

# Semantic and schematic similarities between database objects: a context-based approach

Vipul Kashyap<sup>1,2,\*</sup>, Amit Sheth<sup>2</sup>

<sup>1</sup> Department of Computer Science, Rutgers University, New Brunswick, NJ 08903, USA

<sup>2</sup> LSDIS, Department of Computer Science, University of Georgia, 415 GSRC, GA 30602-7404, USA

Edited by Dennis McLeod. Received 3 November 1993 / Revised 30 March 1994 / Accepted 1 August 1995

**Abstract.** In a multidatabase system, schematic conflicts between two objects are usually of interest only when the objects have some semantic similarity. We use the concept of *semantic proximity*, which is essentially an *abstraction/mapping* between the domains of the two objects associated with the *context of comparison*. An explicit though partial context representation is proposed and the specificity relationship between contexts is defined. The contexts are organized as a meet semi-lattice and associated operations like the greatest lower bound are defined. The context of comparison and the type of abstractions used to relate the two objects form the basis of a semantic taxonomy. At the *semantic level*, the intensional description of database objects provided by the context is expressed using description logics. The terms used to construct the contexts are obtained from *domain-specific ontologies*. *Schema correspondences* are used to store mappings from the semantic level to the data level and are associated with the respective contexts. Inferences about database content at the federation level are modeled as changes in the context and the associated schema correspondences. We try to reconcile the dual (schematic and semantic) perspectives by enumerating *possible semantic similarities* between objects having schema and data conflicts, and modeling schema correspondences as the projection of semantic proximity *with respect to (wrt)* context.

## 1 Introduction

Many organizations face the challenge of interoperating among multiple independently developed database systems to perform critical functions. With high interconnectivity and access to many information sources, the primary issue in the future will not be how to efficiently process the data that is known to be relevant, but to determine which data is relevant [She91]. Thus, the fundamental question in interoperability is that of identifying objects in different databases that are

semantically related, and then resolving the schematic differences among these objects. In this paper, we are interested in the reconciliation of the semantic and schematic perspectives and its use as a step towards *information focusing* and *correlation* across multiple databases.

We characterize the degree of semantic similarity between a pair of objects using the concept of **semantic proximity** [SK92]. It is based on the premise that it is essential to associate the **abstractions/mappings** between the objects with the **context of comparison** for capturing the semantic similarity between them. Other researchers in the field of multidatabases have also made observations that are similar in principle, but different in details [ON93, SSR92, YSDK91]. This association of context with abstractions represents the first step in achieving the reconciliation between the semantic and schematic perspectives.

Inadequacies of purely structural and mapping-based methods are discussed, and explicit representation of context is proposed to resolve some inadequacies. Computational benefits of representing context are also discussed. We propose a partial representation of context as a collection of contextual coordinates and their values. This representation is used to describe objects and the constraints which they must satisfy in an intensional manner. The meaning of the contextual coordinates and their values are informally explained by expressing the context using description logic (DL) expressions [BS85].

In order for a context representation to be useful for semantic interoperability in multidatabases, it is important to have automatic ways of comparing and manipulating them. Based on the proposed representation of context, we define the specificity relationship between two contexts. A definition of the specificity relationship and the *greatest lower bound* (glb) and other operations on contexts are presented. The specificity relationship induces a partial order such that, for any two contexts, there exists a *glb* leading to the organization of the context set as a meet semi-lattice.

The semantic proximity descriptor consists of context and abstraction as its main components. Depending on the values assumed by these two components, we define a data model-independent taxonomy of semantic similarities. The possible values of the first component can be contexts con-

\* Presently at: MCC, 3500 W. Balcones, Center Dr., Austin, Texas, 78759 USA

structured using the various operations mentioned above. Classification or taxonomies of *schematic differences* appear in multidatabase literature. However, purely schematic considerations do not suffice to determine the similarity between objects [FKN91, SG89]. We try to reconcile the two perspectives by enumerating the possible semantic similarities between objects having schematic and data conflicts.

Even though the representation of semantics better enables us to represent the similarities between the various objects, we also need to be able to capture structural similarities in a mathematical formalism for reasoning on the computer. We define the concept of **schema correspondences** to capture the structural similarities between the objects. They are also associated with the context in which the semantic proximity is defined. We reconcile the semantic and schematic perspectives by defining the schema correspondence as a projection of the semantic proximity *with respect to* (wrt) context. The semantics of the projection operation are captured in the rules of the algebra enumerated in Appendix 1.

The overall organization of the paper is as follows. In Sect. 2, we present a model to represent semantic similarities among objects. In Sect. 3, we discuss the rationale for representation of context in a multidatabase environment and propose an explicit, though partial, representation of context. The associated operations for reasoning about and manipulating the context representations are also defined. In Sect. 4, a taxonomy of the various types of possible semantic similarities between the various objects is presented. In Sect. 5, we discuss a broad class of schematic differences and the possible semantic similarities between objects having those differences. In Sect. 6, we define a uniform formalism for representation of structural similarity. It is associated with the context and is defined as the projection of semantic similarity. Examples illustrating the operations from an algebra describing the projection operation (Appendix 1) are presented. A discussion of related work is presented in Sect. 7. Conclusions and future work are presented in Sect. 8.

## 2 Semantic similarities between objects

In this section, we discuss the concept of *semantic proximity* which characterizes **semantic** similarities between objects. We distinguish between the *real world* and the *model world* which is a representation of the real world. As in the work in semantic data modeling [HK87, PM88], we endeavor to capture some of the important semantic information about the real world and represent it in the model world. However, over and above the semantics of the data, we also attempt to capture semantics of queries and applications. This enables us to support semantics-based focusing and correlation of information across multiple databases *with respect to* an application.

Attempts have been made to capture the similarity of objects by using mathematical tools like value mappings between domains and abstractions such as generalization, aggregation, etc. However, it is our belief that the *real-world semantics* (RWS) of an object<sup>1</sup> cannot be captured sufficiently using mathematical formalisms. The term “object”

<sup>1</sup> The term “real-world semantics” distinguishes from the “(model) semantics” that can be captured using the abstractions in a semantic data

in this paper refers to an object in the model world (i.e., a representation or intensional definition in the model world, e.g., an object class definition in object-oriented models or relation in relational models) as opposed to an entity or a concept in the real world. These objects may model information at any level of representation, such as the *attribute* or *entity* level.

We need to understand and represent more knowledge to capture the semantics of relationships between objects. This knowledge should be able to capture the **context** of comparison of the objects and the **abstraction** relating the domains of the two objects. Attempts to partially represent such extra knowledge include the use of meta-attributes [SSR92] and building and partitioning ontologies into micro-theories [Guh90].

Attempts to represent context and abstraction as suggested above have been reflected in the techniques and representational constructs used by various practitioners and researchers in the field of multidatabases. The model for semantic proximity defined in this section has been influenced by these attempts. Some significant attempts are the **semantic proximity** proposal by Sheth and Kashyap [SK92], the **context building** approach by Ouksel and Naiman [ON93], the **context interchange** approach by Sciore et al. [SSR92] and the **common concepts** approach by Yu et al. [YSDK91]. We relate the above attempts to semantic proximity.

### 2.1 Semantic proximity: a model for semantic similarity

Given two objects  $O_1$  and  $O_2$ , the *semantic proximity* between them is defined by the 4-tuple given by [SK92]:

**semPro** ( $O_1, O_2$ )

= < **Context**, **Abstraction**, ( $D_1, D_2$ ), ( $S_1, S_2$ ) > ,

where  $D_i$  is domain of  $O_i$  and  $S_i$  is state of  $O_i$ .

- The first component denotes the context in which the two objects  $O_1$  and  $O_2$  are being compared. This context may be the same, different, or related in some manner to the context(s) in which the objects  $O_1$  and  $O_2$  are defined.
- The second component identifies the abstraction/mapping used to relate the domains of the objects,  $O_1$  and  $O_2$ .
- The third component enumerates the domain definitions of the objects,  $O_1$  and  $O_2$ . The domains may be defined by either enumerating the values as a set or by using existing type definitions in the database.
- The fourth component enumerates the states of the objects, which are the extensions of the objects recorded in their respective databases at a particular time.

In Fig. 1 we have illustrated the definition of the semantic proximity between two objects  $O_1$  and  $O_2$  in the database. Context( $O_1$ ) and context( $O_2$ ) represent the contexts (referred to as *definition contexts* later in the paper) in which the objects  $O_1$  and  $O_2$  are mapped from the real world to the model world. Context refers to the context in which the objects are being compared.

model. Our definition is also intensional in nature and differs from the extensional definition of Elmasri et al. [ELN86] who define the RWS of an object to be the set of real-world objects it represents.

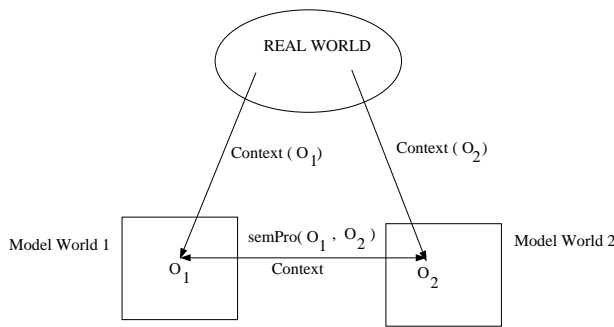


Fig. 1. Semantic proximity between two objects

## 2.2 Context: the semantic component

The context is the key component in capturing the semantics related to an object's definition and its relationships to other objects. Alternatives discussed in the multidatabase literature for representing context are as follows.

- In [ON93], context is defined as the knowledge that is needed to reason about another system, for the purpose of answering a query. It is specified as a set of assertions identifying the correspondences between various schema elements.
- In [SSR92], context is defined as the meaning, content, organization and properties of data. It is modeled using meta-data associated with the data.
- In [YSDK91], *common concepts* are proposed to characterize similarities between attributes in multiple databases.
- When using a well-defined ontology, such as Cyc [Guh90], a well-defined partition (called *Microtheory*) of the ontology is assigned a context.
- A context may be identified or represented using the following [SK92].
  - by association with a database or a group of databases
  - as the *relationship* in which an entity participates
  - from a schema architecture (e.g., the multidatabase or federated schema architecture of [SL90]), a context can be specified in terms of an *export schema* (a context that is closer to the database) or an *external schema* (a context that is closer to the application)
  - at a very elementary level, as a *named collection* of domains of objects

A context may be used in several ways to capture the relevant semantics. A context may be associated with an object to specify the assumptions used in its design and its relationships with other objects. However, the term context in semPro refers to the context in which a particular semantic similarity holds between two objects. As we shall see later, the context in semPro need not be the exactly the same as the contexts associated with the objects.

## 2.3 Abstractions/mappings: the structural component

We use the term abstraction to refer to the relation between the domains of the two objects. Mapping between the domains of objects is the mathematical expression to denote

the abstractions. However, since abstractions by themselves cannot capture semantic similarity, they have to be associated either with the context [KS93] or with extra knowledge in order to capture the RWS. Some of the proposals are as follows.

- In [SK92], abstractions are defined in terms of value mappings between the domains of objects and are associated with the context as a part of the semantic proximity.
- In [ON93], mappings are defined between schema elements called *interschema correspondence assertions* or ISCA's. A set of ISCA's under consideration are a representation of the context for integration of the schemas.
- In [SSR92], mappings called *conversion functions* are associated with the meta-attributes which define the context.
- In [YSDK91], the attributes are associated with “common concepts”. Thus the mappings (relationship) between the attributes are determined through the extra knowledge associated with the concepts.

Some useful and well-defined abstractions are

*Total 1-1 value mapping.* For every value in the domain of one object, there exists a value in the domain of the other object and vice versa.

*Partial many-one mapping.* In this case, some values in the domain of one of the objects might remain unmapped, or a value in one domain might be associated with many values in another domain.

*Generalization/specialization.* One domain can generalize/specialize the other, or domains of both the objects can be generalized/specialized to a third domain.

*Aggregation.* One domain can be an aggregation or a collection of other domains.

*Functional dependencies.* The values of one domain might depend functionally on the other domain.

*ANY.* This is used to denote that any abstraction such as the ones defined above may be used to define a mapping between the domains of two objects.

*NONE.* This is used to denote that there is no mapping defined between the domains of two objects.

## 2.4 Domains of the objects

Domains refer to the sets of values from which the objects can take their values. When using an object-oriented model, the domains of objects can be thought of as types, whereas the collections of objects might themselves be thought of as classes. A domain can be either **atomic** (i.e., cannot be decomposed any further) or composed of other atomic or composite domains. The domain of an object can be thought of as a subset of the cross-product of the domains of the properties of the object (Fig. 2). Analogously, we can have other combinations of domains, such as union and intersection of domains.

An important distinction between a context and a domain should be noted. One of the ways to specify a context is as a named collection of the domains of objects, i.e., it is associated with a group of objects. A domain, on the other

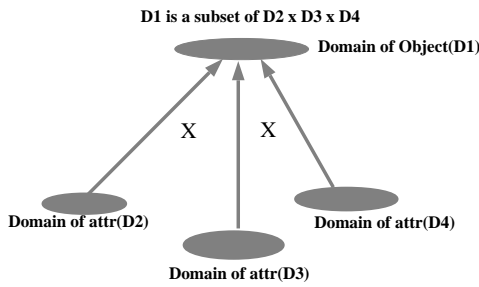


Fig. 2. Domain of an object and its attributes

hand is a property of an object and is associated with the description of that object.

### 2.5 States (extensions) of the objects

The state of an object can be thought of as an extension of an object recorded in a database or databases. However, this extension must not be confused with the actual state of the entity being modeled according to the RWS. Two objects having different extensions can have the same state RWS (and hence be semantically equivalent).

## 3 Explicit context representation in a multidatabase environment

In this section, we discuss the inadequacies of purely structural and mapping-based methods to represent object similarity and how representing context in the model world helps solve some of them. We also discuss computational advantages of representing context in the model world and propose an appropriate representation of context as a collection of contextual coordinates and their values. The contextual coordinates and their values may be chosen from a previously defined ontology of concepts.

We view ontology as the symbolic layer closest to concepts in the real world. An ontology may be defined as the specification of a representational vocabulary for a shared domain of discourse which may include definitions of classes, relations, functions and other objects [Gru93]. Criteria for constructing contexts from an ontology are discussed in [KS95a].

We discuss a partial representation of context, the semantics of which are informally explained using DL expressions. We shall also define operations for automatic ways of comparing (e.g., deciding whether one context is more general than the other) and manipulating contexts (e.g., taking the glb of two contexts). A brief discussion of issues relating to the language for representing contexts and the domain-specific ontologies from which terms to construct contexts are obtained is also presented.

### 3.1 Rationale for context representation

In characterizing the similarity between objects based on the semantics associated with them we have to consider the

RWS of an object. It is not possible to completely define what an object denotes or means in the model world [SG89]. We propose the **context** of an object as the primary vehicle to capture the RWS of the object. The context in which two objects are being compared and the associated abstraction/mapping helps to capture the semantic aspect of the relationship between two objects (Fig. 1). We argue for the need for representing context by showing the inadequacy of purely structural representations. We also discuss the computational benefits of representing context.

#### 3.1.1 Inadequacy of purely structural representations

It has been suggested by Sheth and Gala, Kashyap [SG89, KS94b] and Fankhauser et al. [FKN91] that the ability to represent the structure of an object does not help capture the RWS of the object. It is not possible to provide a structural and hence a mathematical definition of the complex notion of RWS. In [LNE89], a one-to-one mapping is assumed between the attribute definition and the attribute's RWS. They define an attribute in terms of fixed descriptors such as *Uniqueness*, *Lower/Upper Bound*, *Domain*, *Scale*, etc., which are used to generate mappings between two attributes. They are also used to determine the equivalence of attributes. However, what they establish is the structural equivalence of these attributes which is necessary but not sufficient to determine the semantic equivalence of the attributes.

Consider two attributes *person-name* and *department-name*. We may be able to define a mapping between the domains of these two attributes, but we know that they are not semantically equivalent. In order to be able to capture this lack of equivalence, we propose the mappings between the domains of the attributes be made *with respect to* a context. We define two objects to be semantically equivalent if it is possible to define mappings *with respect to* all known and coherent contexts. The respective definition contexts should be coherent *with respect to* each other. Definition contexts and the notion of coherence is defined later in this section. Since the definition contexts of *person-name* and *department-name* are not coherent (one identifies an animate and the other identifies an inanimate object), they are not defined as equivalent attributes.

#### 3.1.2 Computational benefits of representing context

In [Sho91], Shoham discussed the computational benefits that might accrue in modeling and representing context in AI and knowledge-based systems. We believe that there are similarities between AI/knowledge-based and multidatabase systems that suggest context representation in a multidatabase system for a clean and efficient handling of information.

*Economy of representation.* In a manner akin to database views, contexts can act as a *focusing mechanism* when accessing the component databases of a multidatabase system. They can be a *semantic summary* of the information in a database or group of databases and may be able to capture semantic information which cannot be

expressed in the data definition model of the databases. Thus, unnecessary details can be abstracted from the user. Examples detailing this are enumerated in Sect. 6.2.

*Economy of reasoning.* Instead of reasoning with the information present in the database as a whole, reasoning can be performed with the context associated with a database or a group of databases. This approach has been used in [KS94a, MKSI96] for information resource discovery and query processing.

*Handling inconsistent information.* In a multidatabase system, where databases are designed and developed independently, it is not uncommon to have information in one database inconsistent with information in another. As long as information is consistent within the context of the query of the user, inconsistency in information from different databases may be allowed. This is discussed in Sect. 5.3.

*Flexible semantics.* A big fallout of associating abstractions/mappings with the context in the semantic proximity model (Sect. 2.1) is that the same two objects can be related to each other differently in two different contexts. This is because two objects might be semantically closer to each other in one context than in the other.

### 3.2 A partial context representation

There have been attempts to represent the similarity between two objects in databases. In [LNE89], a fixed set of descriptors define essential characteristics of the attribute and are used to generate mappings between them. We have discussed with the help of an example how they do not guarantee semantic similarity. Thus, any representation of context which can be described by a fixed set of descriptors is not appropriate.

The descriptors (or meta-attributes) are not fixed but dynamically chosen to model the characteristics of the application domain in question. It is not possible a priori to determine all possible meta-attributes which would completely characterize the semantics of the application domain. This leads to a *partial* representation of context. We represent context as a collection of contextual coordinates (meta-attributes) as follows:

$$\text{Context} = \langle (C_1, V_1) (C_2, V_2) \dots (C_k, V_k) \rangle$$

We shall informally explain the meaning of the symbols  $C_i$  and  $V_i$  by using examples and by enumerating the corresponding DL expressions (Table 1). Using DL expressions<sup>2</sup>, it is possible to define primitive classes and, in addition, specify classes using intensional descriptions phrased in terms of necessary and sufficient properties that must be satisfied by their instances. The intensional descriptions may be used to express the collection of constraints that make up a context. Also, each  $C_i$  roughly corresponds to a role and each  $V_i$  roughly corresponds to fillers for the role the object must have.

<sup>2</sup> We have proposed a minor addition [ $\langle \text{role-set} \rangle$ ] for  $\langle \text{DL-expression} \rangle$  [MKSI96]. However this is for retrieval only and not used for concept forming.

- $C_i$ ,  $1 \leq i \leq k$ , is a contextual coordinate denoting an aspect of context.
- $C_i$  may model some characteristic of the subject domain and may be obtained from a domain-specific ontology (discussed later in this section).
- $C_i$  may model an implicit assumption in the design of a database.
- $C_i$  may or may not be associated with an attribute  $A_j$  of an object  $O$  in the database.

The value  $V_i$  of a contextual coordinate  $C_i$  can be represented in the following manner:

- $V_i$  can be a variable.
  - It is used only at the highest level of nesting for retrieval of objects/properties.
  - It can be unified (in the sense of Prolog) with another variable, a set of symbols, an object or type defined in the database or another variable.
  - It can be unified with another variable associated with a context.
  - It can be used to impose constraints on the answer.

*Example.* Suppose we are interested in people who are authors and who hold a post. We can represent the query context  $C_q$  (discussed later in this section) as follows:

$$C_q = \langle (\text{author}, X) (\text{designee}, X) \rangle$$

The same thing can be expressed in a DL as follows:

$$C_q = [\text{author}] \text{ for } (\text{SAME-AS author designee})$$

- $V_i$  can be a set.
  - The set may be an enumeration of symbols from a domain-specific ontology.
  - The set may be defined as the extension of an object or as elements from the domain of a type defined in the database.
  - The set may be defined by posing constraints on pre-existing sets.

*Example.* Suppose we want to represent the assumptions implicit in the design of the object EMPLOYEE in a database. We can represent this as the definition context of EMPLOYEE,  $C_{def}(\text{EMPLOYEE})$  as follows:

$$C_{def}(\text{EMPLOYEE})$$

$$= \langle (\text{employer}, [\text{Deptypes} \cup \{\text{restypes}\}]) (\text{article}, \text{PUBLICATION}) \rangle$$

Let Deptconcept = term corresponding to Deptypes in an ontology

The same thing can be expressed in a DL as follows:

$$C_{def}(\text{EMPLOYEE})$$

$$= (\text{AND EMPLOYEE}$$

$$(\text{ALL article PUBLICATION})$$

$$(\text{ALL employer (OR Deptconcept)}$$

$$(\text{ONE-OF research })))$$

**Deptypes** is a type defined in the database. The symbols restypes, employer and article are taken from the ontology. The definition context (defined later in this section) expresses an association between EMPLOYEE and PUBLICATION which may not be captured in the database.

- $V_i$  can be a variable associated with a context.
  - This can be used to express constraints which the result of a query should obey. This is called the constraint context and is defined later in this section.

**Table 1.** Contextual coordinate, value pairs and the corresponding DL expressions

Contextual coordinates and values, $C_{def}(O)$ , $C_q$	DL expressions
$C_{def}(O) = \langle (C_1, V_1) \dots (C_k, V_k) \rangle$	$(\text{AND } O (\text{ALL } C_1 V_1) \dots (\text{ALL } C_k, V_k))$
$C_{def}(O) = \langle (C_i, O_i \circ \langle (C_j, V_j) \rangle) \rangle$	$(\text{AND } O (\text{ALL } C_i (\text{AND } O_j (\text{ALL } C_j V_j))))$
$C_q = \langle (C_i, X) (C_j, X) \rangle$	$[C_i]$ for $(\text{SAME-AS } C_i C_j)$
$C_q = \langle (C_i, X \circ \langle (C_j, V_j) \rangle) \rangle$	$[C_i]$ for $(\text{ALL } C_i (\text{ALL } C_j V_j))$

- The constraints would apply to the set, type or object the variable X would unify with.

*Example.* Suppose we want all the articles whose titles contain the substring “abortion” in them. This can be expressed in the following query context:

$$C_q = \langle (\text{article}, X \circ \langle (\text{title}, \{y | \text{substring}(y) = \text{“abortion”}\}) \rangle) \rangle \\ = \langle (\text{article}, X \circ \text{Cntxt}) \rangle$$

where  $\circ$  denotes association of a context with a variable X and

$$\text{Cntxt} = \langle (\text{title}, \{y | \text{substring}(y) = \text{“abortion”}\}) \rangle$$

*Association* of a variable and a context ensures that the answer satisfies the constraints expressed in the context.

The same thing can be expressed in a DL as follows:

$$\text{Let Extension (AString)} = \{y | \text{substring}(y) = \text{“abortion”}\} \\ C_q = [\text{article}] \text{ for } (\text{ALL } \text{article} (\text{ALL } \text{title AString}))$$

- $V_i$  can be a set, type or an object associated with a context.

- This is called the association context and is defined later in this section.
- This may be used to express semantic dependencies between objects which may not be modeled in the database.

*Example.* Suppose we want to represent information relating publications to employees in a database. Let PUBLICATION and EMPLOYEE be objects in a database. The definition context of HAS-PUBLICATION can be defined as:

$$C_{def}(\text{HAS-PUBLICATION}) \\ = \langle (\text{article}, \text{PUBLICATION}) \\ (\text{author}, \text{EMPLOYEE} \circ \langle (\text{affiliation}, \{\text{research}\}) \rangle) \rangle$$

$$C_{def}(\text{HAS-PUBLICATION}) \\ = \langle (\text{article}, \text{PUBLICATION}) \\ (\text{author}, \text{EMPLOYEE} \circ \text{Cntxt}) \rangle$$

where  $\circ$  denotes association of a context with an object EMPLOYEE, and  $\text{Cntxt} = \langle (\text{affiliation}, \{\text{research}\}) \rangle$

*Association* of a context with an object is similar to defining a view on the object extensions such that only those instances satisfying the constraints defined in the context are exported to the federation. The same thing can be expressed in a DL as follows:

$$C_{def}(\text{HAS-PUBLICATION}) \\ = (\text{AND } \text{HAS-PUBLICATION} \\ (\text{ALL } \text{article PUBLICATION}) \\ (\text{ALL } \text{author } (\text{AND } \text{EMPLOYEE} \\ (\text{ALL } \text{affiliation } (\text{ONE-OF } \text{research}))))))$$

Note that the relationships between EMPLOYEE, PUBLICATION and HAS-PUBLICATION is information represented in the context not modeled in the database.

### 3.2.1 Definition context of an object

Given an object O in a database and a collection of contextual coordinates  $C_i$ s from the ontology, the definition context is denoted as  $C_{def}(O)$  and can be used in the following ways:

- to specify the assumptions used in the design of the object O
- to share only a pre-determined extension of the object O with the federation of databases. This exported object is denoted as  $O_F$

The associations between the objects stored in the database and the objects exported to the federation are expressed using the concepts of **semantic proximity** and **schema correspondences** (defined in Sect. 6.1).

### 3.2.2 Association context of objects

Given objects O and  $O_1$  in a database the dependence of the definition context of O on the context of association between O and  $O_1$ ,  $C_{ass}(O_1, O)$  can be represented as:

$$C_{def}(O) = \langle (C_1, O_1 \circ C_{ass}(O_1, O)) \dots (C_k, V_k) \rangle$$

The association context can be used in the following ways:

- to represent relationships between two objects with reference to an aspect of an application domain. This is done by associating it with the appropriate contextual coordinate
- different relationships between two objects may hold with reference to different aspects of the subject domain. This can be modeled by different association contexts between the two objects associated with different contextual coordinates
- to model the relationships between the object O and different (more than one) objects as a part of the definition context of the same object. Thus, the context of an object would consist of its relationships with other objects

### 3.2.3 Query context

Whenever a query Q is posed to a federation of databases, we associate with it a query context  $C_q$  which makes explicit the partial semantics of the query Q.

- The user can consult ontologies to construct the query context in a semi-automatic manner. Issues of combining and displaying ontologies to enable a user to do this easily are discussed in [MKSI96, MKIS96, KS96].
- Objects and types defined in databases are also available to the user by relating them to some concept in an ontology.

- The query is expressed as a set of constraints which an answer object must satisfy. The constraints expressed in the query context can express incomplete information.

### 3.2.4 Constraint context

The constraint context,  $C_{constr}(X, ANSWER)$  is typically a part of the query context and is used to pose constraints on the answer returned for the query.

$$C_q = \langle (C_1, X \circ C_{constr}(X, ANSWER)) \dots (C_k, V_k) \rangle$$

- It is associated with a variable which may be a placeholder for the answer or a part of the answer. The variable may be instantiated to an object or type definition.
- The context may represent constraints on the object and its attributes or the contextual coordinates associated with an object.
- The constraints which we currently limit to are cardinality constraints on sets and those that may be defined as a predicate on the elements of a set.

### 3.3 Reasoning about and manipulation of contexts

We have proposed a partial representation of context in the previous section. To use this representation meaningfully to focus on relevant information and to correlate information the following needs to be precisely defined:

- the most common relationship between contexts is the “specificity” relationship. Given two contexts  $C_1$  and  $C_2$ ,  $C_1 \leq C_2$  if  $C_1$  is at least as specific as  $C_2$ . This is useful when objects defined in a particular context have to transcend [McC93] to a more specific or general context. This is discussed in detail with examples in [KS95b].
- It is also the case that two contexts may not be comparable to each other, i.e. it may not be possible to decide whether one is more general than the other or not. Thus, the specificity relationship gives us a partial order.
- For every two contexts, we define the glb of two contexts as the most general context which is more specific than each of the two contexts. The set of contexts thus forms a meet semi-lattice.

#### 3.3.1 The specificity relationship

The specificity relationship between two contexts determines which context is more general than the other. We have defined this relationship with the help of specificity rules governing the contextual coordinates and their values.

$$\text{Let } Cntxt_1 = \langle (C_1, V_1) (C_2, V_2) \dots (C_k, V_k) \rangle$$

$$Cntxt_2 = \langle (C'_1, V'_1) (C'_2, V'_2) \dots (C'_m, V'_m) \rangle$$

$$Cntxt_1 \leq Cntxt_2 \text{ if } Cntxt_1 \text{ is at least as specific as } Cntxt_2$$

In the following exposition,  $C, C_1, C_2, C'_1, C'_2, \dots$  denote the contextual coordinates of the contexts under consideration.  $V, V_1, V_2, V'_1, V'_2, \dots$  denote the values of the contextual

coordinates.  $A, A_1, A_2, \dots, S, S_1, S_2, \dots$  stand for sets.  $X, Y, Z, \dots$  stand for variables.

The specificity rules for the values of the contextual coordinates ( $V_i$ s) are as follows:

**variable specificity:**  $V_1 \leq X$ , anything is more specific than a variable

**set specificity:**  $S_1 \leq S_2$  iff  $S_1 \subseteq S_2$

**association context specificity:** these are rules concerning specificity of contextual coordinates when an association context is involved.

- $A_1 \circ Cntxt_i \leq A_2$  if  $A_1 \leq A_2$
- $A_i \circ Cntxt_i \leq A_j \circ Cntxt_j$  if  $A_i \leq A_j \wedge Cntxt_i \leq Cntxt_j$

$Cntxt_1 \leq Cntxt_2$  if the following conditions hold:

- $m \leq k$
- $\forall i, 1 \leq i \leq m, \exists j C_j \leq C'_i \wedge V_j \leq V'_i$

#### 3.3.2 Operations on the context lattice

As observed earlier, the specificity relationship between the contexts induces a partial order among the contexts. Thus, the context can be organized as a meet semi-lattice where every pair of contexts has the glb. In this subsection, we define the *glb* operation and other operations we will use later in the paper.

$overlap(Cntxt_1, Cntxt_2) = \{ C_i \mid C_i \in Cntxt_1 \wedge C_i \in Cntxt_2 \}$   
 $coherent(Cntxt_1, Cntxt_2)$  This operator determines whether the constraints determined by the values of the contextual coordinates are consistent.

*Example.* Let  $Cntxt_1 = \langle (salary, \{x \mid x \leq 10000\}) \rangle$

$Cntxt_2 = \langle (salary, \{x \mid x > 10000\}) \rangle$

Thus,  $coherent(Cntxt_1, Cntxt_2) = FALSE$

##### 3.3.2.1 The glb of two contexts

We now define the glb of two contexts with the help of the rules that determine the glbs of the contextual coordinates and their values. The rules determining  $glb(V_i, V'_j)$  are

**Variable:**  $glb(V_i, X) = V_i$ ;

**sets:**  $glb(S_1, S_2) = S_1 \cap S_2$

**Association contexts.** these are rules concerning the glb of the values of the contextual coordinates when an association context is involved.

- $glb(A_1 \circ Cntxt_i, A_2) = glb(A_1, A_2) \circ Cntxt_i$
- $glb(A_i \circ Cntxt_i, A_j \circ Cntxt_j) = glb(A_i, A_j) \circ glb(Cntxt_i, Cntxt_j)$

The greatest lower bound of the contexts  $glb(Cntxt_1, Cntxt_2)$  can now be defined as:

<sup>3</sup> This specificity relationship between contextual coordinates is determined from the ontology and is beyond the scope of this paper. In defining the various operations on the context lattice we shall use the equality comparison instead.

- $\text{glb}(\text{Cntxt}_1, \text{Cntxt}_2) = \text{Cntxt}_1$ , if  $\text{Cntxt}_2 = \langle \rangle$   
[Empty Context]
- $(C_i, V_i) \in \text{glb}(\text{Cntxt}_1, \text{Cntxt}_2)$ ,  
if  $C_i \notin \text{overlap}(\text{Cntxt}_1, \text{Cntxt}_2)$
- $(C'_i, V'_i) \in \text{glb}(\text{Cntxt}_1, \text{Cntxt}_2)$ ,  
if  $C'_i \notin \text{overlap}(\text{Cntxt}_1, \text{Cntxt}_2)$
- $(C_k, \text{glb}(V_k, V'_j)) \in \text{glb}(\text{Cntxt}_1, \text{Cntxt}_2)$ ,  
if  $C_k = C'_j \in \text{overlap}(\text{Cntxt}_1, \text{Cntxt}_2)$

An alternative equivalent representation of a context (expressed using the  $\text{glb}$  operation) is very useful when there is a need to carry out inferences on the context and information associated with it.

$$\begin{aligned} \text{Cntxt} &= \langle (C_1, V_1)(C_2, V_2) \dots (C_k, V_k) \rangle \\ &= \text{glb}(\langle (C_1, V_1) \rangle, \text{glb}(\langle (C_2, V_2) \rangle, \dots, \\ &\quad \text{glb}(\langle (C_k, V_k) \rangle, \langle \rangle) \dots)) \end{aligned}$$

*Example.* Consider the following two contexts:

$$\begin{aligned} \text{Cntxt}_1 &= \langle (\text{author}, \text{EMPLOYEE} \circ \langle (\text{affiliation}, \{\text{research}\}) \rangle) \\ &\quad (\text{article}, \text{PUBLICATION}) \rangle \end{aligned}$$

$$\begin{aligned} \text{Cntxt}_2 &= \langle (\text{article}, \\ &\quad X \circ \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle) \rangle \end{aligned}$$

It should be noted that

- $\text{article} \in \text{overlap}(\text{Cntxt}_1, \text{Cntxt}_2)$   
 $\Rightarrow (\text{article}, \text{glb}(\text{PUBLICATION},$   
 $X \circ \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle))$   
 $\in \text{glb}(\text{Cntxt}_1, \text{Cntxt}_2)$

- $\text{author} \notin \text{overlap}(\text{Cntxt}_1, \text{Cntxt}_2) \Rightarrow$   
 $(\text{author}, \text{EMPLOYEE} \circ \langle (\text{affiliation}, \{\text{research}\}) \rangle)$   
 $\in \text{glb}(\text{Cntxt}_1, \text{Cntxt}_2)$

$$\begin{aligned} &– \text{glb}(\text{PUBLICATION}, \\ &\quad X \circ \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle) \\ &= \text{glb}(\text{PUBLICATION}, X \circ \\ &\quad \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle) \\ &\quad [\text{Association Contexts}] \\ &= \text{PUBLICATION} \circ \\ &\quad \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle \\ &\quad [\text{glb of a variable and an object}] \end{aligned}$$

$$\begin{aligned} &\text{glb}(\text{Cntxt}_1, \text{Cntxt}_2) \\ &= \langle (\text{author}, \text{EMPLOYEE} \circ \langle (\text{affiliation}, \{\text{research}\}) \rangle) \\ &\quad (\text{article}, \text{PUBLICATION} \circ \\ &\quad \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle) \rangle \end{aligned}$$

### 3.4 Issues of language and ontology in context representation

In this section, we discuss the issues of a language in which the explicit representation discussed above can be best expressed. We also discuss issues of ontology, i.e., the vocabulary used by the language to represent the contexts.

#### 3.4.1 Language for context representation

In Sect. 3.2, we have proposed a context representation as a collection of contextual coordinates and their values. The values themselves may have contexts associated with them. In this section, we enumerate the properties desired of a language to express the context representation.

- The language should be declarative in nature, as the context shall typically be used to express constraints on objects in an intensional manner. Besides, the declarative nature of the language will make it easier to perform inferences on the context.
- The language should be able to express the context as a collection of contextual coordinates, each describing a specific aspect of information present in the database or requested by a query.
- The language should have primitives (for determining the subtype of two types, pattern matching, etc.) in the model world, which might be useful in comparing and manipulating context representations.
- The language should have primitives to perform navigation in the ontology to identify the abstractions related to the ontological objects in the query context or the definition contexts of objects in the databases.

#### 3.4.2 The ontology problem

In constructing the contexts as illustrated in Sect. 3.2, the choice of the contextual coordinates ( $C_i$ s) and the values assigned to them ( $V_i$ s) is very important. There should be *ontological commitments*, i.e., agreements about the ontological objects used between the users and the information system designers. In our case, this corresponds to an agreement on the terms used for the contextual coordinates and their values by a user in formulating the query context and a database administrator for formulating the definition and association contexts. In an example in Sect. 3.2, we have defined  $C_{def}(\text{EMPLOYEE})$  by making use of symbols like *employer*, *affiliation* and *reimbursement* from the ontology for contextual coordinates and *research*, *teaching* etc., for the values of the contextual coordinates.

We assume that each database has available an ontology corresponding to a specific domain. The definition and association contexts of the objects take their terms and values from this ontology. However, in designing the definition contexts and the query context, the issues of combining the various ontologies arise.

We now enumerate various approaches one might take in building ontologies for a federation of information sources. Other than the ontological commitment, a critical issue in designing ontologies is the **scalability** of the ontology as more information sources enter the federation.

#### – The common ontology approach.

- One approach has been to build an extensive global ontology. A notable example of global ontology is Cyc [LG90], consisting of around 100,000 objects. The mapping between each individual information resource and the Cyc global ontology in the Carnot project [CHS91] is accomplished by a set of *articulation axioms* which are used to map the entities of an information resource to the concepts in Cyc's existing ontology [CHS91].
- Another approach has been to exploit the semantics of a single problem domain (e.g., transportation planning) [ACHK93]. The domain model is a declarative description of the objects and activities possible in



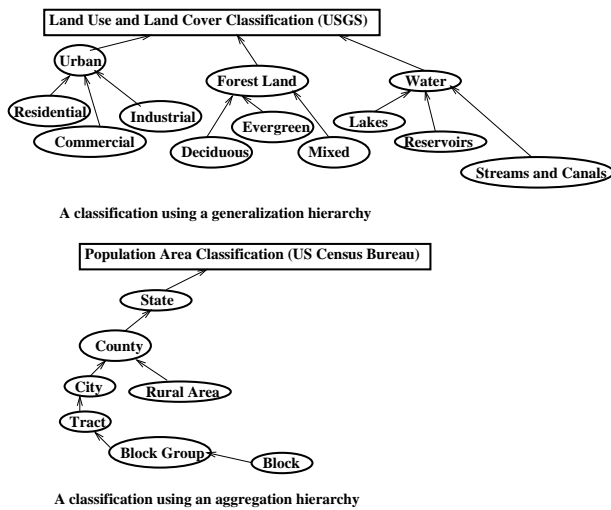


Fig. 3. Examples of generalization and aggregation hierarchies for ontology construction

the application domain as viewed by a typical user. The user formulates queries using terms from the application domain.

- **Re-use of existing ontologies.** Given our assumption that there will be numerous information systems participating in the federation, it is unrealistic to expect any one existing ontology or classification to suffice. We propose a re-use of various existing classifications such as ISBN classification for publications, botanical classification for plants, etc. An example of such a classification is illustrated in Fig. 3. These ontologies can then be combined in different ways and made available to the federation.
  - A critical issue in combining the various ontologies is determining the overlap between them. One possibility [Wie94] is two define the “intersection” and “mutual exclusion” points between the various ontologies. Attempts have been made to use terminological relationships between terms across different ontologies to represent the intersection points. In the OBSERVER system, *synonyms* have been used to represent the intersection points and a proposal to extend the system using *hyponyms* and *hypernyms* has been presented in [MKIS96].
  - Another approach has been adopted in [MS95]. The types determined to be similar by a sharing advisor are classified into a collection called *concept*. A *concept hierarchy* is thus generated modeling super-concept-subconcept relationships. These types may be from different databases and their similarity or dissimilarity is based on heuristics with user input as required.

#### 4 A semantic taxonomy

Our emphasis is on identifying semantic similarity even when the objects have significant representational differences [She91]. Semantic proximity is an attempt to characterize the degree of semantic similarity between two objects using the RWS. It provides a qualitative measure to distinguish

between the terms introduced in [She91], such as *semantic equivalence*, *semantic relationship*, *semantic relevance* and *semantic resemblance*. Two objects can be semantically related in one of the above four ways. Semantic equivalence is *semantically closer* than semantic relationship, and so on.

In this section, we use the concept of semantic proximity defined in Sect. 2 and the context representation discussed above to define a semantic taxonomy consisting of the various types of semantic similarities between objects. The taxonomy thus designed is illustrated in Fig. 5.

##### 4.1 The role of context in semantic classification

The context, which is the pivot on which the semantic proximity depends, plays a key role in this taxonomy. Here we enumerate the possible values for context.

- ALL, i.e., the *semPro* between the objects is being defined *with respect to* all known and *coherent* comparison contexts. There should be coherence between the definition contexts of the objects being compared and between the definition contexts and the context of comparison.
- SOME, i.e., the *semPro* between the objects is being defined *with respect to* some context. This context may be constructed in the following ways.
  - GLB, i.e., the *glb* of the contexts of the two objects. Typically, we are interested in the *glb* of the context of comparison and the definition context of the object.
  - LUB, i.e., the least upper bound (*lub*)<sup>4</sup> of the contexts of the two objects is taken. Typically, we are interested in the *lub* of the definition contexts of the two objects when there does not exist an abstraction/mapping between their domains in the context of comparison.
- SUB-CONTEXTS, we might be interested in the *semPro* between two objects in contexts which are more specific or more general *with respect to* the context of comparison.
- NONE, i.e., there does not exist a context in which a meaningful abstraction or mapping between the domains of the objects may be defined. This is the case when the definition contexts of the objects being compared are *not coherent* with each other.

##### 4.2 Semantic equivalence

This is the strongest measure of semantic proximity two objects can have. Two objects are defined to be *semantically equivalent* when they represent the same real-world entity or concept. Expressed in our model, it means that given two objects  $O_1$  and  $O_2$ , it should be possible to define a total 1-1 value mapping between the domains of these two objects in any known and coherent context. Thus we can write it as:

<sup>4</sup> We have not defined it for the general case. Here, we are only interested in the special case:  
 $(C_k, V_k \cup V'_j) \in \text{lub}(\text{Cntxt}_1, \text{Cntxt}_2)$   
 where  $C_k = C'_j \in \text{overlap}(\text{Cntxt}_1, \text{Cntxt}_2)$

$\text{semPro}(O_1, O_2)$

$= \langle \text{ALL, total 1-1 value mapping, } (D_1, D_2), \_ \rangle^5$ .

The notion of equivalence described above depends on the definition of the domains of the objects and can be more specifically called *domain semantic equivalence*. We can also define a stronger notion of semantic equivalence between two objects, which incorporates the state of the databases to which the two objects belong. This equivalence is called *state semantic equivalence* and is defined as:

$\text{semPro}(O_1, O_2) = \langle \text{ALL, M, } (D_1, D_2), (S_1, S_2) \rangle$ ,

where M is a total 1-1 value mapping between  $(D_1, S_1)$  and  $(D_2, S_2)$ .

For the purposes of this paper we shall use semantic equivalence to mean domain semantic equivalence.

### 4.3 Semantic relationship

This type of semantic similarity is weaker than semantic equivalence. Two objects are said to be *semantically related* when there exists a partial many-one value mapping, or a generalization, or aggregation abstraction between the domains of the two objects. Here, we relax the requirement of a 1-1 mapping in a way that, given an instance  $O_1$ , we can identify an instance of  $O_2$ , but not vice versa. The requirement that the mapping be definable in all the known and coherent contexts is not relaxed. Thus, we define the *semantic relationship* as:

$\text{semPro}(O_1, O_2) = \langle \text{ALL, M, } (D_1, D_2), \_ \rangle$ ,

where M may be a partial many-one value mapping, generalization, or aggregation

### 4.4 Semantic relevance

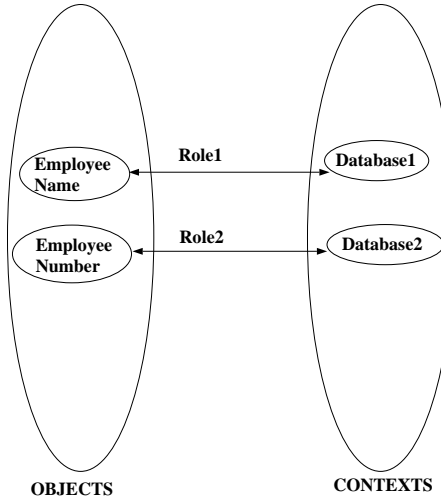
We consider two objects to be *semantically relevant* if they can be related to each other using some *abstraction* in *some context*. Thus the notion of semantic relevance between two objects is context-dependent, i.e., two objects may be semantically relevant in one context, but not so in another. Objects can be related to each other using any abstraction.

$\text{semPro}(O_1, O_2) = \langle \text{SOME, ANY, } (D_1, D_2), \_ \rangle$

### 4.5 Semantic resemblance

This is the weakest measure of semantic proximity, which might be useful in certain cases. Here, we consider the case where the domains of two objects cannot be related to each other by any abstraction in any context. Hence, the exact nature of semantic proximity between two objects is very difficult to specify. In this case, the user may be presented with extensions of both the objects. In order to express this type of semantic similarity, we introduce an aspect of context, which we call **role**, by extending the concept of role defined in [EN89]. Semantic resemblance is defined in detail in the next section.

<sup>5</sup> We use the “\_” sign to denote “don’t care”.



Role1 = role-of(EmployeeName, Database1) = Identifier  
 Role2 = role-of(EmployeeNumber, Database2) = Identifier  
 EmployeeName in Database1.Identifier  
 EmployeeNumber in Database2.Identifier  
 Thus, Role1 = Role2

Fig. 4. Roles played by objects in their contexts

#### 4.5.1 Role played by an object in a context

This refers to the relationship between an object and the semantic context to which it belongs. We characterize this relationship as a binary function, which has the object and its context as the arguments and the name of the role as the value.

role-of : object  $\times$  context  $\rightarrow$  rolename

The mapping defined above may be multivalued, as it is possible for an object to have multiple roles in the same context.

Based on the representation of a context proposed in Sect. 3.2, we can express this by constructing the lub of the contexts. Consider the type **NUMBER** and the type **NAME** defined in the databases.

$C_{def}(\text{Database1}) = \langle (\text{Class}, \{\text{Employee}, \dots\}) \rangle$   
 $(\text{Identifier}, \{\text{Name}, \dots\}) \rangle$

$C_{def}(\text{Database2}) = \langle (\text{Class}, \{\text{Employee}, \dots\}) \rangle$   
 $(\text{Identifier}, \{\text{Number}, \dots\}) \rangle$

$\text{lub}(C_{def}(\text{Database1}), C_{def}(\text{Database2}))$   
 $= \langle (\text{Class}, \{\text{Employee}_1, \text{Employee}_2, \dots\}) \rangle$   
 $(\text{Identifier}, \{\text{Name}, \text{Number}, \dots\}) \rangle$

Thus role-of(Name,  $C_{def}(\text{Database1})$ )  
 $=$  role-of(Number,  $C_{def}(\text{Database2})$ ) = Identifier

Since Name, Number  $\in$  Identifier  
 $\wedge$  Identifier  $\in$   $\text{lub}(C_{def}(\text{Database1}), C_{def}(\text{Database2}))$

This is illustrated in Fig. 4.

#### 4.5.2 Roles and semantic resemblance

Whenever two objects cannot be related to each other by any abstraction in any context, but they are associated with

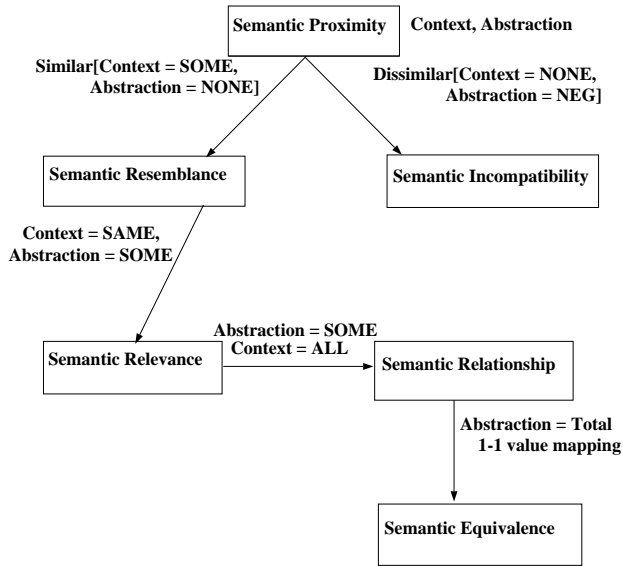


Fig. 5. Semantic classification of object similarities

contexts in which they have the same role and their definition contexts are coherent with respect to each other, they can be said to *semantically resemble* each other. This is a generalization of the DOMAIN-DISJOINT-ROLE-EQUAL concept in [LNE89].

$\text{semPro}(O_1, O_2) = \langle \text{SOME}(\text{LUB}), \text{NONE}, (D_1, D_2), - \rangle$ , where  $\text{coherent}(C_{def}(O_1), C_{def}(O_2))$  and  $\exists \text{Cntxt}_1, \text{Cntxt}_2$  exported by  $DB_1, DB_2$ , respectively and  $\text{SOME}(\text{LUB})$  denotes a context defined as follows:  $\text{context} = \text{lub}(\text{Cntxt}_1, \text{Cntxt}_2)$  and  $D_1 \neq D_2$  and  $\text{role-of}(O_1, \text{context}) = \text{role-of}(O_2, \text{context})$

#### 4.6 Semantic incompatibility

While all the qualitative proximity measures defined above describe semantic similarity, semantic incompatibility asserts semantic dissimilarity. Lack of any semantic similarity does not automatically imply that the objects are semantically incompatible. Establishing semantic incompatibility requires asserting that the definition contexts of the two objects are *incoherent with respect to* each other and there do not exist contexts associated with these objects such that they have the same role.

$\text{semPro}(O_1, O_2) = \langle \text{NONE}, \text{NONE}, (D_1, D_2), - \rangle$  where  $C_{def}(O_1)$  and  $C_{def}(O_2)$  are incoherent with each other and  $D_1$  may or may not be equal to  $D_2$  and  $\nexists$  context such that  $\text{role-of}(O_1, \text{context}) = \text{role-of}(O_2, \text{context})$

### 5 Schematic heterogeneities in multidatabases

In this section, we deal with a broad class of schematic differences and the possible semantic similarities between objects having schematic differences [SK92]. With each type of schematic difference, we enumerate the possible semantic

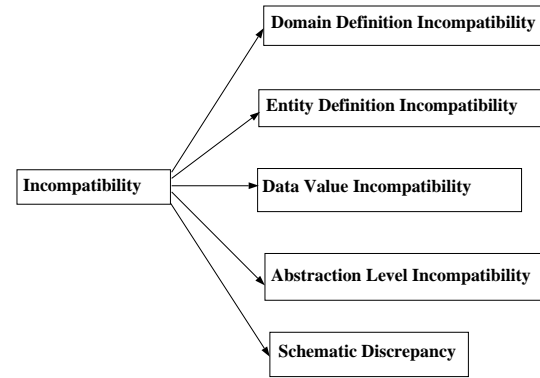


Fig. 6. Schematic heterogeneities

proximity descriptors. The broad classes of schematic heterogeneities we are dealing with are: *domain incompatibility*, *entity definition incompatibility*, *data value incompatibility*, *abstraction level incompatibility* and *schematic discrepancies* (Fig. 6). While the issues of schematic/representational/structural heterogeneity have been dealt with by a number of researchers [DH84, BOT86, CRE87, KLK91, KS91], the unique feature of our work is the strong correlation between the semantic aspects defined above and the structural aspects.

#### 5.1 Domain incompatibility

In this section, we discuss the incompatibilities that arise when two different domain types are used as different definitions of semantically similar attribute domains. We refine the broad definition of this incompatibility given in [CRE87]. We also discuss the possible semantic similarities with each case (Fig. 7).

##### 5.1.1 Naming conflicts

Two attributes that are semantically alike might have different names. They are known as *synonyms*.

*Example.* Consider two databases having the relations:

```

STUDENT(Id#, Name, Address)
TEACHER(SS#, Name, Address)
Id# of STUDENT and SS# of TEACHER are
synonyms.
  
```

Mappings between synonyms can often be established *with respect to* all known and coherent contexts. In such cases, the two domain types may be considered *semantically equivalent*.

Two attributes that are semantically unrelated might have the same names. They are known as *homonyms*.

*Example.* Consider two databases having the relations:

```

STUDENT(Id#, Name, Address)
BOOK(Id#, Name, Author)
Id# of STUDENT and BOOK are homonyms.
  
```

The definition contexts of the two domain types (which are defined in two different databases) may be modeled as follows:

$C_{def}(STUDENT.Id\#) = \langle(\text{identifies}, \text{AnimateObject})\rangle$   
 $C_{def}(BOOK.Id\#) = \langle(\text{identifies}, \text{InAnimateObject})\rangle$

The concepts `AnimateObject` and `InAnimateObject` are obtained from an ontology for the domain.

Since homonyms are semantically unrelated, their definition contexts are modeled in a way that they are *incoherent with respect to* each other. Thus, these two domain types may be considered *semantically incompatible*.

### 5.1.2 Data representation conflicts

Two attributes that are semantically similar might have different data types or representations.

*Example.*

`STUDENT.Id#` is defined as a 9-digit integer.  
`TEACHER.SS#` is defined as an 11-character string.

Conversion mappings or routines between different data representations can often be established *with respect to* all known and coherent contexts. In such cases, these domain types may be considered *semantically equivalent*.

### 5.1.3 Data scaling conflicts

Two attributes that are semantically similar might be represented using different units and measures. There is a 1-1 mapping between the values of the domains of the two attributes. For instance, the salary attribute might have values in \$ and \$.

Typically, mappings between data represented in different scales can be easily expressed in terms of a function or a lookup table, or by using dynamic attributes as in [LA86] and *with respect to* all known and coherent contexts. In such cases, the domain types may be considered *semantically equivalent*.

### 5.1.4 Data precision conflicts

Two attributes that are semantically similar might be represented using different precisions. This case is different from the previous case, because there may not be 1-1 mapping between the values of the domains. There may be a many-one mapping from the domain of the precise attribute to the domain of the coarser attribute.

*Example.*

Let the attribute `Marks` have an integer value from 1 to 100. Let the attribute `Grades` have the values {A, B, C, D, F}.

**Table 2.** Mapping between marks and grades

Marks	Grades
81-100	A
61-80	B
41-60	C
21-40	D
1-20	F

There may be a many-one mapping from `Marks` to `Grades` (Table 2). `Grades` is the coarser attribute. Typically, mappings can be specified from the precise data scale to the coarse data scale *with respect to* all known and coherent contexts. Given a letter grade, determining the precise numerical score is typically not possible. In such cases, the domain types may be considered *semantically related*.

### 5.1.5 Default value conflicts

This type of conflict depends on the definition of the domain of the concerned attributes. The *default value* of an attribute is that value which it is defined to have in the absence of more information about the real world. For instance, the default value for age of an adult might be defined as 18 years in one database and as 21 years in another.

It may not be possible to specify mappings between a default value of one attribute to the default value of another in all known and coherent contexts. However, it is often possible to do so *with respect to* some context. In such cases, the domain types can be considered to be *semantically relevant*, i.e., their *semantic proximity* can be defined as follows:

$\text{semPro}(\text{Age}_1, \text{Age}_2) = \langle\text{SOME}, \text{Abstraction}, (\text{D}_1, \text{D}_2), \text{-}\rangle$   
 $\text{Context} = \langle(\text{default}, \text{DefaultAge})\rangle,$

where the concept `DefaultAge` is obtained from an ontology for the domain. When the `semPro` is *evaluated* with respect to the context, it maps to different ages in the different databases.

### 5.1.6 Attribute integrity constraint conflicts

Two semantically similar attributes might be restricted by constraints which might not be consistent with each other. For instance, in different databases, the attribute `Age` might follow these constraints:

*Example.*

$C1 : \text{Age}_1 \leq 18$

$C2 : \text{Age}_2 > 21$

$C1$  and  $C2$  are inconsistent and hence the integrity constraints on the attribute `Salary` are said to conflict.

If the constraints are captured in the definition contexts of the domain types of `Age1` and `Age2`, then they would be incoherent and can be considered *semantically incompatible*. However, in the case these types are playing the same role in the definition contexts of their respective databases in which they exist, they may be considered to have a *semantic resemblance* to each other.

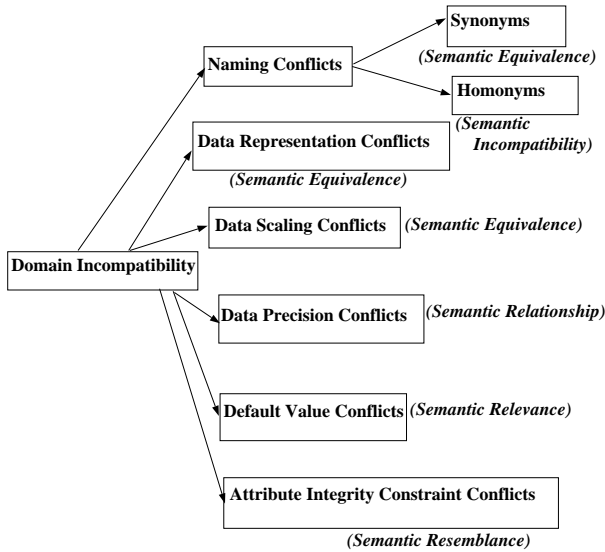


Fig. 7. Domain incompatibility and the likely types of semantic proximities

$$C_{def}(\text{Database}_1) = \langle (\text{timePeriod}, \{\text{Age}, \text{Duration}, \dots\}) \rangle,$$

$$C_{def}(\text{Database}_2) = \langle (\text{timePeriod}, \{\text{Age}, \text{RacePerformance}, \dots\}) \rangle,$$

where  $\text{Age}_1, \text{Age}_2$  denote the attribute Age in Database<sub>1</sub>, Database<sub>2</sub>, respectively

$$\text{semPro}(\text{Age}_1, \text{Age}_2) = \langle \text{SOME}(\text{LUB}), \text{NONE}, (D_1, D_2), \_ \rangle,$$

where SOME(LUB) denotes a context defined as follows: where  $\text{context} = \text{lub}(C_{def}(\text{Database}_1), C_{def}(\text{Database}_2))$  and  $D_1 \neq D_2$  and  $\text{role-of}(\text{Age}_1, \text{context}) = \text{role-of}(\text{Age}_2, \text{context}) = \text{timePeriod}$ .

## 5.2 Entity definition incompatibility

In this section, we discuss the incompatibilities that arise between two objects when the entity descriptors used by the objects are only partially compatible, even when the same type of entity is being modeled. We refine the broad definition of this class of conflicts given in [CRE87]. We also discuss the possible semantic similarities with each case (Fig. 8).

### 5.2.1 Database identifier conflicts

In this case, the entity descriptions in two databases are incompatible, because they use identifier records that are semantically different.

*Example.*

```

STUDENT1(SS#, Course, Grades)
STUDENT2(Name, Course, Grades)
STUDENT1.SS# and STUDENT2.Name are
semantically different keys.
  
```

The semantic proximity of objects having this kind of conflict depends on whether it is possible to define an abstraction to map the keys in one database to another. However, if we

assume that the context(s) of the identifiers are defined in the local schemas, we know that they play the *role of identification* in their respective contexts. Hence, the weakest possible measure of *semantic resemblance* applies, though stronger measures might apply too.

$\text{semPro}(\text{SS}\#, \text{Name}) = \langle \text{SOME}(\text{LUB}), \_ \_, (D_1, D_2), \_ \rangle,$  where  $D_1 = \text{Domain}(\text{SS}\#)$  and  $D_2 = \text{Domain}(\text{Name})$  and where SS# and Name exist in Database<sub>1</sub> and Database<sub>2</sub>, respectively

$$C_{def}(\text{Database}_1) = \langle (\text{Class}, \{\text{STUDENT1}, \dots\}) \rangle$$

$$C_{def}(\text{Database}_2) = \langle (\text{Class}, \{\text{STUDENT2}, \dots\}) \rangle$$

and SOME(LUB) denotes a context defined as follows:

$$\text{and context} = \text{lub}(C_{def}(\text{Database}_1), C_{def}(\text{Database}_2))$$

$$\text{and role-of}(\text{SS}\#, \text{context}) = \text{role-of}(\text{Name}, \text{context}) = \text{Identifiers}$$

### 5.2.2 Naming conflicts

Semantically alike entities might be named differently in different databases. For instance, EMPLOYEE and WORKERS might be two objects describing the same set of entities. They are known as *synonyms*. Typically, mappings between synonyms can often be established *with respect to* all known and coherent contexts. In such cases, the objects may be considered *semantically equivalent*.

On the other hand, semantically unrelated entities might have the same name in different databases. For instance, TICKETS might be the name of a relation which models movie tickets in one database, whereas it might model traffic violation tickets in another database. They are known as *homonyms* of each other. In a manner similar to that demonstrated in Sect. 5.1.1, their definition contexts can be modeled in a way that they are *incoherent with respect to* each other. Thus, these objects may be considered *semantically incompatible*.

### 5.2.3 Schema isomorphism conflicts

Semantically similar entities may have different number of attributes, giving rise to schema isomorphism conflicts.

*Example.*

```

INSTRUCTOR1(SS#, HomePhone, OffPhone)
INSTRUCTOR2(SS#, Phone)
is an example of schema non-isomorphism.
  
```

It should be noted that this can be considered an artifact of the *data precision conflicts* identified in Sect. 5.1.4 of this paper, as the phone number of INSTRUCTOR1 can be considered to be represented in a more precise manner than the phone number of INSTRUCTOR2. However, the conflicts discussed in section 5.1.4 are due to the differences in the domains of the attributes representing the same information and hence are *attribute level conflicts*. Whereas, conflicts in this section arise due to differences in the way the entities

INSTRUCTOR1 and INSTRUCTOR2 are defined in the two databases and hence are *entity level conflicts*.

Since mappings can be established between the objects on the basis of the common and identifying attributes, the two objects may be considered *semantically related*.

$\text{semPro}(\text{Instructor}_1, \text{Instructor}_2)$   
 $= \langle \text{ALL}, \{M_{ID}, M_1\}, \{D_{1,SS\#}, D_{HomePhone}, D_{OffPhone}\}, \{D_{2,SS\#}, D_{Phone}\}, - \rangle,$

where  $M_{ID}$  is a total 1-1 value mapping between  $D_{1,SS\#}$  and  $D_{2,SS\#}$  and represents the mapping between the identifiers of the two objects.

$M_1$  may be a total/partial 1-1/many-one value mapping between  $D_{HomePhone} \cup D_{OffPhone}$  and  $D_{Phone}$ .

#### 5.2.4 Missing data item conflicts

This conflict arises when, of the entity descriptors modeling semantically similar entities, one has a missing attribute. This type of conflict is subsumed by the conflict discussed in the previous section. A special case of the above conflict which satisfies the following conditions:

- the missing attribute is compatible with the entity, and
- there exists an inference mechanism to deduce the value of the attribute.

*Example.*

```
STUDENT(SS#, Name, Type)
GRAD-STUDENT(SS#, Name)
STUDENT.Type can have values "UG"
or "Grad"
GRAD-STUDENT does not have a type
attribute, but that can be implicitly
deduced to be "Grad".
```

In the above example, GRAD-STUDENT can be thought to have a type attribute whose default value is “Grad”. The conflict discussed in this section is different from the *default value* conflict in Sect. 5.1.5, which is an *attribute level conflict*, whereas the conflict discussed here is an *entity level conflict*. The objects may be considered *semantically relevant*, as proposed below.

The definition contexts of the two objects can be defined as:

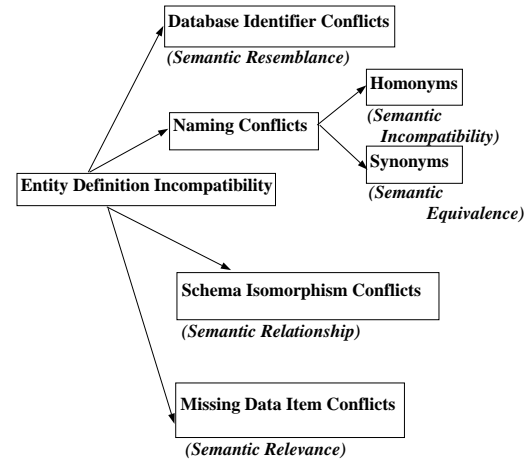
$C_{def}(\text{STUDENT}) = \langle (\text{type}, \{\text{graduate}, \text{undergraduate}\}) \rangle,$   
 $C_{def}(\text{GRAD-STUDENT}) = \langle (\text{type}, \{\text{graduate}\}) \rangle$

The context in which  $\text{semPro}(\text{STUDENT}, \text{GRAD-STUDENT})$  will be defined as:

$\text{glb}(C_{def}(\text{STUDENT}), C_{def}(\text{GRAD-STUDENT}))$   
 $= \langle (\text{type}, \{\text{graduate}\}) \rangle$

The abstraction is then computed by “conditioning” the original student abstraction *with respect to* this new context. Since every abstraction/mapping is associated with a context, the change in the abstraction as a result of the change in the associated context is called conditioning and is discussed in detail in [KS95b].

$\text{semPro}(\text{STUDENT}, \text{GRAD-STUDENT})$



**Fig. 8.** Entity definition incompatibilities and the likely types of semantic proximities

$= \langle \text{SOME}, M, (D_1, D_2), - \rangle,$

where  $M: \text{STUDENT} \rightarrow \text{GRAD-STUDENT}$  is a partial 1-1 value mapping

and Context = SOME =  $\langle (\text{type}, \{\text{graduate}\}) \rangle$

#### 5.3 Data value incompatibility

This class of conflicts covers those incompatibilities that arise due to the values of the data present in different databases [BOT86]. These conflicts are different from default value conflicts (Sect. 5.1.5) and attribute integrity constraint conflicts (Sect. 5.1.6) in that the latter are due to the differences in the definitions of the domain types of the attributes. Here, we refer to the data values already existing in the database. Thus, the conflicts here depend on the database state. Since we are dealing with independent databases, it is not necessary that the data values for the same entities in two different databases be consistent with each other. We also discuss the possible semantic similarities with each case (Fig. 9).

*Example.*

Consider two databases modeling the entity Ship

```
SHIP1(Id#, Name, Weight)
SHIP2(Id#, Name, Weight)
```

Consider an entity represented in both databases as follows:

```
SHIP1(123, USSEnterprise, 100)
SHIP2(123, USSEnterprise, 200)
```

Thus, we have the same entity for which SHIP1.Weight is not the same as SHIP2.Weight, i.e., it has inconsistent values in the database.

##### 5.3.1 Known inconsistency

In this type of conflict, the cause of inconsistency is known ahead of time and hence measures can be initiated to resolve

the inconsistency in the data values. For instance, it might be known ahead of time that one database is more reliable than the other. This information can typically be represented in the query context  $C_q$ . Here, the similarity of objects depends on the state component of semPro and are hence considered *state semantically relevant*.

$$C_q = \langle (\text{class, SHIP}) (\text{dataItem}, \{\text{Id}\#}) \\ (\text{choose-from}, \{\text{DB1}\}) \rangle$$

$\text{semPro}(O_1, O_2) = \langle C_q, M, (D_1, D_2), (S_1, S_2) \rangle$ ,  
where M is a total 1-1 value mapping between  $(D_1, S_1)$  and  $(D_2, S_2)$  (In this case the default is  $(D_1, S_1)$ ).

### 5.3.2 Temporal inconsistency

In this type of conflict, the inconsistency is of a temporary nature. This type of conflict has been identified in [RSK91] and has been expressed as a *temporal consistency predicate*<sup>6</sup>. One of the databases which has conflicting values might have obsolete information. This means that the information stored in the databases is time-dependent. The time lag information ( $\Delta t$ ) can be easily represented in the query context  $C_q$  and hence the objects may be considered *state semantically relevant*. The semPro when evaluated *with respect to context* gives the mapping defined below.

$$C_q = \langle (\text{class, SHIP}) (\text{dataItem}, \{\text{Weight}\}) (\text{timeLag}, \Delta t) \rangle \\ \text{semPro}(O_1, O_2) \\ = \langle C_q, \text{total 1-1 value mapping}, (D_1, D_2), (S_1, S_2) \rangle \\ \text{where } S_2(t + \Delta t) = S_1(t).$$

### 5.3.3 Acceptable inconsistency

In this type of conflict, the inconsistencies between values from different databases might be within an acceptable range. Thus, depending on the type of query being answered, the error in the values of two inconsistent databases might be considered tolerable. The *tolerance* of the inconsistency can be of a numerical or non-numerical nature and can be easily represented in the query context  $C_q$ , and hence the objects may be considered *state semantically relevant*.

*Example.* Numerical inconsistency

QUERY: Find the tax bracket of an employee.

INCONSISTENCY: If the inconsistency in the value of an employee income is up to a fraction of a dollar it may be ignored.

$$C_q = \langle (\text{class, EMPLOYEE}) (\text{dataItem}, \{\text{Salary}\}) \\ (\text{epsValue}, [0, 0.99]) \rangle,$$

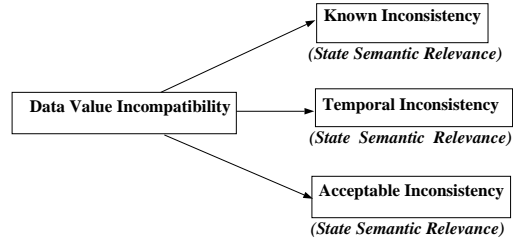
where epsValue is a contextual coordinate which models the level of inconsistency that can be tolerated for the query.

*Example.* Non-numerical inconsistency

QUERY: Find the state of residence of an employee.

INCONSISTENCY: If the employee is recorded as staying in Edison and New Brunswick (both are in New Jersey), then again the inconsistency may be ignored.

<sup>6</sup> Additional information on weaker criteria for consistency can be found in the literature on transaction models (e.g., see [SRK92]).



**Fig. 9.** Data value incompatibilities and the likely types of semantic proximities

$$C_q = \langle (\text{class,EMPLOYEE}) (\text{dataItem}, \{\text{Residence}\}) \\ (\text{epsValue}, \text{sameState}) \rangle$$

$\text{semPro}(O_1, O_2) \\ = \langle C_q, \text{partial many-one value mapping}, (D_1, D_2), (S_1, S_2) \rangle$ ,  
where  $\text{perturb}(S_1, \epsilon) = S_2$  and  $\epsilon$  is the discrepancy in the state of the two objects.

### 5.4 Abstraction level incompatibility

This class of conflicts was first discussed in [DH84] in the context of the functional data model. These incompatibilities arise when two semantically similar entities are represented at differing levels of abstraction. Differences in abstraction can arise due to the different levels of generality at which an entity is represented in the database. They can also arise due to aggregation used both at the entity as well as the attribute level. We also discuss the possible semantic similarities with each case (Fig. 10).

#### 5.4.1 Generalization conflicts

These conflicts arise when two entities are represented at different levels of generalization in two different databases.

*Example.*

Consider the entity "Graduate Students" which may be represented in two different databases as follows:

STUDENT(Id#, Name, Major, Type)

GRAD-STUDENT(Id#, Name, Major)

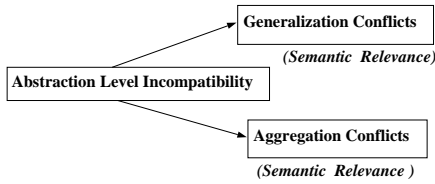
Thus, we have the same entity set being defined at a more general level in the first database.

The definition contexts of the two objects can be defined as:  
 $C_{def}(\text{STUDENT}) = \langle (\text{type}, \{\text{graduate}, \text{undergraduate}\}) \rangle$   
 $C_{def}(\text{GRAD-STUDENT}) = \langle (\text{type}, \{\text{graduate}\}) \rangle$

The context in which  $\text{semPro}(\text{STUDENT}, \text{GRAD-STUDENT})$  is defined is given by:

$$\text{glb}(C_{def}(\text{STUDENT}), C_{def}(\text{GRAD-STUDENT})) \\ = \langle (\text{type}, \{\text{graduate}\}) \rangle$$

The abstraction is then computed by "conditioning" the original student abstraction *with respect to* this new context. Thus, STUDENT and GRAD-STUDENT may be considered *semantically relevant*.



**Fig. 10.** Abstraction level incompatibilities and the likely types of semantic proximities

semPro(STUDENT, GRAD-STUDENT)  
 = <SOME, M, (D<sub>1</sub>, D<sub>2</sub>), ->  
 where M: STUDENT → GRAD-STUDENT is a partial 1-1  
 value mapping  
 and Context = SOME = <(type, {graduate})>

#### 5.4.2 Aggregation conflicts

These conflicts arise when an aggregation is used in one database to identify a set of entities in another database. Also, the properties of the aggregate concept can be an aggregate of the corresponding property of the set of entities.

*Example.*

Consider the aggregation SET-OF which is used to define a concept in the first database and the set of entities in another database as follows:

```

CONVOY(Id#, AvgWeight, Location)
SHIP(Id#, Weight, Location, Captain)
  
```

Thus, CONVOY in the first database is a SET-OF SHIPs in the second database. Also, CONVOY.AvgWeight is the average (aggregate function) of SHIP.Weight of ships that are members of the convoy.

In this case, there is a mapping in only one direction, i.e., an element of a set is mapped to the set itself. In the other direction, the mapping is not precise. When the SHIP entity is known, one can identify the CONVOY entity it belongs to, but not vice versa. Also, the aggregation can be expressed in the definition context of CONVOY using the composition of contextual coordinates as follows:

```

Cdef(CONVOY)
= <(member, SHIP) (weight, ...) (location, ...)>,
Cdef(SHIP) = <(shipweight, ...) (shiplocation, ...)>,
where weight = average(shipweight) and shiplocation = location
are relationships between the various contextual coordinates
obtained from the ontology of the domain.
context = glb(Cdef(CONVOY), Cdef(SHIP))
  
```

```

semPro(CONVOY, SHIP)
= <context, Aggregation, (D1, D2), ->
  
```

Thus, CONVOY and SHIP maybe considered *semantically relevant*.

#### 5.5 Schematic discrepancies

This class of conflicts was discussed in [DAODT85, KLK91]. It was noted that these conflicts can take place within the same data model and arise when data in one database correspond to meta-data of another database. This class of conflicts is similar to that discussed in Sect. 5.3 when the conflicts depend on the database state. We now analyze the problem and identify three aspects with help of an example given in [KLK91]. We also discuss the possible semantic similarities with each case (Fig. 11).

*Example.* Consider three stock databases. All contain the closing price for each day of each stock in the stock market. The schemata for the three databases are as follows:

- **Database DB1 :**  
relation r : {(date, stkCode, clsPrice) ... }
- **Database DB2 :**  
relation r : {(date, stk1, stk2, ... ) ... }
- **Database DB3 :**  
relation stk1 : {(date, clsPrice) ... },  
relation stk2 : {(date, clsPrice) ... },  
:  
:

DB1 consists of a single relation that has a tuple per day per stock with its closing price. DB2 also has a single relation, but with one attribute per stock, and one tuple per day, where the value of the attribute is the closing price of the stock. DB3 has, in contrast, one relation per stock that has a tuple per day with its closing price. Let us consider that the stkCode values in DB1 are the names of the attributes, and in the other databases they are the names of relations (e.g., stk1, stk2).

##### 5.5.1 Data value attribute conflict

This conflict arises when the value of an attribute in one database corresponds to an attribute in another database. Thus, this kind of conflict depends on the *database state*. Referring to the above example, the values of the attribute *stkCode* in the database *DB1* correspond to the attributes *stk1, stk2, ...* in the database *DB2*.

The mappings here are established between sets of attributes ( $\{O_i\}$ ) and values in the extension of the other attribute ( $O_2$ ). This is possible, however only *with respect* to the contexts of the databases they are in. The two objects model data at different levels and hence may be considered to be *meta-semantically relevant* and their *semantic proximity* can be defined as follows:

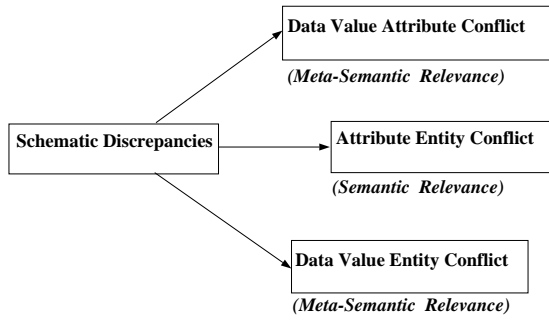
```

semPro({Oi}, O2) = <context, M, (D1, D2), (S1, S2)>,
where context = glb(Cdef(DB1), Cdef(DB2))
and M is a total 1-1 mapping between {Oi} and S2.
  
```

##### 5.5.2 Attribute entity conflict

This conflict arises when the same entity is being modeled as an attribute in one database and a relation in another





**Fig. 11.** Schematic discrepancies and the likely types of semantic proximities

database. This kind of conflict is different from the conflicts defined in the previous and next subsections, because it depends on the *database schema* and not on the *database state*. It can also be considered as a part of the entity definition incompatibility (Sect. 5.2). Referring to the example described in the beginning of this section, the attribute *stk1*, *stk2* in the database *DB2* correspond to relations of the same name in the database *DB3*.

Objects  $O_1$  and  $O_2$  can be considered *semantically relevant*, as 1-1 value mappings can be established between the domains of the attribute ( $O_1$ ) and the domain of the identifying attribute of the entity ( $O_2$ ). It should be noted that  $O_1$  is an attribute (property) and  $O_2$  is an entity (class) and their definition contexts are needed to determine the identifying attribute of the entity ( $O_2$ ).

$\text{semPro}(O_1, O_2)$   
 $= \langle \text{context, total 1-1 value mapping, } (D_1, D_2), \_ \rangle$   
 where  $\text{context} = \text{glb}(C_{\text{def}}(DB2), C_{\text{def}}(DB3))$   
 and  $D_1 = \text{Domain}(O_1)$  and  $D_2 = \text{Domain}(\text{Identifier}(O_2))$ .

### 5.5.3 Data value entity conflict

This conflict arises when the value of an attribute in one database corresponds to a relation in another database. Thus, this kind of conflict depends on the *database state*. Referring to the example described in the beginning of this section, the values of the attribute *stkCode* in the database *DB1* correspond to the relations *stk1*, *stk2* in the database *DB3*.

The mappings here are established between set of entities ( $\{O_i\}$ ) and values in the extension of an attribute ( $O_2$ ). This is possible, however only *with respect to* the contexts of the databases they are in. Thus, the two objects may be considered to be *meta-semantically relevant* and their *semantic proximity* can be defined as follows:

$\text{semPro}(\{O_i\}, O_2) = \langle \text{context, } M, (D_1, D_2), (S_1, S_2) \rangle$ ,  
 where  $\text{context} = \text{glb}(C_{\text{def}}(DB1), C_{\text{def}}(DB2))$   
 and  $M$  is a total 1-1 mapping between  $\{O_i\}$  and  $S_2$ .

## 6 Structural similarity: a component of semantic similarity

In this section, we propose a uniform formalism called **schema correspondences** for representation of structural

similarities between objects. These are associations between objects and types defined in the various databases and can be expressed using operations from a modified object algebra. The schema correspondences so defined are a part of the semantic proximity between the two objects or types and are dependent on the context in which the semantic proximity is defined. **Projection rules** which define the relationship between schema correspondences and semantic proximity are also discussed.

### 6.1 Schema correspondences: a uniform formalism for representation of abstraction

We propose a uniform formalism to represent the mappings which are generated to represent the structural similarities between objects having schematic conflicts and some semantic affinity. This formalism is a generalization of the concept of *connectors* used to augment the relational model in [CRE87].

Given two objects  $O_1$  and  $O_2$ , the *schema correspondence* between them can be represented as

$\text{schCor}(O_1, O_2) = \langle O_1, \text{attr}(O_1), O_2, \text{attr}(O_2), M \rangle$ .

- $O_1$  and  $O_2$  are objects in the model world. They are representations or intensional definitions in the model world. They may correspond to class definitions or type definitions in a database.
- The objects enumerated above may model information at any level of representation (such as the entity or the attribute level). If an object  $O_i$  models information at the entity level, then  $\text{attr}(O_i)$  denotes the representation of the attributes of  $O_i$ . If  $O_i$  models objects at the attribute level, then  $\text{attr}(O_i)$  is an empty set.
- $M$  is a mapping (possibly higher order) expressing the correspondences between objects, their attributes and the values of the objects/attributes. We use object algebra operations enumerated below.

#### 6.1.1 A brief introduction to a limited object algebra

Objects are considered as collections of instances which are homogeneous and have the same type as the abstract data type associated with the object. We list a limited set of operations to manipulate objects in a database; these are very similar to those in object-oriented database literature (e.g., [SZ90])<sup>7</sup>.

*OSelect(p,O)* This operation selects a set of instances of an object  $O$  satisfying a selection predicate,  $p$ .

$O\text{Select}(p,O) = \{o \mid o \in O \wedge p(o)\}$

*makeObject(C,S)* Given a contextual coordinate  $C$  and a set  $S$  (which may be either a set of concepts from an ontology, an object or a type domain), it defines a new object with instances having attribute  $C$  and a value from the set  $S$  as its value.

$\text{makeObject}(C,S) = \{o \mid o.C=s \wedge s \in S\}$

<sup>7</sup> When defining and using these operations, performance issues are ignored in favor of simplicity of description.

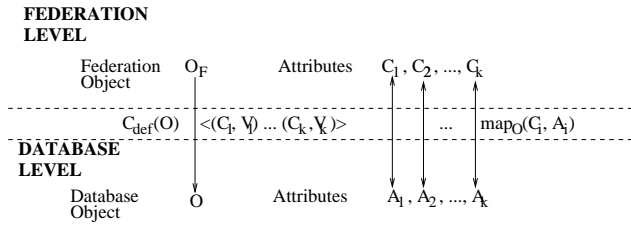


Fig. 12. Schema correspondences: association between federation and database objects

$OProduct(O_1, O_2)$  Given two objects  $O_1$  and  $O_2$ , a new object is created which has the attributes of both  $O_1$ , and  $O_2$ , and for every tuple of values in  $O_1$  has all the tuples of values in  $O_2$  associated with it.

$$OProduct(O_1, O_2) = \{o \mid (o.A_i = o_1.A_i \wedge A_i \in \text{attr}(O_1) \wedge o_1 \in O_1) \vee (o.A_j = o_2.A_j \wedge A_j \in \text{attr}(O_2) \wedge o_2 \in O_2)\}$$

$OJoin(p, O_1, O_2)$  This can be thought of as a special case of the operation  $OProduct$ , except that the instances should satisfy the predicate  $p$ .

$$OJoin(p, O_1, O_2) = \{o \mid o \in OProduct(O_1, O_2) \wedge p(o)\}$$

### 6.1.2 Schema correspondences and context

Each information system exports federation objects  $O_F$  corresponding to the objects  $O$  it manages. The objects  $O_F$  are obtained by applying the constraints in the definition context  $C_{def}(O)$  to the object  $O$ . The user at the federation level sees only the federation objects. The contextual coordinates  $C_i$  of the  $C_{def}(O)$  act as the attributes of  $O_F$ . The exported objects  $O_F$  are associated with the objects and types defined in the database. This association might be implemented in different ways by various component systems. We use schema correspondences to express these associations. This is illustrated in Fig. 12

$$\text{schCor}(O_F, O) = \langle O_F, \{C_i \mid C_i \in C_{def}(O)\}, O, \text{attr}(O), M \rangle$$

- $O_F$  is the exported federation object of an object  $O$  or type  $T$  defined in the database.
- The attributes of the object  $O_F$  are the contextual coordinates of the definition context  $C_{def}(O)$ .
- The mapping operation  $\text{map}_O(C_i, A_i)$  stores the association between contextual coordinate  $C_i$  and attribute  $A_i$  of object  $O$  whenever there exists one.
- The mapping  $M$  between  $O_F$  and  $O$  can be evaluated using the projection rules enumerated and illustrated in Sect. 6.2.

### 6.2 Schema correspondences: projection of semPro with respect to context

We discussed in Sect. 3.1 how representing structural similarities is not enough to capture semantic similarity between two objects. However, for any meaningful operation to be performed on the computer, the semPro descriptor between two objects has to be mapped to a mathematical expression which would essentially express the structural correspondence between two objects. Our approach consists of the following three aspects:

**The semantic aspect:** The semPro descriptor captures the RWS of the data in the database through context and includes intensional descriptions of:

- objects and their attributes
- the relationships between various objects
- the implicit assumptions in the design of the objects
- the constraints which the objects and attributes satisfy

The federation objects are objects obtained by applying the constraints in the intensional descriptions to the database objects.

**The data organization aspect:** This refers to the actual organization of the data in the databases, e.g., the tables and views in a relational database, or the class hierarchy in object-oriented databases.

**The mapping/abstraction aspect:** The schCor descriptor, as defined earlier, captures the association between the federation objects and the database objects. The association uses object algebraic operations to express correspondences between the federation and the database objects. The evaluation of these associations results in the retrieval of database objects which satisfy the constraints specified in the context.

The mapping aspect can be succinctly expressed as

$$\text{schCor}(O_F, O) = \Pi_{Context}(\text{semPro}(O_F, O))$$

In the rest of this section, we explain the mapping aspect with the help of examples. We first define the terminology, operations and the projection rules used to specify the semantics of the associations between the federation and database objects, followed by examples illustrating them.

#### 6.2.1 Relevant terminology and projection rules

We first enumerate the operations used to specify the associations between the exported federation objects and the database objects. We shall use  $\text{Cntxt}$ ,  $\text{Cntxt}_1, \dots$  to refer to contexts and  $C, C_1, \dots$  to refer to contextual coordinates.  $O_1, O_2, \dots$  shall be used to refer to actual database objects whereas  $O_{1F}, O_{2F}, \dots$  shall be used to denote their counterparts exported to the federation.  $O', O'', \dots$  shall be used to denote temporary objects to illustrate each step.

The operations are as follows:

$\text{map}_O(C, A)$  The mapping operation which stores the association between the attribute  $C$  of the exported federation object  $O_F$  (which is essentially a contextual coordinate of the definition context  $C_{def}(O)$  chosen from a domain-specific ontology) and the attribute  $A$  of the object  $O$ .

$\text{semConstrain}(\langle (C_i, V_i) \rangle, \text{semPro}(O', O))$  The exported federation object  $O_F$  is obtained by iteratively applying the constraints in  $C_{def}(O)$  to the database object  $O$ . The  $\text{semConstrain}$  operation models one iteration, i.e., the application of one constraint in  $C_{def}(O)$  to the database object  $O$ . Let

- $\text{semPro}(O_F, O)$  be defined *with respect to*  $C_{def}(O)$
- $C_i$  be a contextual coordinate of  $C_{def}(O)$
- $C_{def}(O) = \text{glb}(\langle (C_i, V_i) \rangle, \text{Cntxt})$  (discussed in Sect. 3.3)

- $\text{semPro}(O', O)$  be defined *with respect to*  $\text{Cntxt}$  and
- $O'$  be a temporary object obtained by applying all the constraints in  $\text{Cntxt}$  on  $O$ ;

then the federation object  $O_F$  may be iteratively defined as

$$\begin{aligned} \text{semPro}(O_F, O) \\ = \text{semConstrain}(\langle (C_i, V_i) \rangle, \text{semPro}(O', O)) \end{aligned}$$

$\text{strConstrain}(\text{map}_O(C_i, A_i), S_i, \text{schCor}(O', O))$   $\text{strConstrain}$  is the structural counterpart of  $\text{semConstrain}$ . It maps the attributes of the federation object to the attributes of the database object. It also recomputes the mappings associated with  $\text{schCor}(O', O)$ . This is done by adding a selection condition to the original mapping as follows:  
 $O_F = O\text{Select}((A_i \in S_i), O')$ ,  
 where there exists a mapping between  $O'$  and  $O$  from  $\text{schCor}(O', O)$

$\text{semCondition}(\text{Cntxt}, \text{semPro}(O_F, O))$  In some cases, a database object  $O$  may be associated with another database object *with respect to* the context  $\text{Cntxt}$ . The  $\text{semCondition}$  operation modifies the semantic proximity descriptor by *lifting* [Guh91] it into a context ( $\text{Cntxt}$ ) different from the one ( $C_{def}(O)$ ) in which it is defined in. This operation can be defined iteratively using the  $\text{semConstrain}$  operation.

$$\begin{aligned} \text{Let } \text{Cntxt} &= \text{glb}(\langle (C_i, V_i) \rangle, \text{Cntxt}_1) \\ \text{semCondition}(\text{Cntxt}, \text{semPro}(O_F, O)) \\ &= \text{semConstrain}(\langle (C_i, V_i) \rangle, \\ &\quad \text{semCondition}(\text{Cntxt}_1, \text{semPro}(O_F, O))) \end{aligned}$$

$\text{semCombine}(C_i, \text{semPro}(O', O), \text{semPro}(O'', O_i))$  In some cases, the definition context of an object  $O$  makes explicit an association between the database objects  $O$  and  $O_i$ . This association is typically *with respect to* the association context between two objects denoted as  $C_{ass}(O_i, O)$ . The  $\text{semCombine}$  operation models the *correlation* of information from objects  $O$  and  $O_i$ , which is then exported as a part of the federation object  $O_F$ . Let

- $\text{semPro}(O_F, O)$  be defined *with respect to*  $C_{def}(O)$
  - $C_{def}(O) = \text{glb}(\langle (C_i, O_i \circ C_{ass}(O_i, O)) \rangle, \text{Cntxt})$
  - $\text{semPro}(O', O)$  is defined *with respect to*  $\text{Cntxt}$
  - $O'$  be a temporary object obtained by applying the constraints in  $\text{Cntxt}$  to  $O$
  - $O''$  be a temporary object obtained by applying the constraints in  $C_{ass}(O_i, O)$  to  $O_i$ ;
- then the  $\text{semPro}(O_F, O)$  can be defined as
- $$\begin{aligned} \text{semConstrain}(\langle (C_i, O_i \circ C_{ass}(O_i, O)) \rangle, \text{semPro}(O', O)) \\ = \text{semCombine}(C_i, \text{semPro}(O', O), \text{semPro}(O'', O_i)) \\ \text{where } \text{semPro}(O'', O_i) \\ = \text{semCondition}(C_{ass}(O_i, O), \text{semPro}(O_i, O_i)) \end{aligned}$$

$\text{strCombine}(\{\text{map}_O(C_i, A_i), \text{map}_{O_i}(C_i, A'_i)\}, \text{schCor}(O', O), \text{schCor}(O'', O_i))$   $\text{strCombine}$  is the structural counterpart of  $\text{semCombine}$ . It maps the contextual coordinate  $C_i$  to the attributes of the database objects  $O$  and  $O_i$ . It correlates instances of the two objects. This results in a join condition used to combine mappings associated with  $\text{schCor}(O', O)$  and  $\text{schCor}(O'', O_i)$ .

$O_F = O\text{Join}((A_i = A'_i), O', O'')$  where there exist mappings between  $O'$  and  $O$  from  $\text{schCor}(O', O)$  and between  $O''$  and  $O_i$  from  $\text{schCor}(O'', O_i)$

## Projection rules

We describe here a set of *projection rules* which specify the semantics of the projection operation discussed earlier in this section. The rules specify an algebra based on the operations discussed above. Here we describe them with the perspective of the role they play in mapping the federation objects to the various database objects. A detailed specification of these rules is presented in the Appendix 1.

*Rule 1.* When the definition context of a database object is empty, i.e., there are no constraints which the object should satisfy, it is exported to the federation as it is without any modifications. This situation is captured by the *Empty Context Rule*.

*Rule 2.* The *Simple Sets Rule* deals with the case when the definition context has simple sets of values associated with each contextual coordinate. Each contextual coordinate is also associated with an attribute of a database object. The effect of this rule can be achieved with repeated applications of *Rule 3.1* but it is used to simplify the evaluation of the projection operation. An example of application of this rule is illustrated in Sect. 6.2.2.

*Rule 3.* The exported federation object  $O_F$  is obtained by iteratively applying the constraints in the definition context to the database object  $O$ . The *Simple Set Constraint Rule* deals with the case where the constraints in the context are applied iteratively to the database objects. The termination condition of the iteration is the case when the context is empty and is covered by the *Empty Context Rule*. This rule deals with the case where the constraint to be applied is of the form  $C \in S$ , where  $C$  is a contextual coordinate and  $S$  is a simple set of symbols from the ontology. This rule may also be used to apply an arbitrary constraint on a federation object.

*Rule 3.1.* This rule deals with the case where the contextual coordinate in the constraint is not present in the definition context in which the  $\text{semPro}$  is defined and there exists an attribute of a database object corresponding to that contextual coordinate.

*Rule 3.2.* This rule deals with the case where the contextual coordinate in the constraint is already present in the definition context and there exists an attribute of a database object corresponding to that contextual coordinate.

*Rule 3.3.* This rule deals with the case where the contextual coordinate in the constraint is not present in the definition context and there does not exist an attribute of a database object corresponding to that contextual coordinate. An example of application of this rule is illustrated in Sect. 6.2.3.

*Rule 3.4.* This rule deals with the case where the contextual coordinate in the constraint is present in the definition context and there does not exist an attribute of a database object corresponding to that contextual coordinate.

*Rule 4.* In some cases, a database object  $O_i$  may be associated with another database object  $O$  *with respect to* an association context. The *Context Conditioning Rule* deals with the case where  $\text{semPro}(O_i, O_i)$  is conditioned *with respect to* the association context. This involves applying

the constraints in the association context to the federation object  $O_{iF}$ .

*Rule 4.1.* The *Empty Context Conditioning Rule* states that when the association context used to condition the semantic proximity is empty, then the semantic proximity is evaluated *with respect to* the definition context. This means that the federation object  $O_{iF}$  is returned as it is without modification.

*Rule 4.2.* The *Constraint Conditioning Rule* deals with the case when the constraints in the association context are applied to the federation object  $O_{iF}$  iteratively. The termination condition of this iteration is when the association context is empty and is covered by the *Empty Context Conditioning Rule*.

*Rule 4.3.* The *Context Conditioning and semCombine Rule* deals with the case when the semantic proximity descriptor is a combination of two semantic proximities combined using the *semCombine* operation. The semantics of the *semCombine Rule* are given by Rule 5.

*Rule 5.* In some cases, the definition context of an object  $O$  makes explicit an association between the database objects  $O$  and  $O_i$ . This association is typically *with respect to* the association context between two objects denoted as  $C_{ass}(O_i, O)$ . The *semCombine Rule* deals with this case and results in the generation and combination of two semantic proximities. An example of application of this rule is illustrated in Sect. 6.2.4.

*Rule 5.1.* This rule maps the contextual coordinate to the attributes in the different objects and performs the correlation of the instances of the two objects. The two attributes may either satisfy the equality predicate or any other well-defined predicate.

*Rule 5.2.* The *Coordinate Composition Rule* deals with the special case where the contextual coordinate in the constraint may be a composition of two contextual coordinates. Each of the contextual coordinate parts may or may not be mapped into attributes of database objects. An example of application of this rule is presented in Sect. 6.2.5.

## 6.2.2 Using ontology for an intensional description of data

In Sect. 3.2, we chose the contextual coordinates  $C_{iS}$  and their values  $V_{iS}$  from an ontology. We illustrate with the help of an example how concepts in an ontology may be mapped to the actual data in the database. Thus, the user at the federation level can view the information in the database with the help of concepts from a domain-specific ontology without being aware of the underlying format of the data.

*Example.* Consider an object EMPLOYEE defined in a database as follows:

EMPLOYEE(SS#,Name,Dept,SalaryType,Affiliation)

The definition context of the object EMPLOYEE may be defined as:

$$C_{def}(EMPLOYEE) = \langle (\text{employer}, [\mathbf{Deptyes} \cup \{\text{restypes}\}]) \\ (\text{affiliation}, \{\text{teaching}, \text{research}, \text{non-teaching}\}) \\ (\text{reimbursement}, \{\text{salary}, \text{honorarium}\}) \rangle$$

- **Deptyes** is a type defined in the database.
- The symbols for the contextual coordinates *employer*, *affiliation* and *reimbursement* are taken from the ontology. The association with the attributes of EMPLOYEE is stored by the  $\text{map}_{EMPLOYEE}(C, A)$  operation.
- The symbols *restypes*, *teaching*, *research*, *non-teaching*, *salary* and *honorarium* may either be taken from the ontology or submitted for inclusion into the ontology by the database administrator.

As discussed in Sect. 6, we associate with definition context an object  $EMPLOYEE_F$  which is exported to the federation of databases.

$\text{semPro}(EMPLOYEE_F, EMPLOYEE)$   
 $= \langle C_{def}(EMPLOYEE), M, (\text{dom}(EMPLOYEE_F), \text{dom}(EMPLOYEE)), - \rangle$ ,  
 where  $M$  is a mapping between the domains of the two objects. The mapping relates information in the ontology to data in the database. The projection is illustrated in Fig. 13.

**Simple Sets Rule**  $\Rightarrow$

$$\begin{aligned} & \Pi_{C_{def}(EMPLOYEE)}(\text{semPro}(EMPLOYEE_F, EMPLOYEE)) \\ & = \text{schCor}(EMPLOYEE_F, EMPLOYEE) \\ & = \langle EMPLOYEE_F, \{\text{employer}, \text{affiliation}, \text{reimbursement}\}, EMPLOYEE, \\ & \quad \{\text{map}_{EMPLOYEE}(\text{employer}, \text{Dept}), \\ & \quad \text{map}_{EMPLOYEE}(\text{affiliation}, \text{Affiliation}), \\ & \quad \text{map}_{EMPLOYEE}(\text{reimbursement}, \text{SalaryType})\}, M \rangle \\ & M \equiv EMPLOYEE_F = O\text{Select}(p, EMPLOYEE) \\ & p \equiv (\text{Dept} \in \{\mathbf{Deptyes} \cup \{\text{restypes}\}\}) \\ & \quad \wedge (\text{Affiliation} \in \{\text{teaching}, \text{research}, \text{non-teaching}\}) \\ & \quad \wedge (\text{SalaryType} \in \{\text{salary}, \text{honorarium}\}) \end{aligned}$$

## 6.2.3 Domain augmentation: representing extra information

In this section, we demonstrate an interesting case where *extra information* can be stored with the intensional descriptions of objects. This extra information is represented as a constraint at the federation level. Consider the constraint: all publications have research areas that are associated with departments. This may be used to make inferences about database content, without actually accessing the database. Consider a query that asks for all publications in a research area not associated with a department. The answer to the query is an empty set which can be determined *without* actually accessing the database.

The constraint involving research areas can be represented in  $C_{def}(\text{PUBLICATION})$  and expressed using the contextual coordinate *researchArea*. However, the information about the research areas of a publication is not modeled by the existing database object PUBLICATION(Id, Title, Journal).

The definition context of the object PUBLICATION is defined as:

$C_{def}(\text{PUBLICATION}) = \langle (\text{researchArea}, \mathbf{Deptyes}) \rangle$  where **Deptyes** is a type defined in the database.

The query discussed above can be processed without accessing the database if the constraint involving research areas is part of the exported federation object. Because the contextual coordinate *researchArea* is not modeled in the database, the projection algorithm creates a new object corresponding to the research areas by using the *makeObject* operation. This new object is then associated with the database object PUBLICATION by using the *OProduct* operation. The above

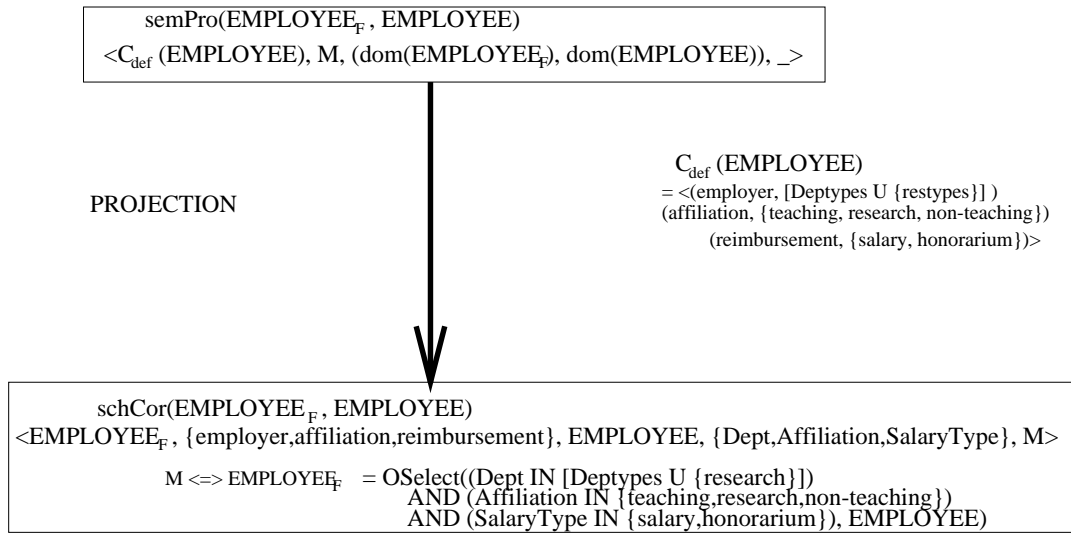


Fig. 13. Mapping  $\text{EMPLOYEE}_F$  to object  $\text{EMPLOYEE}$  in the database

results in the augmentation of the domain of the database object  $\text{PUBLICATION}$  and is expressed in Appendix 1 (Rule 3.3).

$$\text{dom}(\text{PUBLICATION}_F) \subseteq \text{dom}(\text{Id}) \times \text{dom}(\text{Title}) \\ \times \text{dom}(\text{Journal}) \times \mathbf{Deptypes}.$$

The projection operation is diagrammatically illustrated in Fig. 14.

- [A ]  $\text{semPro}(\text{PUBLICATION}_F, \text{PUBLICATION})$  is evaluated *with respect to*  $C_{\text{def}}(\text{PUBLICATION})$ . The definition context expresses extra information about the object  $\text{PUBLICATION}$  not modeled in the database. This step illustrates the augmentation of  $\text{dom}(\text{PUBLICATION})$ . Let:
- $C_{\text{def}}(\text{PUBLICATION}) = \text{glb}(\langle \text{researchArea}, \text{Deptypes} \rangle, \langle \_ \rangle)$
  - $\text{semPro}(\text{PUBLICATION}', \text{PUBLICATION})$  be defined *with respect to*  $\langle \_ \rangle$
  - $\text{PUBLICATION}'$  be a temporary object
- The domain augmentation takes place as follows:
- Simple Set Constraint Rule (New Constraint, Non-existing attribute)  $\Rightarrow$  (step [B])
- $$\text{semPro}(\text{PUBLICATION}_F, \text{PUBLICATION})$$
- $$= \text{semConstrain}(\langle \text{researchArea}, \mathbf{Deptypes} \rangle,$$
- $$\text{semPro}(\text{PUBLICATION}', \text{PUBLICATION}))$$
- Let  $M'$  be the mapping between  $\text{PUBLICATION}'$  and  $\text{PUBLICATION}$  returned by step [C].
  - The constraint about research areas is incorporated in the exported federation object  $\text{PUBLICATION}_F$  by using the mapping  $M$ . The evaluation of the mapping is illustrated in steps [D,E].
  - The resulting augmentation of the domain of the object  $\text{PUBLICATION}$  is reflected in the definition of the modified  $\text{semPro}$  descriptor:
- $$\text{semPro}(\text{PUBLICATION}_F, \text{PUBLICATION})$$
- $$= \langle C_{\text{def}}(\text{PUBLICATION}), M,$$
- $$(\text{dom}(\text{PUBLICATION}_F), \mathbf{dom}(\text{PUBLICATION}) \times \mathbf{Deptypes}), \_ \rangle$$
- [C ] Empty Context Rule  $\Rightarrow$
- $$M' \equiv \text{PUBLICATION}' = \text{PUBLICATION}$$
- [D,E ] Simple Set Constraint Rule (Rule 3.3)  $\Rightarrow$
- $$\text{schCor}(\text{PUBLICATION}_F, \text{PUBLICATION})$$
- $$= \text{strConstrain}(\{\text{researchArea}\}, \mathbf{Deptypes},$$
- $$\text{schCor}(\text{PUBLICATION}', \text{PUBLICATION}))$$
- $$M \equiv \text{PUBLICATION}_F$$
- $$= \text{OProduct}(\text{makeObject}(\text{researchArea}, \mathbf{Deptypes}),$$
- $$\text{PUBLICATION}')$$
- $$= \text{OProduct}(\text{makeObject}(\text{researchArea}, \mathbf{Deptypes}), \text{PUBLICATION})$$

#### 6.2.4 Representing relationships between objects

In this section, we illustrate with the help of an example how context can be used to capture relationships between objects which may not be represented in the database. We illustrate a case where the definition context of the object  $\text{HAS-PUBLICATION}$  captures its relationships with another database object  $\text{EMPLOYEE}$  in an intensional manner. These relationships are *not stored* in the database and the evaluation of the  $\text{semPro}$  descriptor results in *extra information* being associated with the federation object  $\text{HAS-PUBLICATION}_F$ . A naive user will ordinarily not be aware of this relationship.

*Example.* Consider objects  $\text{EMPLOYEE}$  and  $\text{PUBLICATION}$  defined earlier and an object in the same database which represents a relationship between employees and the publications they write,  $\text{HAS-PUBLICATION}(\text{SS}\#, \text{Id})$

$$C_{\text{def}}(\text{HAS-PUBLICATION})$$

$$= \langle (\text{author}, \text{EMPLOYEE} \circ \langle (\text{affiliation}, \{\text{research}\}) \rangle) \rangle$$

This evaluation of the  $\text{semPro}$  descriptor has been diagrammatically illustrated in Fig. 15.

- [A ]  $\text{semPro}(\text{HAS-PUBLICATION}_F, \text{HAS-PUBLICATION})$  is evaluated *with respect to*  $C_{\text{def}}(\text{HAS-PUBLICATION})$ . The definition context makes explicit the relationship between  $\text{HAS-PUBLICATION}$  and  $\text{EMPLOYEE}$ . This step illustrates how the correlation of the instances of  $\text{EMPLOYEE}$  and  $\text{HAS-PUBLICATION}$  is done to satisfy the constraints in the definition context. Let
- $C_{\text{def}}(\text{HAS-PUBLICATION})$
  - $\text{glb}(\langle \text{author}, \text{EMPLOYEE} \circ \langle (\text{affiliation}, \{\text{research}\}) \rangle \rangle,$
  - $\langle C_{\text{ass}}(\text{EMPLOYEE}, \text{HAS-PUBLICATION}) \rangle \rangle, \langle \_ \rangle)$
  - $\text{semPro}(\text{HAS-PUBLICATION}', \text{HAS-PUBLICATION})$  be defined *with respect to*  $\langle \_ \rangle$
  - $C_{\text{ass}}(\text{EMPLOYEE}, \text{HAS-PUBLICATION})$
  - $= \langle (\text{affiliation}, \{\text{research}\}) \rangle$
  - $\text{HAS-PUBLICATION}'$  be a temporary object
  - $\text{EMPLOYEE}'$  be a temporary object obtained by applying the constraints in  $C_{\text{ass}}(\text{EMPLOYEE}, \text{HAS-PUBLICATION})$  to  $\text{EMPLOYEE}_F$
- $\text{semCombine}$  Rule  $\Rightarrow$
- $$\text{semPro}(\text{HAS-PUBLICATION}_F, \text{HAS-PUBLICATION})$$

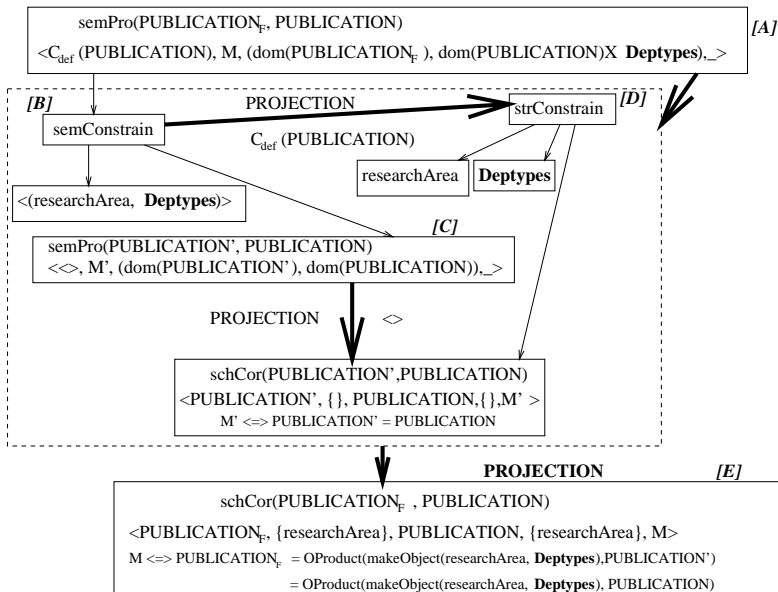


Fig. 14. Domain augmentation: mapping  $PUBLICATION_F$  to object  $PUBLICATION$  in the database

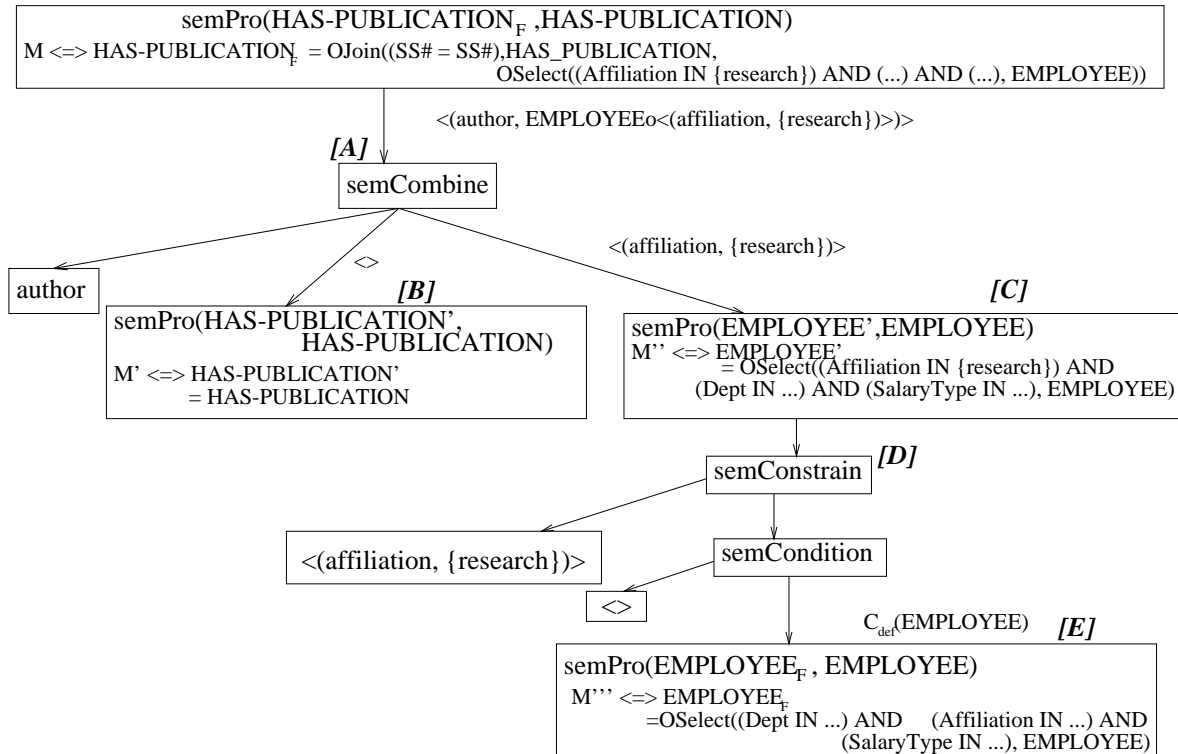


Fig. 15. Correlation of information between  $HAS-PUBLICATION$  and  $EMPLOYEE$

$= semCombine(author,$   
 $semPro(HAS-PUBLICATION', HAS-PUBLICATION),$   
 $semPro(EMPLOYEE', EMPLOYEE))$   
 - Let  $M'$  be the mapping between  $HAS-PUBLICATION'$  and  $HAS-PUBLICATION$  returned by step [B].  
 -  $semPro(EMPLOYEE', EMPLOYEE)$   
 $= semCondition(C_{ass}(EMPLOYEE, HAS-PUBLICATION),$   
 $semPro(EMPLOYEE_F, EMPLOYEE))$   
 Let  $M''$  be the mapping between  $EMPLOYEE'$  and  $EMPLOYEE$  returned by step [C].  
 -  $map_{EMPLOYEE}(author, SS#)$

-  $map_{HAS-PUBLICATION}(author, SS#)$   
 Rule 5.1  $\Rightarrow$   
 $M \equiv HAS-PUBLICATION_F =$   
 $OJoin((SS#=SS#), HAS-PUBLICATION', EMPLOYEE')$   
 $= OJoin((SS#=SS#), HAS-PUBLICATION, EMPLOYEE')$   
 ....  $M'$  From step [B]  
 $= OJoin((SS#=SS#), HAS-PUBLICATION,$   
 $OSelect((Affiliation \in \{research\}) \wedge (...) \wedge (...), EMPLOYEE))$   
 ....  $M''$  From step [C]

[B] Empty Context Rule  $\Rightarrow$

$M' \equiv HAS-PUBLICATION' = HAS-PUBLICATION$

[C ] In this step, we show how the constraints in the association context are applied to the federation object  $EMPLOYEE_F$ . This is done *before* correlation of the instances of  $EMPLOYEE$  and  $HAS-PUBLICATION$ , as only employees who are researchers have publications.

$$C_{ass}(EMPLOYEE,HAS-PUBLICATION)$$

$$= \text{glb}(\langle \text{affiliation}, \{\text{research}\} \rangle, \langle \rangle)$$

Constraint Conditioning Rule  $\Rightarrow$

$$\text{semCondition}(C_{ass}(EMPLOYEE,HAS-PUBLICATION),$$

$$\text{semPro}(EMPLOYEE_F,EMPLOYEE))$$

$$= \text{semConstrain}(\langle \text{affiliation}, \{\text{research}\} \rangle,$$

$$\text{semCondition}(\langle \rangle, \text{semPro}(EMPLOYEE_F,EMPLOYEE)))$$

.... Illustrated in step [D]

$$= \text{semConstrain}(\langle \text{affiliation}, \{\text{research}\} \rangle,$$

$$\text{semPro}(EMPLOYEE_F,EMPLOYEE))$$

.... Empty Context Conditioning Rule

Let  $M''$  be the mapping returned by step [E] between  $EMPLOYEE_F$  and  $EMPLOYEE$ .

Rule 3.2  $\Rightarrow$

$$M'' \equiv EMPLOYEE'$$

$$= \text{OSelect}((\text{Affiliation} \in \{\text{research}\}), EMPLOYEE_F)$$

$$= \text{OSelect}((\text{Affiliation} \in \{\text{research}\}),$$

$$\text{OSelect}((\text{Affiliation} \in \{\text{research}, \text{teaching}, \text{non-teaching}\}) \wedge (\dots) \wedge (\dots),$$

$$EMPLOYEE))$$

....  $M''$  From step [E]

$$= \text{OSelect}((\text{Affiliation} \in \{\text{research}\}) \wedge (\dots) \wedge (\dots), EMPLOYEE)$$

[E ] This step illustrates the association between the federation object  $EMPLOYEE_F$  and the database object  $EMPLOYEE$  and has been discussed in detail in Sect. 6.2.2. The association is given by:

$$M'' \equiv EMPLOYEE_F$$

$$= \text{OSelect}((\text{Affiliation} \in \{\text{research}, \text{teaching}, \text{non-teaching}\}) \wedge (\dots) \wedge (\dots),$$

$$EMPLOYEE)$$

## 6.2.5 Composition of contextual coordinates: representing extra information

In this section, we illustrate an example in which the information that the contextual coordinate *researchInfo* is a composition of two contextual coordinates (*researchArea* and *journalTitle*) is obtained from the ontology of the domain. This is then used to correlate information between the objects  $PUBLICATION$  and  $JOURNAL$ . However, the contextual coordinate *researchArea* has not been modeled for the object  $PUBLICATION$ . Thus, this results in *extra information* about the relevant journals and research areas being associated with the object  $PUBLICATION$ , *even though no information about research areas is modeled for PUBLICATION*.

*Example.* Consider a database containing the following objects:

$$PUBLICATION(\text{Id}, \text{Title}, \text{Journal}) \quad C_{def}(PUBLICATION)$$

$$= \langle (\text{researchInfo}, JOURNAL \circ \langle (\text{researchArea}, \text{Deptypes})$$

$$(\text{journalTitle}, \text{JournalTypes}) \rangle \rangle \rangle$$

$$JOURNAL(\text{Title}, \text{Area}) \quad \text{where } C_{def}(JOURNAL) = \langle \rangle$$

The *correlation of information* is illustrated diagrammatically in Fig. 16.

[A ]  $\text{semPro}(PUBLICATION_F, PUBLICATION)$  is evaluated *with respect to*  $C_{def}(PUBLICATION)$

The definition context makes explicit the relationship between  $PUBLICATION$  and  $JOURNAL$ . This step illustrates the generation of the two  $\text{semPro}$  descriptors, one for applying the remaining constraints in  $C_{def}(PUBLICATION)$  to  $PUBLICATION$  and the other for applying the constraints in  $C_{ass}(JOURNAL, PUBLICATION)$  to  $JOURNAL_F$ . Let

$$- C_{def}(PUBLICATION)$$

$$= \text{glb}(\langle (\text{researchInfo}, JOURNAL \circ \langle (\text{researchArea}, \text{Deptypes})$$

$$(\text{journalTitle}, \text{JournalTypes}) \rangle \rangle, \langle \rangle)$$

$$- \text{semPro}(PUBLICATION', PUBLICATION) \text{ be defined with respect to } \langle \rangle$$

$$- C_{ass}(JOURNAL, PUBLICATION)$$

$$= \langle (\text{researchArea}, \text{Deptypes}) (\text{journalTitle}, \text{JournalTypes}) \rangle$$

$$- PUBLICATION' \text{ be a temporary object}$$

$$- JOURNAL' \text{ be a temporary object obtained by applying the constraints in } C_{ass}(JOURNAL, PUBLICATION) \text{ to } JOURNAL$$

$\text{semCombine Rule } \Rightarrow$

$$\text{semPro}(PUBLICATION_F, PUBLICATION)$$

$$= \text{semCombine}(\text{researchInfo},$$

$$\text{semPro}(PUBLICATION', PUBLICATION),$$

$$\text{semPro}(JOURNAL', JOURNAL))$$

– Let  $M'$  be the mapping between  $PUBLICATION'$  and  $PUBLICATION$  returned by step [B].

–  $\text{semPro}(JOURNAL', JOURNAL)$

$$= \text{semCondition}(C_{ass}(JOURNAL, PUBLICATION),$$

$$\text{semPro}(JOURNAL', JOURNAL))$$

Let  $M''$  be the mapping between  $JOURNAL'$  and  $JOURNAL$  returned by step [C].

[B ] Empty Context Rule  $\Rightarrow$

The mapping  $M'$  associated with  $\text{schCor}(PUBLICATION', PUBLICATION)$  is  $M' \equiv PUBLICATION' = PUBLICATION$

[C ]  $C_{ass}(JOURNAL, PUBLICATION)$

$$= \text{glb}(\langle (\text{researchArea}, \text{Deptypes}) \rangle,$$

$$\text{glb}(\langle (\text{journalTitle}, \text{JournalTypes}) \rangle, \langle \rangle))$$

2 applications of Constraint Conditioning Rule and 1 application of Empty Context Conditioning Rule  $\Rightarrow$

$$\text{semCondition}(C_{ass}(JOURNAL, PUBLICATION),$$

$$\text{semPro}(JOURNAL_F, JOURNAL))$$

$$= \text{semConstrain}(\langle (\text{researchArea}, \text{Deptypes}) \rangle,$$

$$\text{semConstrain}(\langle (\text{journalTitle}, \text{JournalTypes}) \rangle,$$

$$\text{semPro}(JOURNAL_F, JOURNAL))$$

2 applications of Rule 3.2 and  $C_{def}(JOURNAL) = \langle \rangle \Rightarrow$

The mapping  $M''$  associated with  $\text{schCor}(JOURNAL', JOURNAL)$  is  $M'' \equiv JOURNAL'$

$$= \text{OSelect}((\text{Area} \in \text{Deptypes}) \wedge (\text{Title} \in \text{JournalTypes}), JOURNAL)$$

[D ]  $\text{semPro}(PUBLICATION_F, PUBLICATION)$  is evaluated by applying the *Coordinate Composition Rule*. The final result is illustrated in step [E]. This step illustrates how information about the research areas of the publications is propagated to  $PUBLICATION$ , even though there is no information about research areas stored in the object  $PUBLICATION$ . This is achieved by the composition of contextual coordinates obtained from the domain ontology.

$$- \text{researchInfo} = \text{compose}(\text{researchArea}, \text{journalTitle})$$

Coordinate Composition Rule  $\Rightarrow$

$$\text{map}_{PUBLICATION}(\text{researchInfo}, X)$$

$$= \text{compose}(\text{map}_{PUBLICATION}(\text{researchArea}, \text{NA}),$$

$$\text{map}_{PUBLICATION}(\text{journalTitle}, \text{Journal}))$$

$$\text{map}_{JOURNAL}(\text{researchInfo}, Y)$$

$$= \text{compose}(\text{map}_{JOURNAL}(\text{researchArea}, \text{Area}),$$

$$\text{map}_{JOURNAL}(\text{journalTitle}, \text{Title}))$$

– The mapping  $M$  associated with  $\text{schCor}(PUBLICATION_F, PUBLICATION)$  is given by:

$$\text{strCombine}(\{\text{map}_{PUBLICATION}(\text{researchInfo}, X),$$

$$\text{map}_{JOURNAL}(\text{researchInfo}, Y)\},$$

$$\text{schCor}(PUBLICATION', PUBLICATION),$$

$$\text{schCor}(JOURNAL', JOURNAL))$$

$$M \equiv PUBLICATION_F =$$

$$\text{OJoin}((X=Y), PUBLICATION', JOURNAL')$$

$$= \text{OJoin}((\text{researchArea}=\text{Area}) \wedge (\text{Title}=\text{Journal}),$$

$$PUBLICATION', JOURNAL')$$

$$= \text{OJoin}((\text{researchArea}=\text{Area}) \wedge (\text{Title}=\text{Journal}),$$

$$PUBLICATION, JOURNAL')$$

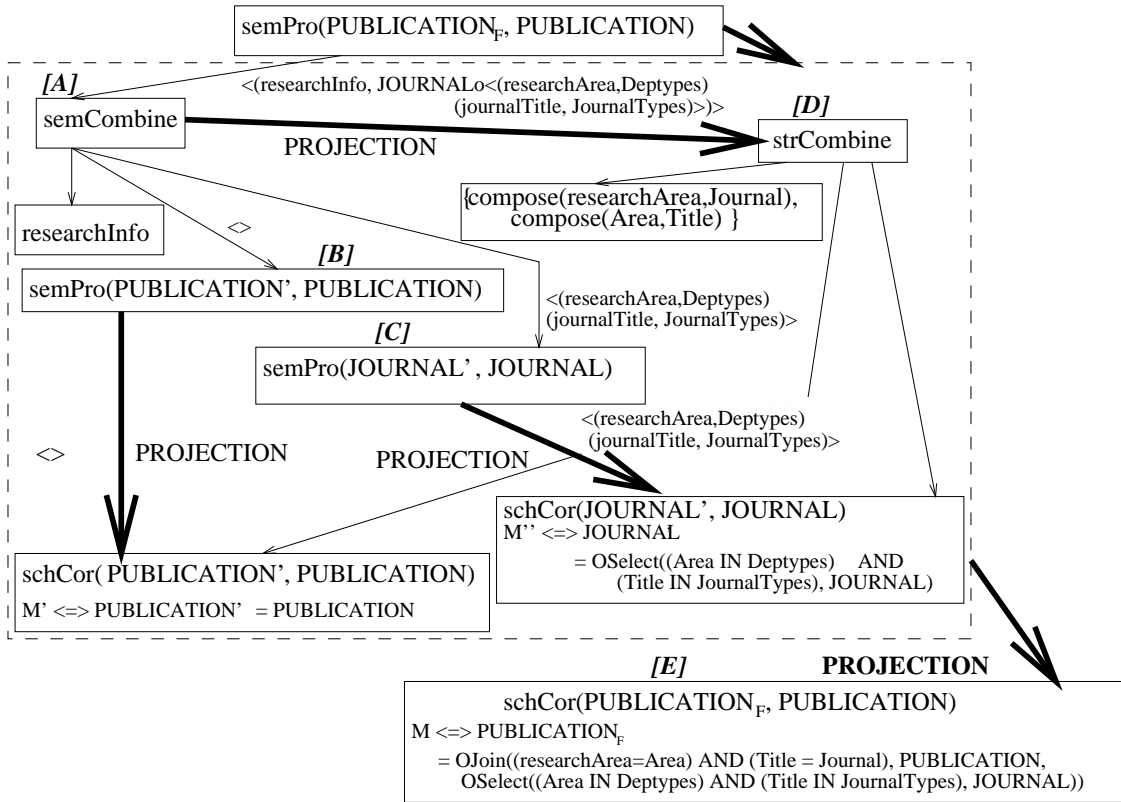


Fig. 16. Correlation between PUBLICATION and JOURNAL due to composition of contextual coordinates

.... mapping  $M'$  from step [B]  
 $= \text{OJoin}(\text{researchArea}=\text{Area}) \wedge (\text{Title}=\text{Journal}), \text{PUBLICATION},$   
 $\text{OSelect}((\text{Area} \in \text{Deptypes}) \wedge (\text{Title} \in \text{JournalTypes}), \text{JOURNAL}))$   
 .... mapping  $M''$  from step [C]

The constraint  $\text{researchArea} \in \text{Deptypes}$  propagates to PUBLICATION. This is because when the correlation takes place between JOURNAL and PUBLICATION (refer to step [E]):

- only journals belong to the research areas corresponding to the departments are selected ( $\text{OSelect}((\text{Area IN Deptypes}) \text{ AND } \dots, \text{JOURNAL}))$ )
- the join condition ( $\text{Title} = \text{Journal}$ ) ensured that only those articles which are from the research areas corresponding to the departments are exported to the federation ( $\text{OJoin}(\text{researchArea}=\text{Area}) \text{ AND } (\text{Title} = \text{Journal}), \dots$ )
- this is achieved despite the attribute Area not being modeled for PUBLICATION. Thus, there is a *selective and implicit domain augmentation* of Deptypes to PUBLICATION through the join condition.

### 6.2.6 Representation of incomplete information

The intensional description of the definition contexts can be easily used to represent incomplete information. Traditional database approaches have used NULL values to represent incomplete information. The semantics of NULL values is not always clear (e.g., a NULL value can mean unknown or not applicable) and this can be a problem while retrieving

incomplete information from the database. We can use intensional descriptions in an attempt to describe incomplete information and to avoid the problems associated with NULL values.

*Example.* Consider the following definition context of the object PUBLICATION.

$C_{def}(\text{PUBLICATION})$   
 $= \langle (\text{title}, \{x \mid \text{substring}(x) = \text{"abortion"}\}) \rangle$

This represents a constraint on the instances of the object PUBLICATION such that all the titles should have the word "abortion" in them. This does not specify the title of each instance of PUBLICATION completely. This information can be represented with the object PUBLICATION<sub>F</sub> at the federation level and can help in querying the database in face of incomplete information.

### 6.3 Applications of context

In Sect. 6.2, we defined and illustrated with examples the relationship between schema correspondences and semantic proximity. We have defined *projection rules* which define schema correspondences as the projection of the semPro descriptor *with respect to* the context. Earlier work on mapping intensional descriptions of concepts to SQL queries on relational databases has been reported in [BB93]. In our approach, however, the mappings expressed using object algebra operations are also associated with the intensional contextual descriptions. Whenever the context changes, we also



keep track of the associated changes in the schema correspondences. Rules modeling the changes in the schema correspondences (and hence the mappings) due to changes in context are presented in [KS95b].

We look at examples in which the semPro descriptors are *lifted* [Guh91] to different contexts. Lifting a semPro to a different context means re-evaluating the semPro in a context which is different from the one it was defined in first. We show in [KS95b] how query processing can be implemented by the comparison of the definition contexts of the objects in the database with the query context. We have illustrated:

- how the modification of schema correspondences due to changes in context lead to *information-focusing*
- how changes in the definition context of one object leads to the modification of schema correspondence of a related object
- how constraints from the query contexts can be applied to an object stored in a database. This results in modification of the schema correspondences and results in information focusing
- how the query context can form the basis of correlation of information across different databases

## 7 Related work

A simple observation made by various researchers in the field of multidatabases, which is also the central premise of this paper is that it is essential to associate abstractions/mappings between objects with the context of comparison to capture semantic similarity. Some significant attempts are the **semantic proximity** proposal by Sheth and Kashyap [SK92], the **context-building** approach by Ouksel and Naiman [ON93], the **context interchange** approach by Sciore et al. [SSR92] and the **common concepts** approach by Yu et al. [YSDK91]. We have related the above attempts to semantic proximity.

There have been attempts to use an attribute-value-based representation for capturing similarities in various areas of research. Larson et al [LNE89] use a set of fixed descriptors to capture similarities between attributes. Sciore et al [SSR92] use meta-attributes to represent context. In linguistics [CMG90], context has been represented using a set of context coordinates subject to certain conditions. Similar attempts have also been made for documents in text retrieval (using thematic roles) [VD92] and for clustering similar objects (using code words) in [ML92]. We have abstracted out the commonalities in these approaches in our representation of context. However, we differ from Sciore et al. [SSR92] and Ouksel et al. [ON93] in the following aspects:

- Sciore et al. [SSR92] represent the context at the extensional level, i.e., at the level of data values and object instances. We represent context at an intensional level, i.e., at the level of the database schema. This gives us an opportunity to represent constraints about objects which cannot be captured at the extensional level. We also view the context of an object as a *collection of constraints on an object* which may not be represented in the database schema

- Ouksel et al. [ON93] represent context as a collection of ISCA's (interschema correspondence assertions), which are essentially structural correspondences between schema elements in different databases. In our approach, schema correspondences are associated with the context and are not considered part of the context. They are used to relate semantic information with the actual data in the database
- the meta-attributes and their values are taken from the ontology of the application domain being modeled by the database. Issues of combining ontologies and scalability are discussed in [MKSI96, MKIS96, KS96]
- we have also defined operations to compare the specificity of contexts, and to manipulate and reason about them. Based on the partial order induced by the specificity relationship, we organize the contexts as a meet semi-lattice. Inferences on a new context *with respect to* the knowledge present in the context set can now be supported by determining its position in the semi-lattice

We have expressed our context descriptions using DL expressions. Well known DL systems are KL-ONE [BS85], LOOM [Mac87], BACK [vLNPS87] and CLASSIC [BBMR89]. We are investigating the use of CLASSIC as the DL system for representing context. The advantage of using CLASSIC is that it is sufficiently expressive and has a polynomial time classification algorithm.

Classification or taxonomies of *schematic differences* appear in [DH84, BOT86, CRE87, KLK91, KS91]. In this paper, we present what we believe is a comprehensive taxonomy of schematic conflicts which subsumes most of the taxonomies found in literature (Table 3 in Appendix 2). We refined the broad definition of domain incompatibility and entity definition incompatibility given in [CRE87]. Our classification consists of conflicts arising out of inconsistencies in the database state [BOT86], conflicts due to representation at differing levels of abstraction [DH84] and conflicts when data in one database corresponds to meta-data in another [DAODT85, KLK91].

## 8 Conclusions and future work

An essential prerequisite to achieving interoperability in a multidatabase environment is to be able to identify semantically similar data in different database systems. Another key issue attracting wide attention with attempts to build a national information infrastructure, is the issue of querying a large number of autonomous databases without prior knowledge of their information content. It is therefore important to capture the semantic content of these databases in as explicit a manner as possible.

We discussed the inadequacy of structural similarity and how semantics cannot be captured by purely mathematical formalisms. This led us to make a case for the explicit identification and representation of context in a multidatabase environment. We define the concept of *semantic proximity*, using which we represent the degrees of semantic similarities between the objects [SK92]. The *context* of comparison of these objects is the fulcrum of the semantic proximity. We propose an explicit though partial representation of context

in a multidatabase environment. We have also defined the concept of *schema correspondences*, using which we represent the structural similarities between objects.

We demonstrate the reconciliation of the dual schematic *vs* semantic perspectives. This is done by associating the mapping/abstraction between objects in different databases with the context of the semantic proximity defined between them. This association enables us to determine qualitative measures of semantic similarity such as *equivalence, relationship, relevance, resemblance and incompatibility* and develop a semantic taxonomy. We also enumerate the various schematic heterogeneities and the possible semantic similarities between them.

We have also defined the concept of *schema correspondences*, using which we represent the structural similarities between objects. Though it is known that representing structural similarities is inadequate to capture semantic similarity between two objects, for any meaningful operation to be performed on the computer, the semPro descriptor between two objects has to be mapped to a mathematical expression which would essentially express the structural correspondence between them. We have defined the schema correspondences as a projection *with respect to* context of the semantic proximity between the objects.

Besides helping to reconcile the semantic and the structural perspectives, it also enables us to represent extra knowledge about the database objects. This includes domain-specific constraints obtained from an ontology and implicit relationships between objects in the databases. We also demonstrate how extra information not modeled for a database object may be associated with it. This enables inferences to be drawn at the federation level without accessing the databases. Some of these inferences involve extra knowledge and would not have been possible, even if the objects in the databases were accessed.

These inferences are modeled as changes in the context and the associated schema correspondences. It enables *information focusing* as some inferences affect the schema correspondences to retrieve only the data relevant to the query. It enables *information correlation*, as one can specify constraints relating different objects in the context. The computation of the resulting schema correspondences enables the correlation of the appropriate instances of the objects. These have been discussed only briefly in the paper due to space constraints. The reader may refer to [KS95b] for details.

The context is the key component in capturing the semantic content of the information present in the various databases. In any attempt to represent the context of objects in a database, issues of language and vocabulary become important. We are looking into the possibility of the knowledge interchange format [GF92] and DL-based languages [BS85, BBMR89, Mac87, PS84, vLNPS87, KBR86] for context representation. In designing the definition context of an object, it is necessary to choose the contextual coordinates and their values in a controlled manner. We are experimenting on using domain-specific ontologies to construct these contexts in a methodical manner. In cases where a domain ontology is not readily available, research is required to enable semi-automatic generation of ontologies. We are looking at clustering and information retrieval techniques for semi-automatic generation of ontologies. We are

also looking into re-using well-established metadata standards and classification taxonomies as domain-specific ontologies as intensional descriptions of information content in the databases.

A complementary problem is that of presenting the ontologies to the user in a methodical manner to enable him/her to construct the query contexts for retrieving information from a federation of databases. Tools to present these ontologies to users and information system designers must be developed to facilitate context design and representation.

There should be an agreement on the meaning of the terms used in the ontologies for construction of the definition contexts on one hand and those used in the ontologies for the construction of the query contexts on the other. Thus, either a common ontology is required, or the correspondence between the terms in the various ontologies needs to be established. We are experimenting with utilization of terminological relationships between terms across ontologies. The OBSERVER system [MKSI96] using *synonym* relationships is a step in this direction. A proposal to extend the system by using *hyponym* and *hypernym* relationships has been presented in [MKIS96]. We plan to extend the system to utilize *knowledge transmutation operators* [Mic93] to express correspondences between terms in the various ontologies in the future.

## References

- [ACHK93] Arens Y., Chee C., Hsu C., Knoblock C (1993) Retrieving and integrating data from multiple information sources. *Int J Intell Coop Inf Syst*, 2
- [BB93] Borgida A., Brachman R. (1993) Loading data into description reasoners. In: *Proceedings of 1993 ACM SIGMOD*.
- [BBMR89] Borgida A., Brachman R., McGuinness D., Resnick L. (1989) CLASSIC: A structural data model for objects. In: *Proceedings of ACM SIGMOD-89*.
- [BOT86] Breitbart Y., Olson P., Thompson G. (1986) Database integration in a distributed heterogeneous database system. In: *Proceedings of the 2nd IEEE Conference on Data Engineering*.
- [BS85] Brachman R., Scmolze J. (1985) An overview of the KL-ONE knowledge representation system. *Cognitive Sci* 9:171–216
- [CHS91] Collet C., Huhns M., Shen W. (1991) Resource integration using a large knowledge base in carnot. *IEEE Comput*
- [CMG90] Chierchia G., McConnell-Ginet S. (1990) *Meaning and grammar: an introduction to semantics*. MIT Press, Cambridge, Mass
- [CRE87] Czejdo B., Rusinkiewicz M., Embley D. (1987) An approach to schema integration and query formulation in federated database systems. In: *Proceedings of the 3rd IEEE Conference on Data Engineering*.
- [DAODT85] Deen S., Amin R., Ofori-Dwumfuo G., Taylor M. (1985) The architecture of a generalised distributed database system PRECI\*. *IEEE Comput* 18
- [DH84] Dayal U., Hwang H. (1984) View definition and generalization for database integration of a multidatabase system. *IEEE Trans Software Eng* 10:628–644
- [ELN86] Elmasri R., Larson J., Navathe S. (1986) Schema integration algorithms for federated databases and logical database design. Technical report, Honeywell Corporate Systems Development Division, Golden Valley, Minn
- [EN89] Elmasri R., Navathe S. (1989) *Fundamentals of database systems*. Benjamin/Cummins, Menlo Park, Calif

- [FKN91] Fankhauser P., Kracker M., Neuhold E. (1991) Semantic vs. structural resemblance of classes. *SIGMOD Record, special issue on Semantic Issues in Multidatabases 20*
- [GF92] Genesereth M., Fikes R. (1992) Knowledge interchange format, version 3.0 reference manual. Technical Report Logic-92-1, Computer Science Department, Stanford University
- [Gru93] Gruber T. (1993) A translation approach to portable ontology specifications. *Knowl Acquis Int J Knowl Acquis Knowledge-Based Syst 5*
- [Guh90] Guha R.V. (1990) Micro-theories and contexts in Cyc. I. Basic issues. Technical Report ACT-CYC-129-90, Microelectronics and Computer Technology Corporation, Austin, Tex
- [Guh91] Guha R. (1991) Contexts: A formalization and some applications. Technical Report STAN-CS-91-1399-Thesis, Department of Computer Science, Stanford University
- [HK87] Hull R., King R. (1987) Semantic database modeling: survey, applications and research issues. *ACM Comput Surv 19:201–260*
- [HM93] Hammer J., McLeod D. (1993) An approach to resolving semantic heterogeneity in a federation of autonomous, heterogeneous, database systems. *Int J Intell Coop Inf Syst 2:51–84*
- [KBR86] Kaczmarek T., Bates R., Robins G. (1986) Recent developments in NIKL. In: *Proceedings AAAI-86*.
- [KCGS93] Kim W., Choi I., Gala S., Schevel M. (1993) On resolving schematic heterogeneity in multidatabase systems. *Distrib Parallel Databases Int J 1*
- [KLK91] Krishnamurthy R., Litwin W., Kent W. (1991) Language features for interoperability of databases with schematic discrepancies. In *Proceedings of 1991 ACM SIGMOD*.
- [KS91] Kim W., Seo J. (1991) Classifying schematic and data heterogeneity in multidatabase systems. *IEEE Comput 24 (12)*
- [KS93] Kashyap V., Sheth A. (1993) Schema correspondences between objects with semantic proximity. Technical Report DCS-TR-301, Department of Computer Science, Rutgers University
- [KS94a] Kashyap V., Sheth A. (1994) Semantics-based information brokering. In *Proceedings of the Third International Conference on Information and Knowledge Management (CIKM)*.
- [KS94b] Kashyap V., Sheth A. (1994) Semantics-based information brokering: a step towards realizing the infocosm. Technical Report DCS-TR-307, Department of Computer Science, Rutgers University
- [KS95a] Kashyap V., Sheth A. (1995) Controlled vocabulary sharing for query processing in global information systems. Technical report, LSDIS Lab, University of Georgia. <http://www.cs.uga.edu/LSDIS/infoquilt>.
- [KS95b] Kashyap V., Sheth A. (1995) Schematic and semantic similarities between database objects: a context-based approach. Technical Report TR-CS-95-001, LSDIS Lab, University of Georgia
- [KS96] Kashyap V., Sheth A. (1996) Semantic heterogeneity: role of metadata, context and ontologies. In: M. Papazoglou, G. Schlageter, (ed) *Cooperative Information Systems: Current Trends and Directions*. 1996.
- [LA86] Litwin W., Abdellatif A. (1986) Multidatabase interoperability. *IEEE Comput 19: 10–18*
- [LG90] Lenat D., Guha R.V. (1990) *Building large knowledge based systems: representation and inference in the Cyc Project*. Addison-Wesley, Reading, Mass.
- [LNE89] Larson J., Navathe S., Elmasri R. (1989) A theory of attribute equivalence in databases with application to schema integration. *IEEE Trans Software Eng 15*
- [vLNPS87] Luck K. von, Nebel B., Peltason C., Schmiedel A. (1987) The anatomy of the BACK system. Technical Report KIT Report 41, Technical University of Berlin
- [Mac87] MacGregor R. (1987) A deductive pattern matcher. In: *Proceedings AAAI-87*.
- [McC93] McCarthy J. (1993) Notes on formalizing context. In: *Proceedings of the International Joint Conference on Artificial Intelligence*.
- [Mic93] Michalski R. (1993) Inferential theory of learning as a conceptual basis for multistrategy learning. *Machine Learning 11*
- [MKIS96] Mena E., Kashyap V., Illarramendi A., Sheth A. (1996) Managing multiple information sources through ontologies: relationship between vocabulary heterogeneity and loss of information. In: *Proceedings of the workshop on Knowledge Representation meets Databases in conjunction with European Conference on Artificial Intelligence*.
- [MKSI96] Mena E., Kashyap V., Sheth A., Illarramendi A. (1996) OBSERVER: An approach for query processing in global information systems based on interoperation across pre-existing ontologies. In *Proceedings of the First IFCIS International Conference on Cooperative Information Systems (CoopIS '96)*.
- [ML92] Myaeng S.H., Li M. (1992) Building term clusters by acquiring lexical semantics from a corpus. In: *Proceedings of the CIKM*.
- [MS95] McLeod D., Si A. (1995) The design and experimental evaluation of an information discovery mechanism for networks of autonomous database systems. In: *Proceedings of the 11th IEEE Conference on Data Engineering*.
- [ON93] Ouksel A., Naiman C. (1993) Coordinating context building in heterogeneous information systems. *J Intell Inf Syst 3:151–183*
- [PM88] Peckham J., Maryanski J. (1988) Semantic data models. *ACM Comput Surv 20:153–190*
- [PS84] Patel-Schneider P. (1984) Small can be beautiful in knowledge representation. In: *Proceedings of the IEEE Workshop on Principle of Knowledge-Based Systems*.
- [RSK91] Rusinkiewicz M., Sheth A., Karabatis G. (1991) Specifying interdatabase dependencies in a multidatabase environment. *IEEE Comput 24:46–53*
- [SG89] Sheth A., Gala S. (1989) Attribute relationships: An impediment in automating schema integration. In: *Proceedings of the NSF Workshop on Heterogeneous Databases*.
- [She91] Sheth A. (1991) Semantic issues in multidatabase systems. *SIGMOD Record, special issue on Semantic Issues in Multidatabases 20 (12)*
- [Sho91] Shoham Y. (1991) Varieties of context.
- [SK92] Sheth A., Kashyap V. (1992) So far (schematically), yet so near (semantically). In: *Proceedings of the IFIP TC2/WG2.6 Conference on Semantics of Interoperable Database Systems, DS-5*, November 1992. In: IFIP Transactions A-25, North Holland
- [SL90] Sheth A., Larson J. (1990) Federated database systems for managing distributed, heterogeneous and autonomous databases. *ACM Comput Surv 22:183–236*
- [SPD92] Spaccapietra S., Parent C., Dupont Y. (1992) Model independent assertions for integration of heterogeneous schemas. *VLDB J 1:81–126*
- [SRK92] Sheth A., Rusinkiewicz M., Karabatis G. (1992) Using polytransactions to manage independent data. In: *Database Transaction Models*
- [SSR92] Sciore E., Siegel M., Rosenthal A. (1992) Context interchange using meta-attributes. In: *Proceedings of the CIKM*.
- [SZ90] Shaw G., Zdonik S. (1990) A query algebra for object-oriented databases. In: *Proceedings of the 6th IEEE Conference on Data Engineering*.
- [VD92] Voss D.A., Driscoll J.R. (1992) Text retrieval using a comprehensive lexicon. In: *Proceedings of the CIKM*.
- [Wie94] Wiederhold G. (1994) Interoperation, mediation and ontologies. In: *FGCS Workshop on Heterogeneous Cooperative Knowledge-Bases*.

[YSDK91] Yu C., Sun W., Dao S., Keirse D. (1991) Determining relationships among attributes for interoperability of multi-database systems. In: *Proceedings of the 1st International Workshop on Interoperability in Multidatabase Systems*.

## Appendix 1 Detailed specification of projection rules

$\text{semPro}(O_{1F}, O_1) = \langle \text{Cntxt}, M, (\text{dom}(O_{1F}), \text{dom}(O_1)), - \rangle$

$\prod_{\text{Cntxt}}(\text{semPro}(O_{1F}, O_1)) = \text{schCor}(O_{1F}, O_1)$   
 $= \langle O_{1F}, \{C_i \mid C_i \in \text{Cntxt}\}, O_1, \text{attr}(O_1), M \rangle$

Rule 1. *Empty Context Rule*, i.e.,  $\text{Cntxt} = \langle \rangle$

$\text{schCor}(O_{1F}, O_1) = \langle O_{1F}, \phi, O_1, \phi, M \rangle \Rightarrow M \equiv O_{1F} = O_1$

Rule 2. *Simple Sets Rule*, i.e.,  $\text{Cntxt} = \langle (C_1, S_1) \dots (C_k, S_k) \rangle$   $\text{schCor}(O_{1F}, O_1)$

$= \langle O_{1F}, \{C_i \mid C_i \in \text{Cntxt}\}, O_1, \{A_i \mid \text{map}_{O_1}(C_i, A_i) \text{ exists} \}, M \rangle$

$M \equiv O_{1F} = O\text{Select}(p, O_1)$ , where  $p \equiv (A_1 \in S_1) \wedge \dots \wedge (A_k \in S_k)$

Rule 3. *Simple Set Constraint Rule*, when  $\text{Cntxt} = \text{glb}(\langle (C_j, S_j) \rangle, \text{Cntxt}_1)$

$\text{semPro}(O_{1F}, O_1) = \text{semConstrain}(\langle (C_j, S_j) \rangle, \text{semPro}(O', O_1))$

where  $\text{semPro}(O', O_1)$  is defined *with respect to*  $\text{Cntxt}_1$  and

$O'$  is a temporary object obtained by applying constraints in  $\text{Cntxt}_1$  on

$O_1$   $\text{schCor}(O_{1F}, O_1) = \prod_{\text{Cntxt}_1}(\text{semConstrain}(\langle (C_j, S_j) \rangle, \text{semPro}(O', O_1)))$

$= \text{strConstrain}(\text{map}_{O_1}(C_j, A_j), S_j, \text{schCor}(O', O_1))$

where the mapping  $M'$  associated with  $\text{schCor}(O', O_1)$  is given by:

$M' \equiv O' = O\text{Select}(p, O_1)$

Rule 3.1. *New Constraint, Existing Attribute*,

i.e.,  $C_j \notin \text{Cntxt}_1$ ,  $\text{map}_{O_1}(C_j, A_j)$  exists.

The Mapping  $M$  associated with  $\text{schCor}(O_{1F}, O_1)$  is given as:  $M$

$\equiv O_{1F} = O\text{Select}((A_j \in S_j), O')$

$= O\text{Select}((A_j \in S_j), O\text{Select}(p, O_1))$

$= O\text{Select}((A_j \in S_j) \wedge p, O_1)$

Rule 3.2. *Existing Constraint, Existing Attribute*,

i.e.,  $C_j \in \text{Cntxt}_1$ ,  $\text{map}_{O_1}(C_j, A_j)$  exists.

Suppose  $(C_j, S'_j) \in \text{Cntxt}_1$ .

Then the mapping  $M'$  associated with  $\text{schCor}(O', O_1)$  may be written as:

$M' \equiv O' = O\text{Select}(p' \wedge (A_j \in S'_j), O_1)$  where  $p' \equiv p \wedge (A_j \in S_j)$ .

The mapping  $M$  associated with  $\text{schCor}(O_{1F}, O_1)$  is then given as:

$M \equiv O_{1F} = O\text{Select}((A_j \in S_j), O')$

$= O\text{Select}((A_j \in S_j), O\text{Select}(p' \wedge (A_j \in S'_j), O_1))$

$= O\text{Select}(p' \wedge (A_j \in S_j \cap S'_j), O_1)$

Rule 3.3. *New Constraint, Non-existing Attribute*,

i.e.,  $C_j \notin \text{Cntxt}_1$ ,  $\text{map}_{O_2}(C_j, A_j)$  does not exist.

The mapping  $M$  associated with  $\text{schCor}(O_{1F}, O_1)$  is given as:

$M \equiv O_{1F} = O\text{Product}(\text{makeObject}(C_j, S_j), O')$

$= O\text{Product}(\text{makeObject}(C_j, S_j), O\text{Select}(p, O_1))$

Rule 3.4. *New Constraint, Non-existing Attribute*,

i.e.,  $C_j \in \text{Cntxt}_1$ ,  $\text{map}_{O_2}(C_j, A_j)$  does not exist

Suppose  $(C_j, S'_j) \in \text{Cntxt}_1$ ,

then the mapping  $M'$  associated with  $\text{schCor}(O', O_1)$  may be written as:

$M' \equiv O' = O\text{Product}(\text{makeObject}(C_j, S_j), O\text{Select}(p', O_1))$

The mapping  $M$  associated with  $\text{schCor}(O_{1F}, O_1)$  can be then given as:

$M \equiv O_{1F} = O\text{Product}(\text{makeObject}(C_j, S_j), O')$

$= O\text{Product}(\text{makeObject}(C_j, S_j),$

$O\text{Product}(\text{makeObject}(C_j, S'_j), O\text{Select}(p', O_1)))$

$= O\text{Product}(\text{makeObject}(C_j, S_j \cap S'_j), O\text{Select}(p', O_1))$

Rule 4. *Context Conditioning Rule*, i.e.,

$\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O_{1F}, O_1))$

Rule 4.1. *Empty Context Conditioning Rule*, i.e.,  $\text{Cntxt}_1 = \langle \rangle$

$\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O_{1F}, O_1)) = \text{semPro}(O_{1F}, O_1)$

Rule 4.2. *Constraint Conditioning Rule*, i.e.  $\text{Cntxt}_1$

$= \text{glb}(\langle (C_j, S_j) \rangle, \text{Cntxt}_2)$

$\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O_{1F}, O_1))$

$= \text{semConstrain}(\langle (C_j, S_j) \rangle,$

$\text{semCondition}(\text{Cntxt}_2, \text{semPro}(O_{1F}, O_1)))$

$\prod_{\text{Cntxt}_1}(\text{semConstrain}(\langle (C_j, S_j) \rangle,$

$\text{semCondition}(\text{Cntxt}_2, \text{semPro}(O_{1F}, O_1))))$

$= \text{strConstrain}(\text{map}_{O_2}(C_j, A_j), S_j,$

$\prod_{\text{Cntxt}_2}(\text{semCondition}(\text{Cntxt}_2, \text{semPro}(O_{1F}, O_1))))$

Rule 4.3. *Context Conditioning and semCombine Rule*, i.e.,

$\text{semCondition}(\text{Cntxt}_1, \text{semCombine}(C_i,$

$\text{semPro}(O', O_1), \text{semPro}(O'', O_i)))$

$= \text{semCombine}(C_i, \text{semCondition}(\text{Cntxt}_1, \text{semPro}(O', O_1)),$

$\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O'', O_i)))$

$\prod_{\text{Cntxt}_1}(\text{semCombine}(C_i,$

$\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O', O_1)),$

$\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O'', O_i))))$

$= \text{strCombine}(\{\text{map}_{O_1}(C_i, A_i), \text{map}_{O_i}(C_i, A'_i)\},$

$\prod_{\text{Cntxt}_1}(\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O', O_1)),$

$\prod_{\text{Cntxt}_1}(\text{semCondition}(\text{Cntxt}_1, \text{semPro}(O'', O_i))))$

Rule 5. *semCombine Rule*, i.e.,

$\text{Cntxt} = \text{glb}(\langle (C_i, O_i \circ C_{\text{ass}}(O_i, O_1)) \rangle, \text{Cntxt}_1)$

$\text{semPro}(O_{1F}, O_1) =$

$\text{semConstrain}(\langle (C_i, O_i \circ C_{\text{ass}}(O_i, O_1)) \rangle, \text{semPro}(O', O_1))$

$= \text{semCombine}(C_i, \text{semPro}(O', O_1),$

$\text{semCondition}(C_{\text{ass}}(O_i, O_1), \text{semPro}(O_{iF}, O_i)))$

where  $\text{semPro}(O', O_1)$  is defined *with respect to*  $\text{Cntxt}_1$  and  $O'$  is a temporary object obtained by applying all the constraints in  $\text{Cntxt}_1$  to  $O_1$

$\prod_{\text{Cntxt}}(\text{semCombine}(C_i, \text{semPro}(O', O_1),$

$\text{semCondition}(C_{\text{ass}}(O_i, O_1), \text{semPro}(O_{iF}, O_i))))$

$= \text{strCombine}(\{\text{map}_{O_1}(C_i, A_i), \text{map}_{O_i}(C_i, A'_i)\},$

$\prod_{\text{Cntxt}_1}(\text{semPro}(O', O_2)),$

$\prod_{C_{\text{ass}}(O_i, O_1)}(\text{semCondition}(C_{\text{ass}}(O_i, O_1), \text{semPro}(O_{iF}, O_i))))$

$= \text{strCombine}(\{\text{map}_{O_i}(C_i, A'_i), \text{map}_{O_1}(C_i, A_i)\},$

$\text{schCor}(O', O_1), \text{schCor}(O'', O_i))$

where  $O''$  is a temporary object obtained by applying all the constraints in  $C_{\text{ass}}(O_i, O_1)$  to  $O_{iF}$

and the mappings  $M'$  and  $M''$  associated with  $\text{schCor}(O', O_1)$  and  $\text{schCor}(O'', O_i)$  are given as:

$M' \equiv O' = O\text{Select}(p', O_1)$   $M'' \equiv O'' = O\text{Select}(p'', O_i)$

Rule 5.1. *New Constraint and Existing Attributes*,

i.e.,  $C_i \notin \text{Cntxt}_1$ ,  $\text{map}_{O_i}(C_i, A'_i)$  and  $\text{map}_{O_1}(C_i, A_i)$  exist.

$M \equiv O_{1F} = O\text{Join}(g(A_j, A'_j), O', O'')$

$= O\text{Join}(g(A_i, A'_i), O\text{Select}(p', O_1), O\text{Select}(p'', O_i))$

Rule 5.2. *Coordinate Composition Rule*, i.e.,  $C_i = \text{compose}(C_{i,1}, C_{i,2})$

The composition of attributes is as follows:

$\text{map}_O(C_i, X) = \text{map}_O(\text{compose}(C_{i,1}, C_{i,2}), \text{compose}(X_1, X_2))$

$= \text{compose}(\text{map}_O(C_{i,1}, X_1), \text{map}_O(C_{i,2}, X_2))$

Let  $\text{map}_{O_1}(C_i, A_i)$

$= \text{compose}(\text{map}_{O_1}(C_{i,1}, A_{i,1}), \text{map}_{O_1}(C_{i,2}, A_{i,2}))$

Let  $\text{map}_{O_i}(C_i, A'_i)$

$= \text{compose}(\text{map}_{O_i}(C_{i,1}, A'_{i,1}), \text{map}_{O_i}(C_{i,2}, A'_{i,2}))$

The mapping  $M$  associated with  $\text{schCor}(O_{1F}, O_1)$  is given as:

$M \equiv O_{1F} = O\text{Join}(g(\langle A_{i,1}, A_{i,2} \rangle, \langle A'_{i,1}, A'_{i,2} \rangle), O', O'')$

$= O\text{Join}(g(\langle A_{i,1}, A_{i,2} \rangle, \langle A'_{i,1}, A'_{i,2} \rangle),$

$O\text{Select}(p', O_1), O\text{Select}(p'', O_i))$

## Appendix 2 Taxonomies of schematic conflicts

In this section, we enumerate the various types of schematic/representational conflicts identified by us in the taxonomy proposed in this paper. We take a representative sample of the multidatabase literature in this area and show the relationship of their work with ours by means of a table (Table 3). We believe this paper provides a more complete enumeration of the various types of conflicts and their definitions.

**Table 3.** Comparison of the types of conflicts. We use the symbol  $\alpha$  to denote that the reference has an informal discussion of the schematic conflict. We use the symbol  $\beta$  to denote that the schematic conflict has been defined formally

<b>Schematic conflicts</b>	[DH84]	[CRE87]	[SPD92]	[SK92]	[KCGS93]	[HM93]
<b>Domain incompatibilities</b>		$\beta$	$\alpha$	$\alpha$		
Naming conflicts	$\beta$	$\alpha$	$\beta$	$\beta$	$\beta$	$\alpha$
Data representation conflicts		$\alpha$		$\beta$	$\beta$	
Data scaling conflicts	$\beta$		$\alpha$	$\beta$	$\beta$	$\beta$
Data precision conflicts				$\beta$	$\beta$	
Default value conflicts				$\beta$	$\beta$	$\alpha$
Attribute integrity constraint conflicts			$\alpha$	$\beta$	$\beta$	$\alpha$
<b>Entity definition incompatibilities</b>		$\beta$		$\alpha$	$\alpha$	
Database identifier conflicts	$\alpha$			$\beta$	$\beta$	
Naming conflicts	$\beta$	$\alpha$		$\beta$	$\beta$	$\beta$
Schema isomorphism conflicts	$\alpha$	$\alpha$	$\beta$	$\beta$	$\beta$	$\alpha$
Missing data item conflicts	$\beta$			$\beta$	$\beta$	$\alpha$
<b>Data value incompatibilities</b>	$\alpha$			$\alpha$	$\alpha$	
Known inconsistency	$\beta$			$\beta$	$\beta$	
Temporary inconsistency	$\beta$			$\beta$	$\beta$	
Acceptable inconsistency				$\beta$		
<b>Abstraction level incompatibilities</b>	$\alpha$			$\alpha$	$\alpha$	
Generalization conflicts	$\beta$		$\beta$	$\beta$	$\beta$	$\beta$
Aggregation conflicts	$\beta$		$\alpha$	$\beta$	$\beta$	$\beta$
<b>Schematic discrepancies</b>				$\alpha$		
Data value attribute conflict				$\beta$		
Attribute entity conflict	$\alpha$		$\beta$	$\beta$	$\beta$	
Data value entity conflict				$\beta$		