc.* For those rare cases of unnormalized schemas, it is straightforward to normalize them (as informally suggested in the above requirements). The full machinery for normalizing schemas and their instance trees is given in [2].

Note that our framework currently does not deal with schemas with recursive grammars, e.g., section  $\rightarrow$  title section, where section derives itself. Although in principle it should be feasible, supporting grammar recursion would require more “powerful” path patterns to represent recursively generated paths, which will greatly complicate the mechanism. In fact, we believe that non-recursive schemas cover most practical cases, for instance, document metadata and structured information.

## 4.2 Constraints as Path-Value Pairs

While for flat data constraints correspond naturally to attribute-value pairs, we need to generalize the notion of constraints for hierarchical data, which is in turn based on extending attributes to paths.

Intuitively, a natural generalization is to “flatten” a hierarchical structure (e.g., Fig. 5(c)) into *path-value* pairs, each representing a constraint (e.g., Fig. 7). More formally, we generalize a constraint to be of the form  $[P \mathcal{V}(P)]$ , where P is a *path* expression, and  $\mathcal{V}(P)$  the *value* associated with P. We next discuss what paths are, and their values.

In the simplest case, if a path P (starting at the root) leads to a leaf node, then  $\mathcal{V}(P)$  is the value assigned to the leaf node that P ends with. We call such P a *terminal path*, and  $[P \mathcal{V}(P)]$  a *terminal constraint*. For instance, in  $I_{ref}$  (Fig. 5) the path

date.year leads to the value 1998 (*i.e.*,  $\mathcal{V}(\text{date.year}) = 1998$ ), thus representing the constraint [date.year 1998]. Note that, as (for translation) we are only interested in paths starting from the root, we omit the root label (*e.g.*, ref) in the paths.

Note that, because each terminal path is unique (as guaranteed by normalized schemas [2]), an instance tree can be flattened into (as well as reconstructed from) a set of terminal constraints, which is a generalization of the constraint set representation for flat data. Thus, as we have intuitively observed, the constraint set in Fig. 7 is indeed a valid representation for the instance tree  $I_{ref}$  in Fig. 5(c),

Furthermore, as the translation is on tree structures, it may be desirable to have mapping functions that process whole “sub-structures”, *e.g.*, subtrees, or lists of terminal paths. In other words, we need *macro constraints* to conveniently represent sets of terminal constraints in the sub-structures. In the following, we discuss such constraints, which in turn are based on the notions of structural paths (for subtrees), and iteration paths (for list construction).

First, to represent a subtree, a *structural path* is a path that ends with some interior node, *e.g.*, in  $I_{ref}$ : authors.name[1], date, keywords, *etc.* The value of a structural path is the set of terminal constraints contained in the (sub-)tree rooted at the path. For example, in  $I_{ref}$   $\mathcal{V}(\text{authors.name}[1])$  is

$$\{[\text{lname "Chang"}], [\text{fname "Kevin"}]\},$$

and  $\mathcal{V}(\text{date})$  is

$$\{[\text{month "January"}], [\text{year "1998"}]\}.$$

Formally, the value of a structural path  $p_s$  *w.r.t.* an instance  $I$  is defined as  $\mathcal{V}(p_s) =$

$$\{[p' \ v] \mid p_s.p' \text{ is a terminal path in } I, v = \mathcal{V}(p_s.p')\}.$$

Second, we use iteration paths to construct lists of (terminal or structural) paths. An *iteration path*  $p_x$  is a sequence of labels,  $l_1.l_2.\dots.l_n$ , where some  $l_i$  is of the form  $a_i[k_1:k_2]$ . For instance, kws.term[:] denotes the (whole) list of keywords, and its value  $\mathcal{V}(\text{kws.term}[:])$

$$\begin{aligned} &= \langle \mathcal{V}(\text{kws.term}[1]), \mathcal{V}(\text{kws.term}[2]) \rangle \\ &= \langle \text{"query ..."}, \text{"data ..."} \rangle, \end{aligned}$$

*i.e.*, a list of terms (we use  $\langle \dots \rangle$  to denote a list).

As another example, authors.name[1:2] represents the list of subtrees rooted at authors.name[ $k$ ],  $k \in 1:2$ . Therefore, we have  $\mathcal{V}(\text{authors.name}[1:2])$

$$\begin{aligned} &= \langle \mathcal{V}(\text{authors.name}[1]), \mathcal{V}(\text{authors.name}[2]) \rangle \\ &= \langle \{[\text{lname "Chang"}], [\text{fname "Kevin"}]\}, \\ &\quad \{[\text{lname "Garcia..."}], [\text{fname "Hector"}]\} \rangle \end{aligned}$$

Formally, the value of an iteration path  $p_x$  is defined recursively as follows. Assuming  $p_x = p_1.a[k_1:k_2].p_2$ , *s.t.*  $a[k_1:k_2]$  is the first label with an iteration range,  $\mathcal{V}(p_x) :=$

$$\langle \mathcal{V}(p_1.a[k_1].p_2), \mathcal{V}(p_1.a[k_1+1].p_2), \dots, \mathcal{V}(p_1.a[k_2].p_2) \rangle$$

Finally, note that a macro constraint actually represents a set of terminal constraints. For instance, a constraint specified

with a structural path  $p_s$ ,  $[p_s \ \mathcal{V}(p_s)]$ , represents the set of terminal constraints contained in the subtree rooted at  $p_s$ . To see what terminal constraints are actually covered, as the translation algorithm requires, we can “flatten” a constraint with the following operational definition.

**Definition 5 (Flattening Function):** Let  $c = [p \ \mathcal{V}(p)]$  be a constraint. The flattening function  $\mathcal{F}(c)$  returns the set of terminal constraints contained in  $c$ :

1. if  $p$  is a terminal path (*i.e.*,  $c$  is a terminal constraint), then

$$\mathcal{F}([p \ \mathcal{V}(p)]) := \{[p \ \mathcal{V}(p)]\};$$

2. otherwise, if  $p$  is a structural path, then

$$\mathcal{F}([p \ \mathcal{V}(p)]) = \{[p.p' \ v] \mid [p' \ v] \in \mathcal{V}(p)\};$$

3. otherwise, if  $p$  is an iteration path  $p_1.a[k_1:k_2].p_2$ , where  $a[k_1:k_2]$  is the first label with an iteration range, then

$$\mathcal{F}([p \ \mathcal{V}(p)]) = \bigcup_{j=k_1}^{k_2} \mathcal{F}([p_1.a[j].p_2 \ \mathcal{V}(p_1.a[j].p_2)]) \quad \blacksquare$$

For example, to find the terminal constraints contained in the macro constraint with the iteration path author.name[1:2], we evaluate the flattening function as follows:

$$\begin{aligned} &\mathcal{F}([\text{author.name}[1:2] \ \mathcal{V}(\text{author.name}[1:2])]) \\ &= \bigcup_{j=1}^2 \mathcal{F}([\text{author.name}[j] \ \mathcal{V}(\text{author.name}[j])]) \\ &= \mathcal{F}([\text{author.name}[1] \\ &\quad \{[\text{lname "Chang"}], [\text{fname "Kevin"}]\}] \cup \\ &\quad \mathcal{F}([\text{author.name}[2] \\ &\quad \{[\text{lname "Garcia..."}], [\text{fname "Hector"}]\}]) \\ &= \{[\text{author.name}[1].\text{lname "Chang"}], \\ &\quad [\text{author.name}[1].\text{fname "Kevin"}]\} \cup \{\dots\} \\ &= \{c_4, c_5, \dots\} \cup \{c_6, c_7, \dots\} = \{c_4, c_5, c_6, c_7, \dots\} \end{aligned}$$

### 4.3 Constraint Mapping for Hierarchical Data

With the notion of constraints, we can translate hierarchical data using the constraint mapping machinery discussed in Section 3. Specifically, in our example of translating Ref (Fig. 5) to TechReport (Fig. 6), to apply Algorithm 1, the input will be a source instance tree (*e.g.*,  $I_{ref}$  in Fig. 5(c)), and a mapping specification consisting of rules (*e.g.*,  $K_{ref}$  in Fig. 8) that direct translation to the target context. Note that the source instance tree conceptually represents a set of terminal constraints, as the algorithm requires (*e.g.*,  $I_{ref} = \{c_1, c_2, \dots, c_{12}\}$  in Fig. 7). However, it is not necessary to actually “materialize” the constraint set representation; the matching and evaluation of rules can be done directly on the tree representation. As its output, the algorithm generates conjunctive constraints that represent the target instance tree (*e.g.*,  $I_{tr}$  in Fig. 6(c)).

The evaluation and semantics of mapping rules are essentially the same as what we discussed in Section 3, although the matching of constraint patterns is slightly more complex. In particular, a pattern can represent a macro constraint and thus match a set of terminal constraints, while in Section 3 a pattern may match only a single attribute-value pair. To see

$K_{ref} \equiv \{\mathcal{R}_{r.1}, \mathcal{R}_{r.2}, \mathcal{R}_{r.3}, \mathcal{R}_{r.4}, \mathcal{R}_{r.5}\}$	
$\mathcal{R}_{r.1}$	[A1 X]; SimpleMapping(A1) $\mapsto$ A2 = AttrNameMapping(A1); emit: [A2 X]
$\mathcal{R}_{r.2}$	[authors.name[:] L] $\mapsto$ S = ConsolidateNames(L); emit: [author S]
$\mathcal{R}_{r.3}$	[date D] $\mapsto$ M, Y = MonthYear(D); emit: [month M] AND [year Y]
$\mathcal{R}_{r.4}$	[kws.term[:K2] L]; K2 $\leq$ 3 $\mapsto$ S = MergeString(L); emit: [keywords S]
$\mathcal{R}_{r.5}$	NoConstraint([tr-no N]) $\mapsto$ emit: [number "to be assigned"]

Figure 8: Mapping specification  $K_{ref}$  for translating instances from  $G_{ref}$  to  $G_{tr}$ .

how data translation works, in the remaining of this section, we explain the evaluation of the rules in  $K_{ref}$  with respect to the source instance  $I_{ref}$ , which will generate  $I_{tr}$  as the translation.

For instance, the rule  $\mathcal{R}_{r.3}$  in  $K_{ref}$  converts the source date (subtree) into the target month and year. Given  $I_{ref}$ , the evaluation will consider all the paths (that start from the root), and find that `date` is a matching path for the pattern `[date D]`. That is, the pattern matches the macro constraint `[date  $\mathcal{V}(\text{date})$ ]` (denoted  $C_d$ ), resulting in binding `D` to  $\mathcal{V}(\text{date})$ , which is `{[month "January"], [year "1998"]}`. Note that  $C_d$  is a macro constraint representing the set of terminal constraints  $\mathcal{F}(C_d) = \{c_9, c_{10}\}$ , thus  $\mathcal{M}(\mathcal{R}_{r.3}, I_{ref}) = \{\{c_9, c_{10}\}\}$ . The function `MonthYear( $\cdot$ )` then takes `D` as input, in whatever data structure (appropriate for the function’s programming language) encoding the value of `D`, and returns the month and year values. Finally, the rule emits  $\mathcal{E}(\mathcal{R}_{r.3}, \{c_9, c_{10}\}) = [\text{month "January"}] \text{ AND } [\text{year 1998}]$  as the target constraints.

In general, a constraint pattern in mapping rules is of the form `[P X]`, where `P` is a *parameterized path* and `X` is a variable or constant representing  $\mathcal{V}(P)$ . A parameterized path `P` is a path expression that may contain variables to be assigned values such that `P` matches a path in the source instance. For instance, let’s consider  $\mathcal{R}_{r.1}$ , which handles those “simple” translations that require only attribute name mapping. In particular, `A1` is a parameterized path that can match any path in  $I_{ref}$  consisting of a single attribute (such as `id`, `tr-no`, `title`, `authors`, *etc.*), and `X` will be bound to whatever  $\mathcal{V}(A1)$  is. Further restricted by the condition `SimpleMapping(A1)`, which we assume returns `True` for the “simple” attributes `id`, `tr-no`, `title`, and `inst`, the matchings are thus  $\mathcal{M}(\mathcal{R}_{r.1}, I_{ref}) = \{\{c_1\}, \{c_2\}, \{c_3\}, \{c_8\}\}$ . For each matching,  $\mathcal{R}_{r.1}$  then emits the target constraint, *e.g.*,  $\mathcal{E}(\mathcal{R}_{r.1}, \{c_1\}) = [\text{key "CG98"}]$ .

A parameterized path can also contain variables that represent integers in the index ranges of an iteration path, *e.g.*, `kws.term[:K2]` in rule  $\mathcal{R}_{r.4}$ . Note that the rule emits the target keywords from at most the first 3 source terms (assuming the target context imposes this restriction). In evaluation, because the list `kws.term[:]` is of length 2 in  $I_{ref}$ , only the bind-

ings of `K2` to either 1 or 2 will result in valid (iteration) paths representing sub-lists. Furthermore, the binding of `K2 = 1` represents the matching

$$\mathcal{F}([\text{kws.term}[1:1] \mathcal{V}(\text{kws.term}[1:1])]) = \{c_{11}\},$$

while that of `K2 = 2`

$$\mathcal{F}([\text{kws.term}[1:2] \mathcal{V}(\text{kws.term}[1:2])]) = \{c_{11}, c_{12}\}.$$

(This example does not have the matchings of `K2 = 3` and above.) Therefore, the former “sub-matching” will be suppressed by the latter (step 2, Algorithm 1). This suppression mechanism, together with the condition `K2  $\leq$  3`, will select the “largest” matching within the restriction. Note that for `K2 = 2`, `L` is bound to the list `<"query ...", "data ...">`, thus generating the target constraint `[keywords "query ...", "data ..."]`.

As another example of iteration paths, rule  $\mathcal{R}_{r.2}$  emits the author constraint. The expression `authors.name[:]` matches the whole list (*i.e.*, `authors.name[:2]`), resulting in the binding of `L` to  $\mathcal{V}(\text{authors.name}[:2])$  (given in Section 4.2). The function `ConsolidateNames(L)` then merge the names in the required target format.

Finally, rule  $\mathcal{R}_{r.5}$  shows the special case of empty set as a matching to generate “default” constraints (see Definition 3), which is similar to  $\mathcal{R}_{c.7}$  (Fig. 4) and is not specific to hierarchical data. Note that a source `Ref` instance may not have the `tr-no` constraint (which is optional in  $G_{ref}$ , Fig. 5). The condition `NoConstraint([tr-no N])` checks if the pattern (in which `N` is simply a dummy variable) does not find any matching in the source tree. If so, the rule emits the default constraint `[tr-no "to be assigned"]`. Otherwise, the translation of `tr-no` will be handled by  $\mathcal{R}_{r.1}$  instead, as we just discussed.

## 5 Related Work

Information integration has long been recognized as a central problem of modern information systems [1, 8]. In this paper we have presented our data translation framework with the goal of coping with semantic or schematic inconsistency in data exchange.

Fully automatic semantic integration is extremely difficult [1], if not impossible. Thus, most related efforts [4, 5, 6] advocate, like we do, using human-specified rule systems (*e.g.*, datalog rules [7]). The common idea behind all these approaches is that a rule essentially specifies how target “patterns” are to be generated from matching source “patterns.”

Our work differs from other rule translation efforts in its data representations and its supporting theoretic framework, which in turn impacts how mapping rules are specified and evaluated. In particular, while we abstract data objects as “flattened” sets of constraints, other systems perform translation directly on the object structures (trees or graphs). Furthermore, because we view information as constraints, we address the data translation problem with the theoretic framework of conjunctive query mapping, and specifically using

the notion of query subsumption. To the best of our knowledge, no other existing work uses these notions for data translation. We believe that our approach leads to the following advantages:

- *Simple Rules.* Because of our unique data representation, our rules use patterns for constraints, instead of *object patterns* that match the object layouts. Consequently, our rules are more modular and declarative; independent constraints are separated into different rules, each declaring a correspondence between some source and target constraints. (We could call this a “divide and conquer” approach.) In contrast, with object patterns, rules must be composed in a way that mirror the structure of objects. In fact, object rules are typically not independent of each other, because they together direct the bottom-up (or top-down) construction of target structures, in a more procedural fashion.
- *Minimal Translations.* Since we model information as queries, we can generate minimal translations that use the most source constraints, when there are multiple options (as supported by sub-matching suppression, Section 3). In contrast, in other systems, first, it is cumbersome to enumerate all such options in object patterns (*e.g.*, the combination of the absence or presence of some attributes with respect to the whole layout). Second, it would then require an additional mechanism, *e.g.*, rule ordering, to resolve the multiple choices, which introduces more complications.
- *Duality of Data and Queries.* Although in this paper we mainly focused on data translation, our same mechanism can be used for the translation of user queries as well [2]. This unification can be very beneficial in the integration of heterogeneous digital libraries that support Boolean query languages.

Some data translation proposals do not rely on rules, as we do, but instead rely on explicit metadata that encodes the context for automated mediation. For instance, the COIN framework [3] uses the *semantic-value* model, where each attribute is annotated with a *property* list that defines its context (*e.g.*, unit = “usd” for attribute revenue). This approach assumes one-to-one mappings and thus cannot handle the general case of many-to-many mappings for interrelated attributes. In addition, it only deals with attribute “value” conversion, and is not applicable to hierarchical data. Lastly, it relies on implicit agreement on what the properties for different attributes are, as well as what conversion functions are applicable for translation.

Finally, we note that the framework described in this paper is closely related to and complements our earlier work on query translation for Boolean IR systems [9, 17]. In particular, in [9] we discussed the theoretic framework for optimal

query mapping; however, it assumes that constraints are independent in translation (see Section 2). The framework in this paper does consider the dependence. Furthermore, our earlier work does not incorporate the ability to call arbitrary functions to translate values, nor does it deal with hierarchical attributes. Actually, we can view the algorithms of [17] as providing translation functions for the particular case of “text search” constraints (*e.g.*, ones that search for words near others, or words matching wildcards).

## 6 CONCLUSION

In a heterogeneous environment, autonomous information systems often structure and represent their information inconsistently. To ensure semantic interoperability, information must be mapped appropriately when exchanged across interacting contexts. In this paper, we have presented a rule-system framework for information translation. We view information as a set of conjunctive constraints that identify the represented object(s). Consequently, our machinery provides “minimal” translation, as guaranteed by the supporting query mapping formalism [2].

We have also shown how to generalize the constraint mapping framework for hierarchical data, such as SGML documents. The generalization is based on extending attributes to paths in trees, thus flattening the hierarchy into a set of path-value pairs. In addition, we have also presented macro constraints that can be used to process a whole substructure in the trees, *e.g.*, a subtree, or a list of children.

Although our translation rules are modular and declarative, it may still be a tedious task to develop mapping specifications like the ones in Figures 4 and 8. Therefore, it may be useful to develop an interactive graphical interface, where a librarian can browse through metadata specifications for the source and target contexts, and through libraries of available value transformation functions. Through this front-end system the librarian can generate the specifications by selecting options on the screen, as opposed to writing text as shown in our figures. We believe that such a front-end design system can significantly enhance the effective data translation machinery we have developed.

## REFERENCES

1. Sandra Heiler. Semantic interoperability. *Computing Surveys*, 27(2):271–273, June 1995.
2. Chen-Chuan K. Chang and Héctor García-Molina. Constraint mapping for query and data translation. Technical report, Dept. of Computer Science, Stanford Univ., Stanford, California, 1998. In preparation.
3. Edward Sciore, Michael Siegel, and Arnon Rosenthal. Using semantic values to facilitate interoperability among heterogeneous information systems. *Transactions on Database Systems*, 19(2):254–290, June 1994.
4. Serge Abiteboul, Sophie Cluet, and Tova Milo. Cor-

- responsiveness and translation for heterogeneous data. In *Proceedings of ICDT '97*, Delphi, Greece, 1997. Springer.
5. Yannis Papakonstantinou, Héctor García-Molina, and Jeffrey Ullman. Medmaker: A mediation system based on declarative specifications. In *Proceedings of ICDE '96*, New Orleans, Louisiana, 1996.
  6. P. Buneman, S.B. Davidson, K. Hart, C. Overton, and L. Wong. A data transformation system for biological data sources. In *Proceedings of the 21st VLDB Conference*, Zurich, Switzerland, 1995.
  7. Jeffrey D. Ullman. *Principles of Database and Knowledge-Base Systems*. Computer Science Press, Rockville, M.D., 1988.
  8. Jeffrey D. Ullman. Information integration using logical views. In *Proceedings of ICDT '97*, Delphi, Greece, 1997.
  9. Chen-Chuan K. Chang, Héctor García-Molina, and Andreas Paepcke. Boolean query mapping across heterogeneous information sources. *IEEE Transactions on Knowledge and Data Engineering*, 8(4):515–521, August 1996.
  10. Stuart Weibel, Jean Godby, Eric Miller, and Ron Daniel, Jr. OCLC/NCSA metadata workshop report. Accessible at [http://purl.org/metadata/-dublin\\_core](http://purl.org/metadata/-dublin_core), March 1995.
  11. Charles F. Goldfarb. *The SGML Handbook*. Oxford University Press, New York, 1990.
  12. W3C. Extensible markup language (XML). Accessible at <http://www.w3.org/XML/>, 1997.
  13. Network Development and MARC Standards Office. *USMARC Format for Bibliographic Data: Including Guidelines for Content Designation*. Library of Congress, Cataloging Distribution Service, Washington, D.C., 1994.
  14. Carl Lagoze, Clifford A. Lynch, and Ron Daniel, Jr. The Warwick Framework: A container architecture for aggregating sets of metadata. Technical Report TR96-1593, Cornell University, Computer Science Dept., June 1996.
  15. Niklaus Wirth. What can we do about the unnecessary diversity of notation for syntactic definitions. *Communications of the ACM*, 20(11):822–823, November 1977.
  16. Marc Gyssens, Jan Paredaens, and Dirk Van Gucht. A grammar-based approach towards unifying hierarchical data models. In *Proc. of the 1989 ACM SIGMOD Conference*, 1989.
  17. Chen-Chuan K. Chang, Héctor García-Molina, and Andreas Paepcke. Predicate rewriting for translating boolean queries in a heterogeneous information system. Technical Report SIDL-WP-1996-0028, Dept. of Computer Science, Stanford Univ., Stanford, California, January 1996. Accessible at <http://www-diglib.stanford.edu>.