# **Reference Model and Perspective schemata inference for Enterprise Data Integration**

Valéria Magalhães Pequeno<sup>1</sup>, João Carlos Moura Pires<sup>1</sup>

<sup>1</sup>CENTRIA, Departamento de Informática, Faculdade de Ciências e Tecnologia, FCT, Universidade Nova de Lisboa 2829-516, Caparica, Portugal

vmp@fct.unl.pt, jmp@di.fct.unl.pt

**keywords:** ETL process, conceptual data model, data integration, correspondence assertions, object-relational database, rewriting-rules

#### Abstract

One of leading issues in database research is to develop flexible mechanisms for providing integrated access to distributed, heterogeneous databases and other information sources. A wide range of techniques has been developed to address this problem, the main drawback being the difficulty in developing a single (global) database schema that captures all the nuances of diverse data types, and expresses a unified view of the enterprise. We deal with this problem by taking a declarative approach, which is based on the creation of a reference model and perspective schemata. The former provides a common semantic, while the latter connects schemata. This paper focus on deduction of new perspective schemata using a proposed inference mechanism.

### 1. Introduction

One of the leading issues in database research is to develop flexible mechanisms for providing integrated access to multiple, distributed, heterogeneous databases and other information sources. A wide range of techniques has been developed to address this problem, including approaches based on creation of <u>Data Warehouses</u> (DWs), and <u>Federated Database Systems</u> (FDBs). DWs are highly specialised database systems which contain the unified history of an enterprise at a suitable level of detail for decision support. All data are integrated into, usually, a single repository, with a generalised and global schema. A FDB enables a unified virtual view of one or more autonomous sources of information to hide data heterogeneity from the applications and users. Tightly coupled FDB, those that occur in DWs, provides a global schema expressed in a common, "canonical" data model. Unlike a DW, a FDB leaves data at the source.

One of the main drawbacks of these approaches is the difficulty in developing a single (global or common) database schema that captures all the nuances of diverse data types, and expresses a unified view of the enterprise. The designer should usually deal with incompatible data models, characterised by subtle differences in structure and semantic. Besides, he/she should define mappings between the global schema and the schemata of the source information. These problems are hardest to deal with because of the rapid growing of the data volume and the data model complexity (both in sources and in global schema), which implies the rise of the difficulty of managing and understanding these models [1].



Figure 1: Proposed architecture.

In order to deal with this problem, it is proposed to take a declarative approach, based on the creation of a reference model and perspective schemata. A Reference Model is an abstract framework that provides a common semantic that can be used to guide the development of other models and help with data consistency [1]. A perspective schema describes a data model, part or whole (target schema), in terms of other data models (base schemata). In Fig. 1, for instance,  $P_{s'1|RM}$ , ...,  $P_{s'n|RM}$  are perspective schemata that map the reference model (**RM**) in terms of the source schemata.

In the proposed approach, the relationship between the base schemata and the target schema is made explicitly and declaratively through correspondence assertions. An advantage of the proposed approach is that by using the reference model the designer does not need to map schemata each other. This effort is theoretically reduced since schemata (source or global) must only align with the reference model, rather than each participating schema. Besides, the designer does not need to have a deep knowledge of all schemata involved in the DW system or in the federation system. Thus, the designer can describe the global system without concerns about where the sources are or how they are stored. Furthermore, the mapping between the global schema and its sources is automatically generated by an inference mechanism. This paper focuses on the deduction of new perspective schemata using a proposed inference mechanism.

The remainder of this paper is laid out as follows. Section 2 presents an overview of the reference model-based framework proposed in [2]. Section 3 briefly describes the language to define the perspective schemata. Section 4 details the process to infer new perspective schemata. Section 5 concisely mentions representative works in data integration area. The paper ends with Section 6, which points out the new features of the approach presented here and in ongoing or planned future work on this topic.

# 2. The framework

The proposal presented in [2] offers a way to express the existing data models (source, reference model, and global/integrated schema) and the relationship between them. The approach is based on *Schema language* ( $L_S$ ) and *Perspective schema language* ( $L_{PS}$ ).

Schema language  $(L_S)$  is used to describe the actual data models (source, reference model, and global/integrated schema). The formal framework focuses on an object-relational paradigm, which includes definitions adopted by the main concepts of object and relational models as they are widely ac-



Figure 2: Motivating example.

cepted in literature - cf. [3, 4].

Perspective schema language ( $L_{PS}$ ) is used to describe perspective schemata. A perspective schema is a special kind of schema that describes a data model (part or whole) (*target schema*) in terms of other data models (*base schemata*).  $L_{PS}$  mainly extends  $L_S$  with two components: Correspondence Assertions (CAs) and Matching Functions (MFs). Correspondence Assertions formally specify the relationship between schema components. Matching functions indicate when two data entities represent the same instance of the real world.  $L_{PS}$  includes data transformations, such as names conversion and data types conversion.

Fig. 1 illustrates the basic components of the proposed architecture and their relationships. The schemata **RM**, **S**<sub>1</sub>,...,**S**<sub>n</sub> and **G** are defined using the language L<sub>S</sub> and represent, respectively, the reference model, the source schemata **S**<sub>1</sub>,...,**S**<sub>n</sub>, and a global schema. The schemata **S**'<sub>1</sub> and **S**'<sub>2</sub> are defined using the language L<sub>PS</sub>. They are special kinds of perspective schemata (called *view schema*), since the target schema is described in the scope of a perspective schema, instead of just referring to an existing schema. **S**'<sub>1</sub> and **S**'<sub>2</sub> represent, respectively, the view schemata **S**'<sub>1</sub> (a viewpoint of schema **S**<sub>1</sub>), and **S**'<sub>2</sub> (an integrated viewpoint of schemata **S**<sub>2</sub> and **S**<sub>3</sub>). The relationships between the target schema and the base schemata are shown through the perspective schemata **P**<sub>s'1|RM</sub>,..., **P**<sub>s'n|RM</sub>,**P**<sub>RM|G</sub>, and **P**<sub>s'1,s'2,...,s'n|G</sub> (denoted by arrows). In the current research, the perspective schema **P**<sub>s'1,s'2,...,s'n|G</sub> can be automatically deduced by the proposed inference mechanism. The next Section illustrates, through the examples, the language L<sub>PS</sub>, and the Section 4 presents the proposed inference mechanism. For a more detailed and formal description of L<sub>S</sub> and L<sub>PS</sub> languages, the reader is referred to [5, 6, 2].

### 3. Perspective Schema Language

The remainder of the paper, considers a simple sales scenario comprising two data sources  $S_1$  and  $S_2$ , a reference model **RM**, and a global schema **G**. The schemata are shown in Fig. 2. All properties that are key to a relation (or class) are shown in Fig. 2 using "#" before their names.

The language  $L_{PS}$  is used to define perspective schemata. A perspective schema describes a data

model, part or whole (*target schema*), in terms of other data models (*base schemata*). Usually, a perspective schema is formed by the following components:

- Name is a schema name with the notation: P<sub>S|T</sub>, being S the name of one or more base schemata and T the name of the target schema. In Fig. 2, for instance, P<sub>s'1|RM</sub> is a name of a perspective schema whose base schema is S'<sub>1</sub> and the target schema is RM;
- 2. '*Require' declarations* express the subset of the components of the target schema (classes, relations, keys, and foreign keys) that will be necessary in the perspective schema;
- 3. <u>Matching Function signatures</u> indicate which matching functions must be implemented to determine when two objects/tuples are distinct representations of the same object in the real-world;
- 4. <u>Correspondence Assertions</u> establish the semantic correspondence between schemata's components.

The target schema may have much more information than is required to represent in a perspective schema, namely when the target schema is the Reference Model. Hence, it is required to clearly indicate which elements of the target schema are in the scope of the perspective schema. This is done in  $L_{PS}$  using 'require' declarations. For instance, consider the perspective schema  $P_{S_2|RM}$  between the schemata **RM** (the target schema) and  $S_2$  (the base schema), both as presented in Fig. 2. For this perspective schema, four relations from **RM** are needed (PRODUCT, CUSTOMER, SALE, and SALE\_ITEM). The 'require' declaration to relation CUSTOMER, for example, would be as follows:

# require(CUSTOMER, { $cid_{RM}$ , cname<sub>RM</sub>, cphone<sub>RM</sub>})

Note that, for instance, the properties  $cregion_id_{RM}$  and  $caddress_{RM}$  from RM.CUSTOMER are not declared as being required.

## 3.1. Matching Functions

From a conceptual viewpoint, it is essential to provide a way to identify instances of different schemata that represent the same entity in the real-world. The proposal presented in [2] is to use matching functions, which can include various techniques for matching instances, including some of those used in data cleaning, such as lookup tables, user-defined functions, heuristics and past matching. These functions, as occur in [7], define a 1:1 correspondence between the objects/tuples in families of corresponding classes/relations. In particular, the work shown in [2] is based on the following matching function signature:

$$match: ((\mathbf{S}_1 [\mathbf{R}_1], \tau_1) \times (\mathbf{S}_2 [\mathbf{R}_2], \tau_2)) \to Boolean,$$
(1)

being  $\mathbf{S}_i$  schema names,  $R_i$  class/relation names, and  $\tau_i$  the data type of the instances of  $R_i$ , for  $i \in \{1,2\}$ . When both arguments are instanced, **match** verifies whether two instances are semantically equivalent or not. If only one argument is instanced, e.g.  $\mathbf{S}_1.R_1$ , then it obtains the semantically equivalent  $\mathbf{S}_2.R_2$ instance of the given  $\mathbf{S}_1.R_1$  instance, returning true when it is possible, and false when nothing is found or when there is more than one instance to match.

In some scenarios one-to-many correspondence between instances are common, e.g. when historical data is stored in the DW. In this case, a variant of **match** should be used, which has the following form:

$$match: ((\mathbf{S}_1[\mathbf{R}_1], \tau_1) \times (\mathbf{S}_2[\mathbf{R}_2(\mathsf{predicate})], \tau_2)) \to \text{Boolean}.$$
(2)

**predicate** is a boolean condition that determines the context in which the instance matching must be applied in  $S_2$ .R<sub>2</sub>.

An example of a matching function signature involving schemata of Fig. 2 is presented in Fig. 3. The implementation of the matching functions shall be externally provided, since their implementation is very close to the application domain. However, in order to make easer the implementation of a simple prototype, a new variants of **match** is introduced in  $L_{PS}$ :

match :  $((\mathbf{S}_1[\mathbf{R}_1], \tau_1, \{\mathbf{p}'_1 : \tau'_1, ..., \mathbf{p}'_n : \tau'_n\}) \times (\mathbf{S}_2[\mathbf{R}_2], \tau_2, \{\mathbf{p}''_1 : \tau''_1, ..., \mathbf{p}''_n : \tau''_n\})) \rightarrow \text{Boolean}$  (3)

being that  $p'_i:\tau'_i \in type(R_1)$ ; and  $p''_i:\tau''_i \in type(R_2)$ ,  $1 \le i \le n$ .<sup>1</sup> This variant of the matching function is automatically generated by the system and indicates that the matching is done by simple attribute comparison, i.e. each property  $p'_i$  of  $R_1$  will be compared with the property  $p''_i$  of  $R_2$ , for  $1 \le i \le n$ .

match:((**RM**[CUSTOMER], $\tau_1$ )×(**G**[CUSTOMER], $\tau_2$ ))→Boolean

Figure 3: Example of a matching function signature.

#### **3.2.** Correspondence Assertions

The semantic correspondence between schemata's components is declared in the proposal presented in [2] through the <u>C</u>orrespondence <u>A</u>ssertions (CAs), which are used to formally assert the correspondence between schema components in a declarative fashion. CAs are classified in four groups: <u>P</u>roperty <u>C</u>orrespondence <u>A</u>ssertion (PCA), <u>Extension C</u>orrespondence <u>A</u>ssertion (ECA), <u>S</u>ummation <u>C</u>orrespondence <u>A</u>ssertion (SCA), and <u>Aggregation C</u>orrespondence <u>A</u>ssertion (ACA). Examples of CAs are shown in Fig. 4 and explained in this Section.

Property Correspondence Assertions (PCAs)		
$\psi_1$ :	$\mathbf{P}_{\mathbf{RM} \mathbf{G}} [\text{CUSTOMER}] \bullet idcard_{G} \rightarrow numberTOtext (\mathbf{RM} [\text{CUSTOMER}] \bullet cid_{RM})$	
$\psi_2$ :	$\mathbf{P}_{\mathbf{RM} \mathbf{G}}[\text{CUSTOMER}] \bullet \mathbf{contact}_{G} \rightarrow \mathbf{RM}[\text{CUSTOMER}] \bullet \mathbf{cphone}_{RM}$	
Extension Correspondence Assertions (ECAs)		
$\psi_3$ :	$P_{RM G}[\mathrm{CUSTOMER}]  o RM[\mathrm{CUSTOMER}]$	
$\psi_4$ :	$\mathbf{S}_{v}$ [CUSTOMER] $\rightarrow \mathbf{S}_{1}$ [CUSTOMER] $\supset \mathbf{S}_{2}$ [CUSTOMER]	
Summation Correspondence Assertion (SCA)		
$\psi_5$ :	$\mathbf{P}_{\mathbf{S}_{3} \mathbf{R}\mathbf{M}} [PRODUCT] (\mathbf{pid}_{RM}) \rightarrow normalise (\mathbf{S}_{3} [PRODUCT\_SALES] (\mathbf{product\_number}_{S3}))$	

Figure 4: Examples of correspondence assertion.

Property CAs relate properties of a target schema to the properties of base schemata. They allow dealing with several kinds of semantic heterogeneity such as: *naming conflict* (for instance synonyms and homonyms properties), *data representation conflict* (that occur when similar contents are represented by different data types), and *encoding conflict* (that occur when similar contents are represented by different formats of data or unit of measures). For example, the PCAs  $\psi_1$  and  $\psi_2$  (see Fig. 4) deal with, respectively, *data representation conflict* and *naming conflict*.  $\psi_1$  links the property **idcard**<sub>G</sub> to

<sup>&</sup>lt;sup>1</sup>*type*() is a function defined in language  $L_S$  that returns the structural type of a relation/class.

The Extension CAs are used to describe which objects/tuples of a base schema should have a corresponding semantically equivalent object/tuple in the target schema. For instance, the relation **G**.CUSTOMER is linked to relation **RM**.CUSTOMER through the ECA  $\psi_3$  presented in Fig. 4.  $\psi_3$  determines that **G**.CUSTOMER and **RM**.CUSTOMER are equivalent, i.e., for each tuple of CUSTOMER of the schema **RM** there is one semantically equivalent tuple in CUSTOMER of the schema **G**, and vice-versa.

There are five different kinds of ECAs: equivalence, selection, difference, union, and intersection, being the ECA of union similar to the *natural outer-join* of the usual relational models. For instance, consider the view schema  $S_v$  with the relation CUSTOMER, which is related to the relation CUSTOMER of the schema  $S_1$  and to the relation with the same name of  $S_2$  through the ECA  $\psi_4$  shown in Fig. 4.  $\psi_4$ determines that CUSTOMER in  $S_v$  is the union/join of CUSTOMER in  $S_1$  and CUSTOMER in  $S_2$ , i.e., for each tuple of CUSTOMER of the schema  $S_1$  there is one semantically equivalent tuple in CUSTOMER of the schema  $S_v$ , or for each tuple of CUSTOMER of the schema  $S_2$  there is one semantically equivalent tuple in CUSTOMER of the schema  $S_v$ , and vice-versa. In an ECA, any relation/class can appear with a selection condition, which determines the subset of instances of the class/relation being considered. This kind of ECA is especially important to the DW because through it the current instances of the DW can be selected and related to the instances of their sources (which usually do not have historical data).

The Summation CAs are used to describe the summary of a class/relation whose instances are related to the instances of another class/relation by breaking them into logical groups that belong together. They are used to indicate that the relationship between classes/relations involve some type of aggregate functions (called SCA of groupby) or a normalisation process (called SCA of normalisation)<sup>2</sup>. For example, consider the source schema  $S_3$  (not presented in any figure), which contains a denormalised relation PRODUCT\_SALES(**product\_number**<sub>S3</sub>, **product**<sub>S3</sub>, **quantity**<sub>S3</sub>, **price**<sub>S3</sub>, **purchase\_order**<sub>S3</sub>) and the schema **RM** presented in Fig. 2. PRODUCT\_SALES holds information about sold items in a purchase order as well as information logically related to products themselves, which could be in another relation, occurring in schema **RM**. The SCA  $\psi_5$ , displayed in Fig. 4, determines the relationship between PRODUCT\_SALES and **RM**.PRODUCT when a normalisation process is involved, i.e., it determines that **RM**.PRODUCT is a normalisation of **S**<sub>3</sub>.PRODUCT\_SALES based on distinct values of property **product\_number**<sub>S3</sub>.

The Aggregation CAs link properties of the target schema to the properties of the base schema when a SCA is used. ACAs associated to SCAs of groupby contains aggregation functions supported by most of the queries languages, like SQL-99 [8], i.e. *summation, maximum, minimum, average* and *count*. The ACAs, similar to the PCAs, allow for the description of several kinds of situations; therefore, the aggregate expressions can be more detailed than simple property references. Calculations performed can include, for example, ordinary functions (such as sum or concatenate two or more properties' values before applying the aggregate function), and Boolean conditions (e.g. count all male students whose grades are greater or equal to 10).

<sup>&</sup>lt;sup>2</sup>This research also deal with denormalisations, which is defined using *path expressions* (component of the language  $L_S$ ).

## 4. Inference Mechanism

This proposal provides an inference mechanism to automatically infer a new perspective schema (see Fig. 5(c)), given:

- 1. a set of *origin* schemata and their associated perspective schemata, which take the *origin* schemata as *base* and the reference model as *target* (see Fig. 5(a));
- 2. a *destination* schema and its associated perspective schema, which take the reference model as *base* and the *destination* schema as *target* (see Fig. 5(b)).

In context of the Fig. 1, the perspective schema  $\mathbf{P}_{s'1,s'2,...,s'n|G}$  can be inferred taking as *origin* the schemata  $\mathbf{S}_{1,...,\mathbf{S}_{n}}$  as well as the perspective schemata  $\mathbf{P}_{s'1|RM},...,\mathbf{P}_{s'n|RM}$ , and as *destination* the schema **G** as well as the perspective schema  $\mathbf{P}_{RM|G}$ .

The inferred perspective schema will have as *base* a subset of *origin* schemata, and as *target* the *destination* schema. Its '*require*' *declarations* will be the same '*require*' *declarations* present in the perspective schema associated to the *destination* schema. The *MF signatures* and *CAs* of the inferred perspective schema will be automatically generated using a rule-based rewritten system.



Figure 5: Sketch of the inference mechanism.

The rule-based rewriting system is formed by a set of rules having the general form:

**Rule:** 
$$\frac{X \Rightarrow Y}{Z}$$
 (read X is rewritten in Y if Z is valid), (4)

In (4), **Rule** is the name of the rule. *X* and *Y* can be formed by any of the following expressions: a CA pattern expression, a MF pattern signature, or a component pattern expression. CA pattern expressions and MF pattern signatures are expressions conforming to the  $L_{PS}$  syntax to declare, respectively, CAs and MF signatures, being that some of their elements are variables to be used in a unification process. Component pattern expressions, functions with n-ary arguments, values, or conditions of selection (predicates), being that some of their elements are variables to be used in a unification process. *Z* is a condition formed by a set of CA pattern expressions, or expressions of the forms: a)  $A \Rightarrow B$  such that *A* and *B* are component pattern expressions; b) (*FK*, C, ., C', .) such that *FK* is a foreign key name of a class/relation C' that refers to a class C.<sup>3</sup> CA pattern expressions, MF pattern signatures, and component pattern expressions are formally defined in the following text.

 $<sup>{}^{3}</sup>type()$  is a function defined in language L<sub>S</sub> that returns the structural type of a relation/class.

**Definition 1** (*CA pattern expression*) Let  $\hat{A}$  be a set of correspondence assertions defined in language  $L_{PS}$ . A CA pattern expression is an expression having the general form:

 $\underline{K} \to \underline{L}$ 

with K and L being variables that can be instatiated with, respectively, the left-side and the right-side of a correspondence assertion  $\hat{\psi} \in \hat{A}$ .

**Definition 2** (*MF* pattern signature) Let T a set of data types and  $\hat{\mathcal{L}}$  a set of schemata. A MF pattern signature is an expression having one of the following forms:

$$\begin{array}{l} \textit{match} : \left( \left( \underline{\mathbf{S}}_{1} \left[ \underline{\mathbf{C}}_{1} \right], \underline{\tau}_{1} \right) \times \left( \underline{\mathbf{S}}_{2} \left[ \underline{\mathbf{C}}_{2} \right], \underline{\tau}_{2} \right) \right) \rightarrow \textit{Boolean}; \\ \textit{match} : \left( \left( \underline{\mathbf{S}}_{1} \left[ \underline{\mathbf{C}}_{1} \right], \underline{\tau}_{1} \right) \times \left( \underline{\mathbf{S}}_{2} \left[ \underline{\mathbf{C}}_{2} \left( \underline{\textit{pred}} \right) \right], \underline{\tau}_{2} \right) \right) \rightarrow \textit{Boolean}; \\ \textit{match} : \left( \left( \underline{\mathbf{S}}_{1} \left[ \underline{\mathbf{C}}_{1} \right], \underline{\tau}_{1}, \{ \underline{p'}_{1} : \underline{\tau'}_{1}, \dots, \underline{p'}_{n} : \underline{\tau'}_{n} \} \right) \times \left( \underline{\mathbf{S}}_{2} \left[ \underline{\mathbf{C}}_{2} \right], \underline{\tau}_{2}, \{ \underline{p''_{1}} : \underline{\tau''_{1}}, \dots, \underline{p''_{n}} : \underline{\tau''_{n}} \} \right) \right) \rightarrow \textit{Boolean}; \\ \textit{match} : \left( \left( \underline{\mathbf{S}}_{1} \left[ \underline{\mathbf{C}}_{1} \right], \underline{\tau}_{1}, \{ \underline{p'}_{1} : \underline{\tau'}_{1}, \dots, \underline{p'_{n}} : \underline{\tau''_{n}} \} \right) \times \left( \underline{\mathbf{S}}_{2} \left[ \underline{\mathbf{C}}_{2} \left( \underline{\textit{pred}} \right) \right], \underline{\tau}_{2}, \{ \underline{p''_{1}} : \underline{\tau''_{1}}, \dots, \underline{p''_{n}} : \underline{\tau''_{n}} \} \right) \right) \rightarrow \\ \rightarrow \textit{Boolean} \end{array}$$

With  $\mathbf{S}_1$  and  $\mathbf{S}_2$  being variables that can be instatiated with any of the schemata belonging to  $\hat{\mathcal{L}}$ ;  $\mathbf{C}_1$  and  $\mathbf{C}_2$  being variables that can be instantiated with any class/relation of the schemata belonging to  $\hat{\mathcal{L}}$ ; **pred** being a variable that can be instantiated with a predicate (as defined in  $L_{PS}$ );  $\mathbf{p}'_i$  and  $\mathbf{p}''_i$ , for  $1 \le i \le n$ , are variables that can be instantiated with any property of a class/relation of the schemata belonging to  $\hat{\mathcal{L}}$ ;  $\tau_1$ ,  $\tau_2$ ,  $\tau'_i$ ,  $\tau''_i$ , for  $1 \le i \le n$ , are variable that can be instantiated with any property of a class/relation of the schemata belonging to  $\hat{\mathcal{L}}$ ;  $\tau_1$ ,  $\tau_2$ ,  $\tau'_i$ ,  $\tau''_i$ , for  $1 \le i \le n$ , are variable that can be instantiated with any property of a class/relation of the schemata belonging to  $\hat{\mathcal{L}}$ ).

**Definition 3** (*Component pattern expression*) Let  $\hat{\mathcal{L}}$  be a set of schemata. A Component pattern expression is a expression formed by a single variable that can be instantiated with a predicate **pred**; or is an expression having one of the following forms:

$$\underline{\underline{S}} [\underline{\underline{C}}_1] \bullet \underline{p}$$

$$\underline{\underline{S}} [\underline{\underline{C}}_1] \bullet \underline{\varrho}$$

$$\underline{\varphi} (\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n)$$

$$\underline{\underline{S}} [\underline{\underline{C}}_1] \bullet \underline{p} \{\underline{p}''\}$$

$$\underline{\ell}_1 : \underline{\underline{C}}_1 \to \underline{\underline{C}}_2$$

With **S** being a variable that can be instatiated with any of the schemata belonging to  $\hat{\mathcal{L}}$ ;  $C_1$  and  $C_2$  being variables that can be instantiated with any class/relation of the schemata belonging to  $\hat{\mathcal{L}}$ ;  $\varrho$  being a variable that can be instantiated with a (value or reference) path expression as defined in  $L_S$ ;  $\ell_1$  being a variable that can be instantiated with a link of a path expression;  $\mathbf{p}$  and  $\mathbf{p}''$  being variables that can be instantiated with a link of the schemata belonging to  $\hat{\mathcal{L}}$ , being that  $\mathbf{p}''$  is part of the structural type of  $\mathbf{p}$ ;  $\varphi$  being a variable that can be instantiated with  $n \ge 1$  arguments that returns a value; and  $X_i$ ,  $1 \le i \le n$ , being variables that can be instantiated with other component pattern expressions.

A condition Z is valid when all of its expressions are valid: a) the CA pattern expression is valid if there is a CA, which is declared in one of the perspective schemata associated to the *origin* schemata or the *destination* schema, that unifies with it; b) the expression of the form  $A \Rightarrow B$ , such that A and B are component pattern expressions, is valid if there is a rule which unifies with it and which is recursively applied; c) the expression of the form  $(FK,C,\_,C',\_)$  is valid if there is a foreign key declaration declared as required in one of the perspective schemata associated to the *origin* schemata that unifies with it; d) the expression of the form  $\mathbf{p}: \natural C \in type(C')$  is valid if there are both a class/relation and a property, both declared as required in one of the perspective schemata associated to the *origin* schemata, that unify with them. A formal definition of a valid condition is as follows:

**Definition 4** (Valid condition (in an inference rule)) Let  $\hat{\mathcal{L}}_p$  be a set of perspective schemata and  $\hat{\mathcal{A}}_i$ a set of correspondence assertions declared in some perspective schema belonging to  $\hat{\mathcal{L}}_p$ . Let also  $\{X_1, X_2, ..., X_n\}$  be a condition Z in an inference rule. Z is a valid condition iff for each  $X_i$ ,  $1 \le i \le n$ , if:

- $X_i$  is a CA pattern expression, then there is a  $\hat{\psi} \in \hat{A}_i$  that unifies with  $X_i$ ;
- $X_i$  is an expression the form  $A \Rightarrow B$ , such that A and B are component pattern expressions, then there is an inference rule that unifies with  $X_i$  whose condition is a valid condition,
- X<sub>i</sub> is an expression of the form (FK,C,\_,C',\_) then there is a foreign key declaration, as defined in L<sub>S</sub>, declared as required in some perspective schema belonging to L̂<sub>p</sub> that unifies with X<sub>i</sub>;
- $X_i$  is an expression of form  $p: \natural C \in type(C')$ , then there are both a class/relation and a property, both declared as required in some perspective schema belonging to  $\hat{\mathcal{L}}_p$ , such that  $p: \natural C \in type(C')$ .

When X and Y are CA pattern expressions, the rule are rewritten-rules that rewrite CAs in other CAs (RR-CAs). When X and Y are MF pattern expressions, the rule are rewritten-rules that rewrite MFs in other MFs (RR-MFs). When X and Y are component pattern expressions, the rule are substitution-rules that rewrite components in other components (RR-Cs). The latter are used as an intermediary process by the RR-CAs and RR-MFs.

An example of a RR-CA is as follows:

$$\mathbf{RR-CA1:} \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \to \underline{\mathbf{K}}^{S}}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \to \underline{\mathbf{K}}^{S}}.$$
(5)

In (5) all variables are indicated by an underline. **D** is the *destination* schema, **RM** is the reference model schema, and **S** is a variable that will be instantiated with some of the *origin* schemata.  $C^D$  is a variable that will be instantiated with a class/relation of the schema **D**; mutatis mutandis to  $C^{RM}$ . **K** is a variable that will be instantiated with the right side of a CA pattern expression of extension. The letter S in **K**<sup>S</sup> means that all elements in that expression belong to schema **S**. The value of **S** and **K** will depend on which CA, that is declared in the perspective schema associated to some *origin* schemata, will unify with the condition of the rule. The notation in (5) will be used through the paper to explain examples of rules.

The rule **RR-CA1** rewrites an ECA of equivalence, which connects a class/relation  $C^D$  of the *destination* schema to a class/relation  $C^{RM}$  of the reference model schema, into an ECA, which connect the

class/relation  $C^D$  to a class/relation  $C^S$  of some *origin* schema; when is provided an ECA that connect the class/relation  $C^{RM}$  to  $C^S$ .

An example of a RR-MF is as follows:

$$\mathbf{RR-MF1:} \frac{\mathbf{match} : \left( \left( \mathbf{RM} \begin{bmatrix} \underline{\mathbf{C}^{RM}} \end{bmatrix}, \underline{\tau^{RM}} \right) \times \left( \mathbf{D} \begin{bmatrix} \underline{\mathbf{C}^{D}} \end{bmatrix}, \underline{\tau^{D}} \right) \right) \to \text{Boolean} \Rightarrow}{\mathbf{RR-MF1:}} \frac{\mathbf{match} : \left( \left( \underline{\mathbf{S}} \begin{bmatrix} \underline{\mathbf{C}^{S}} \end{bmatrix}, \underline{\tau^{S}} \right) \times \left( \mathbf{D} \begin{bmatrix} \underline{\mathbf{C}^{D}} \end{bmatrix}, \underline{\tau^{D}} \right) \right) \to \text{Boolean}}{\mathbf{P}_{\underline{S}|RM} \begin{bmatrix} \underline{\mathbf{C}^{RM}} \end{bmatrix} \to \underline{\mathbf{S}} \begin{bmatrix} \underline{\mathbf{C}^{S}} \end{bmatrix}}.$$
(6)

In (6)  $\tau$  is a data type. The rule **RR-MF1** rewrites a match function signature, which matches a class/relation  $C^{RM}$  of the reference model schema to a class/relation  $C^D$  of the *destination* schema, in a match function signature that matches a class/relation  $C^S$  of some *origin* schema to the class/relation  $C^D$ , when is provided an ECA of equivalence that connects the class/relation  $C^{RM}$  to  $C^S$ .

An example of RR-C is as follows:

**RR-C1:** 
$$\frac{\operatorname{RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{\boldsymbol{p}}^{RM} \Rightarrow \underline{\mathbf{A}}^{S}}{\operatorname{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{\boldsymbol{p}}^{RM} \to \underline{\mathbf{A}}^{S}}.$$
 (7)

In (7)  $p^{RM}$  is a variable that will be instantiated with a property of a class/relation and **A** is a variable that will be instantiated with a component pattern expression. Similar to **K**<sup>S</sup> in (5), the letter **S** in **A**<sup>S</sup> means that all elements into that expression belong to schema **S**. The value of **S** and **A** will depend on which CA declared in the perspective schema associated to some *origin* schemata will unify with the condition of the rule.

The rule **RR-C1** rewrites a property  $p^{RM}$  of a class/relation of the reference model schema in a property, a path expression, or a function of some *origin* schema, when is provided an PCA that connects the property  $p^{RM}$  to this property, path expression, or function. The whole set of proposed rules can be found in appendix A.

1: procedure INFER\_CAS( $A^{\mathbf{G}} \to A^{\mathbf{RM}}, CAs$ ) 2: repeat 3: find  $A^{\mathbf{G}} \to A^{\mathbf{Si}}$  applying the inference rule R: 4:  $R: \frac{A^{\mathbf{G}} \to A^{\mathbf{Si}} \to A^{\mathbf{Si}}}{\text{conditions}};$ 5: add  $A^{\mathbf{G}} \to A^{\mathbf{Si}}$  to CAs; 6: until all rules for rewriting CAs have been tested 7: end procedure

Figure 6: The pseudo-code to the inference mechanism to generate new CAs.

A pseudo-code detailing as new CAs are deduced is shown in Fig. 6. In Fig. 6, **G** is the *destination* schema, **RM** a reference model schema, and  $S_i$ ,  $i \ge 1$ , *origin* schemata. The algorithm tries to find, for each CA  $A^{\mathbf{G}} \rightarrow A^{\mathbf{RM}}$  assigning the global schema to the reference model schema, one or more CAs  $A^{\mathbf{G}} \rightarrow A^{\mathbf{Si}}$  as a result of applying to  $A^{\mathbf{G}} \rightarrow A^{\mathbf{RM}}$  some rule for rewriting CAs. Notice that, in the condition of the rule can exists expressions of the form  $A \Rightarrow B$ . In this case, the recursivity will be present. For instance, a new ECA:

```
P_{S1|G} [CUSTOMER] \rightarrow S1 [CUSTOMER]
```

can be created based on  $\psi_3$  (see Fig. 4), using the rule **RR-CA1** since that the CA  $\psi_6$  is defined in perspective schema  $\mathbf{P}_{S1|RM}$  (as shown in Fig. 7).

Exte	nsion Correspondence Assertion (ECA)
$\psi_6$ :	$\boldsymbol{P_{S1 RM}}\left[\text{CUSTOMER}\right] \to \boldsymbol{S1}\left[\text{CUSTOMER}\right]$
$\psi_7$ :	$P_{S2 RM}$ [CUSTOMER] $\rightarrow$ S2 [CUSTOMER]

Figure 7: More examples of correspondence assertions.

A pseudo-code detailing as new MF signatures are deduced is shown in Fig. 8. In Fig. 8 K and L are pairs (classes/relations, data type) of the reference model schema or of the *destination* schema, while K' and L' are pairs (classes/relations, data type) of some *origin* schemata or of the *destination* one. For each MF M that is declared in the perspective schema associated to the *destination* schema, the algorithm tries to find one or more MFs as a result of applying to M some rule for rewriting MFs. For instance, two new MF signatures:

**match**((**S1**[CUSTOMER], $\tau_1$ )×(**G**[CUSTOMER], $\tau_2$ ))→Boolean **match**((**S2**[CUSTOMER], $\tau_1$ )×(**G**[CUSTOMER], $\tau_2$ ))→Boolean

can be created based on MF signature presented in Fig. 3, using the rule **RR-MF1** twice, since as the CAs  $\psi_6$  and  $\psi_7$  are defined, respectively, in perspective schemata  $\mathbf{P}_{S1|RM}$  and  $\mathbf{P}_{S2|RM}$  (as shown in Fig. 7).

```
    procedure INFER_MFS(match(K×L)→Boolean ,MFs)
    repeat
    find match(K'×L')→Boolean applying the inference rule R:
    R:match(K×L)→Boolean⇒match(K'×L')→Boolean;
conditions
    add match(K'×L')→Boolean to MFs;
    until all rules for rewriting MFs have been tested
    end procedure
```

Figure 8: The pseudo-code to the inference mechanism to generate new MFs.

A pseudo-code with the iteration of the process to generate a new perspective is shown in Fig. 9. In Fig. 9  $P_T$  is a perspective schema from the reference model to the global schema;  $P_j$ ,  $1 \neq j \neq n$ , are perspective schemata from source schemata to the reference model; and  $P_I$  is the inferred perspective schema from source schemata to the global schema. All elements of the perspective schemata are grouped in lists: classList, relationList, keyList, caList, and mfList. The three first lists hold 'require' declarations of, respectively, classes, relations, and keys and foreign keys. caList contains correspondence assertion declarations, and mfList has match function signatures.

This mechanism has been developed as part of a proof-of-concept prototype using a Prolog language. Beside the inference mechanism module, the prototype consist of more five modules, such as the *schema* manager, and the *ISCO translator*. The *schema manager* module is employed by the designer to manage the schemata (in language  $L_S$ ) as well as the perspective schemata (in language  $L_{PS}$ ). The *ISCO translator* performs the mapping between schemata written in  $L_S$  or  $L_{PS}$  languages to schemata defined in a language programming called Information Systems COnstruction language (ISCO) [9]. ISCO is based

```
1: procedure GENERATENEWPERSPECTIVE(P_T, P_1, ..., P_n, P_I)
        for each CA \ A^{\mathbf{G}} \to A^{\mathbf{RM}} in P_T.caList do
2:
            \texttt{infer_CAs}(A^{\mathbf{G}} \to A^{\mathbf{RM}}, \{A^{\mathbf{G}} \to A^{\mathbf{Si}}\});
3:
             add CAs \mathbb{A}^{\mathbf{G}} \to \mathbb{A}^{\mathbf{Si}} to P_I.caList;
4:
5:
        end for
        for each MF m \text{ in } P_T.mfList do
6:
7:
             infer_MFs(m, \{m'_i\});
             add MFs m'_i to P_I.mfList
8:
9:
        end for
10:
        for each E in classList/relationList/keyList do
             create a require declaration to P_I;
11:
12:
             add it, appropriately, to P<sub>I</sub>.classList/
             P_I.relationList/P_I.keyList
13:
14:
        end for
15: end procedure
```

Figure 9: The pseudo-code to the creation of inferred perspective schemata.

on a contextual constraint logic programming that allows the construction of information systems. It can define (object) relational schemata, represent data, and transparently access data from various heterogeneous sources in a uniform way, like a mediator system [10]. Thus, it is possible to access data from information sources using the perspective schema in ISCO. Furthermore, once the perspective schema from source schemata to the global schema has been inferred, as well as the new match functions have been implemented, it can be translated to ISCO language and so the data of the global schema can be queried.

### 5. Related work

The database community has been for many years engaged with the problem of data integration. Researches on this area have developed in several important directions: schema matching, data quality, to cite a few (see [11] for a survey), which can cover different architectures (e.g. FDBSs and DWs), representation of data and involved data models (e.g. relational and non-structured). Recent research in Federated Database Systems (FDBSs) has included: behaviour integration [12], integration of non-traditional data (e.g biomedical [13, 14], intelligence data [15], and web source [16]), interactive integration of data [17, 18], and federated data warehouse systems [19]. All these approaches use a global schema, but do not deal with a reference model schema. Similarly the authors' research of the current paper, [16] uses correspondence assertions (in this case, for specifying the semantics of XML-based mediators). However, their CAs only deal with part of the semantic correspondence managed here. Furthermore, they assume that there is a universal key to determine when two distinct objects are the same entity in the real-world, which is a supposition often unreal.

Researches in Data Warehouses (DWs) have focused on technical aspects such as multidimensional data models (e.g. [20, 21, 22, 23, 24, 25]) as well as the materialised view definition and maintenance (e.g. [26]). In particular, the most conceptual multidimensional models are extensions to the Entity-Relationship model (e.g. [27, 28, 29, 30]) or extensions to UML (e.g. [31, 32, 33]).

Reference in [34] focuses on an ontology-based approach to determine the mapping between at-

tributes from the source schemata and the DW schema, as well as to identify the transformations required for correctly moving data from source information to the DW. Their ontology, based on a common vocabulary as well as a set of data annotations (both provided by the designer), allows formal and explicit description of the semantic of the sources and the DW schemata. However, their strategy requires a deep knowledge of all schemata involved in the DW system, in what is usually not an usual task. In the proposed research of the present paper, it is dispensable, since each schema (source or DW) needs to be related only to the reference model one. Additionally, in [34] there is nothing about the matching of instances.

The approach closest to authors' research is described in [35]. Similar to this study, their proposal includes a reference model (cited as "enterprise model") designed using an Enriched Entity-Relationship (EER) model. However, unlike the authors' research, all their schemata, including the DW schema, are formed by relational structures, which are defined as views over the reference model. Their proposal provides the user with various levels of abstraction: conceptual, logical, and physical. In their conceptual level, they introduce the notion of intermodel assertions that precisely capture the structure of an EER schema or allow for the specifying of the relationship between diverse schemata. However, any transformation (e.g. restructuring of schema and values) or mapping of instances is deferred for the logical level, unlike the current work. In addition, they did not deal with complex data, integrity constraints, and path expressions, as this research does.

### 6. Conclusions and future works

In this paper, the authors have presented a proposal to automatically connect a global schema to its sources by using an inference mechanism taking into account a reference model. In proposed approach the relationship between the global schema and the source schemata is made explicitly and declaratively through correspondence assertions. This approach is particularly useful in data integration systems that define a common or canonical schema, such as in Data Warehouse (DW) systems and in Federated Database (FDB) systems. An advantage of the proposed approach is that by using the reference model the designer user does not need to have a deep knowledge of all schemata involved in the DW system or in the federation system, since that each schema (source or global) needs to be related only to the reference model one. Thus, the designer user can describe the global system without concerns about where the sources are or how they are stored. Besides, the process of data integration can be incrementally done in two sense:

- View schemata can be created as an intermediary process to relate portions of data that have been integrated (those view schemata, in turn, are related to the reference model). Thus the data integration process can be divided in small parts, instead of being seen as a whole, turning the integration task easiest.
- 2. New source schemata can be added or actual source schemata can eventually change. It is completely transparent to the DW systems or FDB systems since the relationship between the global schema and its source schemata is automatically created by the inference mechanism.

A prototype Prolog-based has been developed to allow the description of schemata and perspective schemata in the proposed language as well as to infer new perspective schemata based on other ones. The matching functions can be implemented using Prolog itself or external functions. In addition, the prototype include translators from the proposed language to the ISCO one. ISCO [9] allows access to

heterogeneous data sources and to perform arbitrary computations. Thus, user-queries can be done, in a transparent way, to access the information sources, like occurs in mediator systems [10].

For future work, investigations will be made into how the perspective schemata can be used to automate the materialisation of the data in the DWs or in other repository of a data integration environment. Another important direction for future work is the development of a graphical user-friend interface to declare the schemata in the proposed language, and thus, to hide some syntax details.

#### References

- [1] Claudia Imhoff, Nicholas Galemmo, and Jonathan G. Geiger, *Mastering Data Warehouse Design Relational and Dimensional Techniques*, Wiley Publishing, 2003.
- [2] Valéria Magalhães Pequeno and João Carlos Gomes Moura Pires, "Using perspective schemata to model the ETL process", in *ICMIS 2009 :Intl. Conf. on Management Information Systems*, France, June 2009, to appear.
- [3] Edgar F. Codd, "A relational model of data for large shared data banks", in *Communications of the ACM*, 1970, pp. 377–387.
- [4] Rick G.G. Cattell and D. Barry, Eds., *The Object Database Standard ODMG 3.0*, Morgan Kaufmann Publishers, 2000.
- [5] Valéria Magalhães Pequeno and João Carlos Gomes Moura Pires, "A formal object-relational data warehouse model", Tech. Rep., Universidade Nova de Lisboa, November 2007.
- [6] Valéria Magalhães Pequeno and João Carlos Gomes Moura Pires, "Using perspective schemata to model the ETL process", Tech. Rep., Universidade Nova de Lisboa, 2009.
- [7] G. Zhou, Richard Hull, and Roger King, "Generating data integration mediators that use materialization", J. Intell. Inf. Syst., vol. 6(2/3), pp. 199–221, May 1996.
- [8] Ramez Elmasri and Shamkant B. Navathe, *Fundamentals of database systems*, Pearson Education, 5th edition, 2006.
- [9] Salvador Abreu and Vitor Nogueira, "Using a logic programming language with persistence and contexts", in *INAP'05: 16th Intl. Conf. on applications of declarative programming and knowledge management.* 2006, vol. 4369 of *Lecture Notes in Computer Science*, pp. 38–47, Springer, (Revised Selected Papers).
- [10] G. Wiederhold, "Mediators in the architecture of future information systems", in *IEEE Computer*, 1992, vol. 25(3), pp. 38–49.
- [11] Alon Y. Halevy, Anand Rajaraman, and Joann J. Ordille, "Data integration: The teenage years.", in *VLDB*, 2006, pp. 9–16.
- [12] Markus Stumptner, Michael Schrefl, and Georg Grossmann, "On the road to behavior-based integration", in *APCCM: First Asia-Pacific Conf. on Conceptual Modelling*, 2004, pp. 15–22.
- [13] Brenton Louie, Peter Mork, Fernando Martin-Sanchez, Alon Halevy, and Peter Tarczy-Hornoch, "Data integration and genomic medicine", *Journal of Biomedical Informatics*, vol. 40, pp. 5–13, 2007.

- [14] Pavithra G. Naidu, Mathew J. Palakal, and Shielly Hartanto, "On-the-fly data integration models for biological databases", in SAC'07: Proceedings of the 2007 ACM symposium on Applied computing, USA, 2007, pp. 118–122, ACM.
- [15] S. Yoakum-Stover and T. Malyuta, "Unified architecture for integrating intelligence data", in *DAMA: Europe Conf.*, UK, 2008.
- [16] Vânia Maria Ponte Vidal, Bernadette Farias Lóscio, and Ana Carolina Salgado, "Using correspondence assertions for specifying the semantics of XML-based mediators", in *Workshop on Information Integration on the Web*, 2001, vol. 3(11).
- [17] Zachary G. Ives, Craig A. Knoblock, Steven Minton, Marie Jacob, Partha Pratim Talukdar, Rattapoom Tuchinda, Jos Luis Ambite, Maria Muslea, and Cenk Gazen, "Interactive data integration through smart copy & paste.", in *CIDR:4th Biennial Conference on Innovative Data Systems Research*. 2009, www.crdrdb.org.
- [18] Robert Mccann, Anhai Doan, Vanitha Varadarajan, and Er Kramnik, "Building data integration systems via mass collaboration", in *WebDB: Intl. Workshop on the Web and Databases*, USA, 2003.
- [19] Stefan Berger and Michael Schrefl, "From federated databases to a federated data warehouse system", in *HICSS '08: 41st Annual Hawaii Intl. Conf. on System Sciences*, USA, 2008, p. 394, IEEE Computer Society.
- [20] Dov Dori, Roman Feldman, and Arnon Sturm, "From conceptual models to schemata: An object-process-based data warehouse construction method", *Inf. Syst.*, vol. 33, no. 6, pp. 567–593, 2008.
- [21] E. Malinowski and E. Zimányi, "A conceptual model for temporal data warehouses and its transformation to the ER and the object-relational models", *Data knowl. eng.*, vol. 64, no. 1, pp. 101–133, 2008.
- [22] Juan Manuel Pérez, Rafael Berlanga, María José Aramburu, and Torben Bach Pedersen, "A relevance-extended multi-dimensional model for a data warehouse contextualized with documents", in *DOLAP'05: Proc. of the 8th ACM Intl. Workshop on Data Warehousing and OLAP*, USA, 2005, pp. 19–28, ACM.
- [23] Matteo Golfarelli, Vittorio Maniezzo, and Stefano Rizzi, "Materialization of fragmented views in multidimensional databases", *Data Knowl. Eng.*, vol. 49, no. 3, pp. 325–351, 2004.
- [24] Bodo Husemann, Jens Lechtenborger, and Gottfried Vossen, "Conceptual data warehouse modeling", in *Design and Management of Data Warehouses*, 2000, p. 6.
- [25] S. Rizzi, "Conceptual modeling solutions for the data warehouse", In Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications, vol. Information Science Reference, pp. 208–227, 2008.
- [26] Robert Wrembel, "On a formal model of an object-oriented database with views supporting data materialisation", in *Proc. of the Conf. on Advances in Databases and Information Systems*, 1999, pp. 109–116.
- [27] Enrico Franconi and Anand Kamble, "A data warehouse conceptual data model", *Proc. of the Int. Conf. on Scientific and Statistical Database Management*, vol. 00, pp. 435–436, 2004.
- [28] Anand S. Kamble, "A conceptual model for multidimensional data", in APCCM'08: Proc. of the 15th on Asia-Pacific Conf. on Conceptual Modelling, Australia, 2008, pp. 29–38, Australian Computer Society, Inc.

- [29] Carsten Sapia, Markus Blaschka, Gabriele Höfling, and Barbara Dinter, "Extending the E/R model for the multidimensional paradigm", in *Proc. of the Workshops on Data Warehousing and Data Mining*, 1999, pp. 105–116.
- [30] Nectaria Tryfona, Frank Busborg, and Jens G. Borch Christiansen, "starER: a conceptual model for data warehouse design", in *DOLAP '99: Proc. of the 2nd ACM Intl. Workshop on Data warehousing and OLAP*, USA, 1999, pp. 3–8, ACM.
- [31] Sergio Luján-Mora, Juan Trujillo, and Il-Yeol Song, "A UML profile for multidimensional modelling in data warehouses", *Data Knowl. Eng.*, vol. 59, no. 3, pp. 725–769, 2005.
- [32] T. B. Nguyen, A. Min Tjoa, and Roland Wagner, "An object oriented multidimensional data model for OLAP", in *Web-Age Inf. Management*, 2000, pp. 69–82.
- [33] Juan Trujillo, Manuel Palomar, and Jaime Gomez, "Applying object-oriented conceptual modeling techniques to the design of multidimensional databases and OLAP applications", WAIM'00. Lecture Notes in Computer Science (LNCS), vol. 1846, pp. 83–94, 2000.
- [34] Dimitrios Skoutas and Alkis Simitsis, "Designing ETL processes using semantic web technologies", in *DOLAP'06: Proceedings of the 9th ACM international workshop on Data warehousing and OLAP*, USA, 2006, pp. 67–74, ACM.
- [35] Diego Calvanese, Luigi Dragone, Daniele Nardi, Riccardo Rosati, and Stefano M. Trisolini, "Enterprise modeling and data warehousing in TELECOM ITALIA", *Inf. Syst.*, vol. 31, no. 1, pp. 1–32, 2006.

### APPENDIX

## A. Inference Rules

Hereafter, condider the following notation:

- $\mathcal{L}$  is a set of schema names.
- $\mathcal{L}_p$  is a set of perspective schema names.
- $\hat{\mathcal{L}}_p$  is a set of perspective schemata.
- $\mathcal{W}$  is a set of typed values.
- $\mathcal{T}$  is a set of data types.
- **D** is the *destination* schema, **RM** is the reference model schema, and **S** is a variable that can be instantiated with any of the *origin* schemata.
- C<sup>D</sup> is a variable that can be instantiated with any class/relation of the schema **D**; mutatis mutandis to C<sup>RM</sup> and C<sup>S</sup>.
- P<sub>RM|D</sub> ∈ L<sub>p</sub> is a perspective schema name, with **RM** being the *base* schema and **D** the *target* schema, mutatis mutandis to P<sub>S|D</sub> and P<sub>S|RM</sub>.
- All variables are indicated by an underline.

# A.1. Substituition-rules

The following notation will be used in this Section and through the paper to explain rules:

- A is a variable that can be instantiated with a component pattern expression, being that the letter S in A<sup>S</sup> means that all elements into that expression belong to schema S, mutatis mutandis to A<sup>RM</sup>.
- $p^{RM}$  is a variable that can be instantiated with any property of a class/relation of the schema **RM**, mutatis mutandis to  $p^S$ .
- $\rho$  is a variable that can be instantiated with a (value or reference) path expression as defined in L<sub>S</sub>, being that the letter S in  $\rho^S$  means that all elements into this expression belong to schema **S**, mutatis mutandis to  $\rho^{RM}$ .
- \$\ell\_i\$ are variables that can be instantiated with a links of a path expression, being that the letter S in \$\ell\_i^S\$ means that this element belong to schema \$\mathbf{S}\$, mutatis mutandis to \$\ell\_i^{RM}\$.
- pred is a predicate as defined in L<sub>PS</sub>, being op operands in pred such that op ∈ {<</li>
  ,>,≤,≥, \$\$,=} and B is an expression in pred such that B = A or B = w, with w ∈ W. The letter S in pred<sup>S</sup> and B<sup>S</sup> means that all elements into those expressions belong to schema S, mutatis mutandis to pred<sup>RM</sup> and B<sup>RM</sup>.
- $\varphi$  is a variable that can be instatiated with any function with  $n \ge 1$  arguments that returns a value.
- *FK* is a foreign key name.
- The expression p:  $\natural$ C means that the type of the property p is a reference to the class C.
- *type*() is a function defined in language L<sub>S</sub> that returns the structural type of a relation/class.

The substituition-rules are formed by 12 rules as follows:

$$\mathbf{RR-C1}: \frac{\mathrm{RM}\left[\underline{\mathbf{C}^{RM}}\right] \bullet \underline{p}^{RM} \Rightarrow \underline{\mathbf{A}}^{S}}{\mathrm{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}^{RM}}\right] \bullet \underline{p}^{RM} \to \underline{\mathbf{A}}^{S}}$$

$$\mathbf{RR-C2}: \frac{\mathrm{RM}\left[\underline{\mathbf{C}^{RM}}\right] \bullet \underline{\varrho^{RM}} \Rightarrow \underline{\mathbf{S}}\left[\underline{\mathbf{C}^{S}}\right] \bullet \underline{\varrho^{S}}}{\underline{\ell_{i+1}^{RM}} \Rightarrow \ell_{i+1}^{S}, \text{ for } 0 \le i \le n-1,} \\ \overline{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}_{n}^{RM}}\right] \bullet \underline{p}^{RM}} \to \underline{\mathbf{S}}\left[\underline{\mathbf{C}_{n}^{S}}\right] \bullet \underline{p}^{S}}$$

$$\mathbf{RR-C3}: \frac{\mathrm{RM}\left[\underline{C^{RM}}\right] \bullet \underline{\varrho^{RM}} \Rightarrow \underline{S}\left[\underline{C^{s}}\right] \bullet \underline{\varrho^{S}}}{\underline{\ell_{i+1}^{RM}} \Rightarrow \underline{\ell_{i+1}^{S}}, \text{ for } 0 \le i \le n-1}$$

$$\mathbf{RR-C4}: \frac{\underline{\varphi}\left(\underline{\mathbf{A}_{1}^{RM}}, \underline{\mathbf{A}_{2}^{RM}}, \dots, \underline{\mathbf{A}_{n}^{RM}}\right) \Rightarrow \underline{\varphi}\left(\underline{\mathbf{A}_{1}^{S}}, \underline{\mathbf{A}_{2}^{S}}, \dots, \underline{\mathbf{A}_{n}^{S}}\right)}{\underline{\mathbf{A}_{i}^{RM}} \Rightarrow \underline{\mathbf{A}_{i}^{S}}, \text{for } 1 \leq i \leq n}$$

$$\mathbf{RR-C5}: \frac{\mathrm{RM}\left[\underline{\mathbf{C}^{RM}}\right] \bullet \underline{\boldsymbol{p}^{RM}}\{\underline{\boldsymbol{p}_i}\} \Rightarrow \underline{\mathbf{A}_i^S}, \text{ for } 1 \leq i \leq n}{\mathrm{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}^{RM}}\right] \bullet \underline{\boldsymbol{p}^{RM}}\{\boldsymbol{p}_1, \boldsymbol{p}_2, ..., \boldsymbol{p}_n\} \rightarrow \left(\underline{\mathbf{A}_1^S}, \underline{\mathbf{A}_2^S}, \dots, \underline{\mathbf{A}_n^S}\right)}$$

**RR-C6** : 
$$\underline{w} \Rightarrow \underline{w}$$

$$\mathbf{RR-C7}: \frac{\underline{\mathbf{A}^{RM}}}{\underline{\mathbf{A}^{RM}}} \xrightarrow{\mathbf{op}} \underline{\mathbf{B}^{RM}} \Rightarrow \underline{\mathbf{A}^{S}} \xrightarrow{\mathbf{op}} \underline{\mathbf{B}^{S}}$$
$$\frac{\underline{\mathbf{A}^{RM}}}{\underline{\mathbf{B}^{RM}}} \Rightarrow \underline{\mathbf{A}^{S}},$$
$$\underline{\mathbf{B}^{RM}} \Rightarrow \underline{\mathbf{B}^{S}}$$

$$\mathbf{RR}\text{-}\mathbf{C8}: \frac{\underline{\mathbf{A}^{RM}}}{\underline{\mathbf{p}}} \quad \underline{\mathbf{p}} \quad \underline{\mathbf{B}^{RM}} \text{ and } \underline{\mathbf{pred}^{RM}} \Rightarrow \underline{\mathbf{A}^{S}} \quad \underline{\mathbf{op}} \quad \underline{\mathbf{B}^{S}} \text{ and } \underline{\mathbf{pred}^{S}}$$
$$\frac{\underline{\mathbf{A}^{RM}}}{\underline{\mathbf{B}^{RM}}} \Rightarrow \underline{\mathbf{A}^{S}},$$
$$\frac{\underline{\mathbf{B}^{RM}}}{\underline{\mathbf{pred}^{RM}}} \Rightarrow \underline{\mathbf{pred}^{S}}$$

$$\mathbf{RR-C9}: \underbrace{\frac{\mathbf{A}^{RM}}{\mathbf{p}} \quad \underline{\mathbf{op}} \quad \underline{\mathbf{B}}^{RM} \text{ or } \underline{\mathbf{pred}}^{RM} \Rightarrow \underline{\mathbf{A}}^{S} \quad \underline{\mathbf{op}} \quad \underline{\mathbf{B}}^{S} \text{ or } \underline{\mathbf{pred}}^{S}}_{\mathbf{pred}^{RM}} \Rightarrow \underline{\mathbf{A}}^{S}, \\ \underbrace{\frac{\mathbf{A}^{RM}}{\mathbf{p}} \Rightarrow \underline{\mathbf{B}}^{S}}_{\mathbf{pred}^{RM}} \Rightarrow \underline{\mathbf{pred}}^{S}}_{\mathbf{pred}^{RM}}$$

$$\mathbf{RR-C10}: \frac{\underline{p^{RM}}: \underline{\mathbf{C}^{RM}} \to \underline{\mathbf{C}_1^{RM}} \Rightarrow \underline{p^S}: \underline{\mathbf{C}^S} \to \underline{\mathbf{C}_1^S}}{\underline{\mathbf{P}_{\underline{S}|RM}} \left[ \underline{\mathbf{C}^{RM}} \right] \bullet \underline{p^{RM}} \to \underline{\mathbf{S}} \left[ \underline{\mathbf{C}^S} \right] \bullet \underline{p^S}}, \\ \underline{p^S}: \natural \underline{\mathbf{C}_1^S} \in type(\mathbf{C}^S)$$

$$\mathbf{RR-C11}: \frac{\ell^{RM}}{\mathbf{P}_{\underline{S}|RM}} \stackrel{\longrightarrow}{\to} \underbrace{\mathbf{C}_{1}^{RM}}_{\mathbf{C}} \Rightarrow \underbrace{\mathbf{FK}^{S}}_{\mathbf{K}}: \underbrace{\mathbf{C}^{S}}_{\mathbf{C}} \rightarrow \underbrace{\mathbf{C}_{1}^{S}}_{\mathbf{C}}$$
$$\underbrace{\mathbf{P}_{\underline{S}|RM}}_{(\mathbf{FK}^{S}, \mathbf{C}^{S}, \neg, \mathbf{C}_{1}^{S}, \neg)} \Rightarrow \underbrace{\mathbf{S}}_{\mathbf{C}} \underbrace{\left[\underline{\mathbf{C}^{S}}\right]}_{\mathbf{C}},$$

$$\mathbf{RR-C12}: \frac{\underline{F\underline{K}^{RM}}: \underline{\mathbf{C}^{RM}} \to \underline{\mathbf{C}_{1}^{RM}} \Rightarrow \underline{\underline{p}^{S}}: \underline{\mathbf{C}^{S}} \to \underline{\mathbf{C}_{1}^{S}}}{\mathbf{P}_{\underline{S}|RM} \left[\underline{\mathbf{C}^{RM}}\right] \to \underline{\mathbf{S}} \left[\underline{\mathbf{C}^{S}}\right],}$$
$$\mathbf{P}_{\underline{S}|RM} \left[\underline{\mathbf{C}_{1}^{RM}}\right] \to \underline{\mathbf{S}} \left[\underline{\mathbf{C}_{1}^{S}}\right],$$
$$\underline{\underline{p}^{S}}: \natural \underline{\mathbf{C}_{1}^{S}} \in type(\underline{\mathbf{C}^{S}})$$

### A.2. rewritten-rules to rewrite CAs

The rules to rewrite CAs are subdivided in four groups in accordance to kind of CA involved. Thus, there are rules for rewriting PCAs, ECAs, SCAs and ACAs, which are presented in following text.

#### A.2.1. rewritten-rules to rewrite PCAs

Consider the following notation to describe the RR-PCAs:

•  $\mathbf{G}^S$  is a variable that can be instantied with the right side of a CA pattern expression of property consisting of one of two forms:  $(\mathbf{A}_1^S, \mathbf{A}_2^S, \dots, \mathbf{A}_n^S)$  or  $(\mathbf{B}_1^S, \mathbf{pred}_1^S)$ ,  $(\mathbf{B}_2^S, \mathbf{pred}_2^S)$ ,  $\dots$ ,  $(\mathbf{B}_{n-1}^S, \mathbf{pred}_{n-1}^S)$ ,  $\mathbf{B}_n^S$ .

The RR-PCAs are formed by five rules as follows:

$$\mathbf{RR}\text{-}\mathbf{PCA1}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{A}}^{RM} \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{A}}^{S}}{\underline{\mathbf{A}}^{RM} \Rightarrow \underline{\mathbf{A}}^{S}}$$

$$\mathbf{RR}\text{-}\mathbf{PCA2}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{p}^{RM} \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{G}}^{\underline{S}}}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{p}^{RM} \to \underline{\mathbf{G}}^{\underline{S}}}$$

$$\mathbf{RR}\text{-}\mathbf{PCA3}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D}\{\underline{p}_{1}, \underline{p}_{2}, ..., \underline{p}_{n}\} \to \left(\underline{\mathbf{A}}_{1}^{RM}, \underline{\mathbf{A}}_{2}^{RM}, \dots, \underline{\mathbf{A}}_{n}^{RM}\right) \Rightarrow}{\underline{\mathbf{P}}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D}\{\underline{p}_{1}, \underline{p}_{2}, ..., \underline{p}_{n}\} \to \left(\underline{\mathbf{A}}_{1}^{S}, \underline{\mathbf{A}}_{2}^{S}, \dots, \underline{\mathbf{A}}_{n}^{S}\right)}$$

$$\frac{\underline{\mathbf{A}}_{i}^{RM} \Rightarrow \underline{\mathbf{A}}_{i}^{S}, \text{ for } 1 \leq i \leq n}{\underline{\mathbf{A}}_{i}^{RM} = \underline{\mathbf{A}}_{i}^{S}}$$

$$\mathbf{RR} \cdot \mathbf{PCA4}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \left(\underline{\mathbf{A}}_{1}^{RM}, \underline{\mathbf{A}}_{2}^{RM}, \dots, \underline{\mathbf{A}}_{n}^{RM}\right) \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \left(\underline{\mathbf{A}}_{1}^{S}, \underline{\mathbf{A}}_{2}^{S}, \dots, \underline{\mathbf{A}}_{n}^{S}\right)}{\underline{\mathbf{A}}_{i}^{RM} \Rightarrow \underline{\mathbf{A}}_{i}^{S}, \text{ for } 1 \leq i \leq n}$$

$$\mathbf{RR}\text{-}\mathbf{PCA5}: \frac{\mathbf{P}_{RM|D} \Big[ \underline{\mathbf{C}}^{D} \Big] \bullet \underline{p}^{D} \to \Big( \underline{\mathbf{B}}_{1}^{RM}, \underline{\mathbf{pred}}_{1}^{RM} \Big) ; \Big( \underline{\mathbf{B}}_{2}^{RM}, \underline{\mathbf{pred}}_{2}^{RM} \Big) ; ...; \Big( \underline{\mathbf{B}}_{n-1}^{RM}, \underline{\mathbf{pred}}_{n-1}^{RM} \Big) ; \underline{\mathbf{B}}_{n}^{RM} \Rightarrow \\ \frac{\mathbf{P}_{\underline{S}|D} \left[ \underline{\mathbf{C}}^{D} \right] \bullet \underline{p}^{D} \to \Big( \underline{\mathbf{B}}_{1}^{S}, \underline{\mathbf{pred}}_{1}^{S} \Big) ; \Big( \underline{\mathbf{B}}_{2}^{S}, \underline{\mathbf{pred}}_{2}^{S} \Big) ; ...; \Big( \underline{\mathbf{B}}_{n-1}^{S}, \underline{\mathbf{pred}}_{n-1}^{S} \Big) ; \underline{\mathbf{B}}_{n}^{S}} \\ \frac{\mathbf{B}_{i}^{RM} \Rightarrow \underline{\mathbf{B}}_{i}^{S}, \underline{\mathbf{pred}}_{1}^{RM} \Rightarrow \underline{\mathbf{pred}}_{i}^{RM} \Rightarrow \underline{\mathbf{pred}}_{i}^{S}, \text{ for } 1 \leq i \leq n \end{aligned}$$

# A.2.2. rewritten-rules to rewrite ECAs

Consider the following notation to describe the RR-ECAs:

- $\mathbf{K}$  is a variable that can be instantiated with the right side of a CA pattern expression of extension, being that the letter S in  $\mathbf{K}^S$  means that all elements in that expression belong to schema **S**.
- $\diamond$  is any operand appearing in a ECA, i.e.  $-, \cap,$  or  $\square$ .
- $C_i$ , for  $1 \le i \le n$ , are class/relation names in some schema belonging to  $\hat{\mathcal{L}}$ .

The RR-ECAs are formed by four rules as follows:

$$\mathbf{RR}\text{-}\mathbf{ECA1}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \to \underline{\mathbf{K}}^{S}}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \to \underline{\mathbf{K}}^{S}}$$

$$\begin{aligned} \mathbf{RR}\text{-}\mathbf{ECA2}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\left(\underline{\mathbf{pred}}^{RM}\right)\right] \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \to \underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\left(\underline{\mathbf{pred}}^{S}\right)\right] \\ & \\ \mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \to \underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\right], \\ & \\ & \\ \underline{\mathbf{pred}}^{RM} \Rightarrow \underline{\mathbf{pred}}^{S} \end{aligned}$$

$$\mathbf{RR}\text{-}\mathbf{ECA3}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\left(\underline{\mathbf{pred}}^{RM}\right)\right] \Rightarrow}{\mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \to \underline{\mathbf{S}}\left[\underline{\mathbf{C}}_{1}^{S}\left(\underline{\mathbf{pred}}_{1}^{S}\right)\right] \diamond \underline{\mathbf{S}}\left[\underline{\mathbf{C}}_{2}^{S}\left(\underline{\mathbf{pred}}_{2}^{S}\right)\right] \diamond \ldots \diamond \underline{\mathbf{S}}\left[\underline{\mathbf{C}}_{n}^{S}\left(\underline{\mathbf{pred}}_{n}^{S}\right)\right]}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \to \underline{\mathbf{S}}\left[\underline{\mathbf{C}}_{1}^{S}\right] \diamond \underline{\mathbf{S}}\left[\underline{\mathbf{C}}_{2}^{S}\right] \diamond \ldots \diamond \underline{\mathbf{S}}\left[\underline{\mathbf{C}}_{n}^{S}\left(\underline{\mathbf{pred}}_{n}^{S}\right)\right]}{\mathbf{pred}^{RM}} \Rightarrow \underline{\mathbf{pred}}_{i}^{S}, \text{ for } 1 \leq i \leq n$$

$$\mathbf{RR}\text{-}\mathbf{ECA4}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}^{D}}\right] \to \mathrm{RM}\left[\underline{\mathbf{C}_{1}^{RM}}\right] \diamond \ldots \diamond \mathrm{RM}\left[\underline{\mathbf{C}_{j}^{RM}}\left(\underline{\mathbf{pred}_{j}^{RM}}\right)\right] \diamond \ldots \diamond \mathrm{RM}\left[\underline{\mathbf{C}_{n}^{RM}}\right] \Rightarrow}{\mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}^{D}}\right] \to \underline{\mathbf{K}_{1}^{S}} \diamond \underline{\mathbf{K}_{2}^{S}} \diamond \ldots \diamond \underline{\mathbf{C}_{j}^{S}}(\underline{\mathbf{pred}_{j}^{S}}) \diamond \ldots \diamond \underline{\mathbf{K}_{n}^{S}}}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}_{i}^{RM}}\right] \to \underline{\mathbf{K}_{i}^{S}}, \text{ for } 1 \leq i \leq j-1, j+1 \leq i \leq n,} \\ \mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}_{i}^{RM}}\right] \to \underline{\mathbf{K}_{i}^{S}}, \mathbf{pred}_{j}^{RM} \Rightarrow \underline{\mathbf{pred}_{j}^{S}}$$

# A.2.3. rewritten-rules to rewrite SCAs

Consider the following notation to describe the RR-SCAs:

x is a variable that can be instantiated with a property p or a path expression ρ or a property p' into another property p (a structured type) (notation: p{p'}).

- **Q** is a variable that can be instantiated with the right side of a CA pattern expression of summation, being that the letter S in **Q**<sup>S</sup> means that all elements in that expression belong to schema **S**.
- $\vartheta$  is a variable that can be instantiated with the keywords *groupby* or *normalize*, the two possible kinds of SCA.
- *p<sub>i</sub>*, for 1 ≤ i ≤ n, are property names belonging to classes/relations in some schema belonging to L̂.

The RR-SCAs are formed by four rules as follows:

$$\mathbf{RR}\text{-}\mathbf{SCA1}: \frac{\mathsf{P}_{RM|D}\left[\underline{\mathsf{C}}^{D}\right]\left(\underline{p}_{1}^{D}, \dots, \underline{p}_{n}^{D}\right) \to \underline{\vartheta}(\mathsf{RM}\left[\underline{\mathsf{C}}^{RM}\right](\underline{\mathbf{X}}_{1}^{RM}, \dots, \underline{\varphi}(\underline{\mathbf{A}}_{1}^{RM}, \dots, \underline{\mathbf{A}}_{n}^{RM}), \dots, \underline{\mathbf{X}}_{n}^{RM})) \Rightarrow}{\mathsf{P}_{\underline{S}|D}\left[\underline{\mathsf{C}}^{D}\right]\left(\underline{p}_{1}^{D}, \dots, \underline{p}_{n}^{D}\right) \to \underline{\vartheta}(\underline{\mathsf{S}}\left[\underline{\mathsf{C}}^{S}\right](\underline{\mathbf{X}}_{1}^{S}, \dots, \underline{\varphi}(\underline{\mathbf{A}}_{1}^{S}, \dots, \underline{\mathbf{A}}_{n}^{S}), \dots, \underline{\mathbf{X}}_{n}^{S}))} \\ \xrightarrow{\mathsf{P}_{\underline{S}|RM}\left[\underline{\mathsf{C}}^{RM}\right] \to \underline{\vartheta}(\underline{\mathsf{S}}\left[\underline{\mathsf{C}}^{S}\right](\underline{\mathbf{X}}_{1}^{S}, \dots, \underline{\varphi}(\underline{\mathbf{A}}_{1}^{S}, \dots, \underline{\mathbf{X}}_{n}^{S}))}{\mathsf{RM}\left[\underline{\mathsf{C}}^{RM}\right] \to \underline{\mathsf{S}}\left[\underline{\mathsf{C}}^{S}\right], \\ \operatorname{\mathsf{RM}}\left[\underline{\mathsf{C}}^{RM}\right] \bullet \underline{\mathbf{X}}_{i}^{RM} \Rightarrow \underline{\mathsf{S}}\left[\underline{\mathsf{C}}^{S}\right] \bullet \underline{\mathbf{X}}_{i}^{S}, \text{ for } 1 \leq i \leq n, \\ \underline{\varphi}(\underline{\mathsf{A}}_{1}^{RM}, \dots, \underline{\mathsf{A}}_{n}^{RM}) \Rightarrow \underline{\varphi}(\underline{\mathsf{A}}_{1}^{S}, \dots, \underline{\mathsf{A}}_{n}^{S})}{\mathsf{P}_{RM|D}\left[\underline{\mathsf{C}}^{D}\right]\left(\underline{p}_{1}^{D}, \dots, \underline{p}_{n}^{D}\right) \to \underline{\vartheta}(\mathsf{RM}\left[\underline{\mathsf{C}}^{RM}\left(\underline{\mathsf{pred}}^{RM}\right)\right](\underline{\mathsf{X}}_{1}^{RM}, \dots, \underline{\varphi}(\underline{\mathsf{A}}_{1}^{RM}, \dots, \underline{\mathsf{A}}_{n}^{RM}), \dots, \mathbf{X}_{n}^{RM})}$$

$$\mathbf{RR}\operatorname{-}\mathbf{SCA2}: \xrightarrow{\mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{\underline{D}}\right]\left(\underline{p}_{1}^{D},\ldots,\underline{p}_{n}^{D}\right)} \xrightarrow{\rightarrow \underline{\vartheta}(\underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\left(\underline{\mathsf{pred}}^{S}\right)\right](\underline{\mathbf{x}}_{1}^{S},\ldots,\underline{\varphi}(\underline{\mathbf{A}}_{1}^{S},\ldots,\underline{\mathbf{A}}_{n}^{S}),\ldots,\underline{\mathbf{x}}_{n}^{S}))}{\operatorname{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \rightarrow \underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\right],} \\ \xrightarrow{\mathbf{pred}^{RM} \Rightarrow \mathbf{pred}^{S},} \\ \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{\mathbf{x}}_{i}^{RM} \Rightarrow \underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\right] \bullet \underline{\mathbf{x}}_{i}^{S}, \text{ for } 1 \leq i \leq n,} \\ \xrightarrow{\underline{\varphi}(\underline{\mathbf{A}}_{1}^{RM},\ldots,\underline{\mathbf{A}}_{n}^{RM}) \Rightarrow \underline{\varphi}(\underline{\mathbf{A}}_{1}^{S},\ldots,\underline{\mathbf{A}}_{n}^{S})}$$

$$\mathbf{RR}\text{-}\mathbf{SCA3}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right]\left(\underline{p}_{1}^{D}, \dots, \underline{p}_{n}^{S}\right) \to \underline{\mathbf{Q}}^{S}}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right]\left(\underline{p}_{1}^{RM}, \dots, \underline{p}_{n}^{RM}\right) \to \underline{\mathbf{Q}}^{S}}, \\ \mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}_{i}^{D} \to \mathbf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{p}_{i}^{RM}, \text{ for } 1 \leq i \leq n$$

$$\mathbf{RR}\text{-}\mathbf{SCA4}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \to \mathrm{RM}\left[\underline{\mathbf{C}}^{RM}\left(\underline{\mathbf{pred}}^{RM}\right)\right] \Rightarrow}{\mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right]\left(\underline{\mathbf{p}}_{1}^{D}, \dots, \underline{\mathbf{p}}_{n}^{D}\right) \to \underline{\vartheta}\left(\underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\left(\underline{\mathbf{pred}}^{S}\right)\right]\left(\underline{\mathbf{x}}_{1}^{S}, \dots, \underline{\varphi}\left(\underline{\mathbf{A}}_{1}^{S}, \dots, \underline{\mathbf{A}}_{n}^{S}\right), \dots, \underline{\mathbf{x}}_{n}^{S}\right)\right)}{\mathbf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right]\left(\underline{\mathbf{p}}_{1}^{RM}, \dots, \underline{\mathbf{p}}_{n}^{RM}\right) \to \underline{\vartheta}(\underline{\mathbf{S}}\left[\underline{\mathbf{C}}^{S}\right](\underline{\mathbf{x}}_{1}^{S}, \dots, \underline{\varphi}(\underline{\mathbf{A}}_{1}^{S}, \dots, \underline{\mathbf{A}}_{n}^{S}), \dots, \underline{\mathbf{x}}_{n}^{S}))}{\mathbf{pred}^{RM} \Rightarrow \mathbf{pred}^{S}}, \\ \frac{\mathbf{p}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{\mathbf{p}}_{\underline{i}}^{D} \to \mathrm{RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{\mathbf{p}}_{\underline{i}}^{RM}, \text{ for } 1 \leq i \leq n$$

## A.2.4. rewritten-rules to rewrite ACAs

Consider the notation used to define the rules RR-PCAs. Also consider the following notation:

•  $\gamma$  is a variable that can be instantiated with one of the aggregation functions (sum, count, min, max, avg) used in SCAs.

• **sca** is a variable that can be instantiated with the name of the respective SCA asigned to an ACA.

The RR-ACAs are formed by eight rules as follows (the first six rules are to rewrite ACAs related to SCA of normalisation, while the last two rules are to rewrite ACAs related to SCA of group by):

$$\mathbf{RR}\text{-}\mathbf{ACA1}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}^{D}}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\mathbf{A}^{RM}} \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}^{D}}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\mathbf{A}^{S}}}{\underline{\mathbf{A}^{RM}} \Rightarrow \underline{\mathbf{A}^{S}}}$$

$$\mathbf{RR}\text{-}\mathbf{ACA2}: \frac{\mathsf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \mathsf{RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{p}^{RM} \Rightarrow \mathsf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\mathbf{G}}^{S}}{\mathsf{P}_{\underline{S}|RM}\left[\underline{\mathbf{C}}^{RM}\right] \bullet \underline{p}^{RM}} \to \underline{\mathbf{G}}^{S}$$

$$\mathbf{RR}\text{-}\mathbf{ACA3}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D}\{\underline{p}_{1}, \underline{p}_{2}, \dots, \underline{p}_{n}\} \to \underline{\mathbf{sca}}, \left(\underline{\mathbf{A}_{1}^{RM}}, \underline{\mathbf{A}_{2}^{RM}}, \dots, \underline{\mathbf{A}_{n}^{RM}}\right) \Rightarrow}{\underline{\mathbf{P}_{S|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D}\{\underline{p}_{1}, \underline{p}_{2}, \dots, \underline{p}_{n}\} \to \underline{\mathbf{sca}}, \left(\underline{\mathbf{A}_{1}^{S}}, \underline{\mathbf{A}_{2}^{S}}, \dots, \underline{\mathbf{A}_{n}^{S}}\right)}{\underline{\mathbf{A}_{i}^{RM}} \Rightarrow \underline{\mathbf{A}_{i}^{S}}, \text{ for } 1 \leq i \leq n}$$

$$\mathbf{RR}\text{-}\mathbf{ACA4}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \left(\underline{\mathbf{A}}_{1}^{RM}, \underline{\mathbf{A}}_{2}^{RM}, \dots, \underline{\mathbf{A}}_{n}^{RM}\right) \Rightarrow}{\underline{\mathbf{P}}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \left(\underline{\mathbf{A}}_{1}^{S}, \underline{\mathbf{A}}_{2}^{S}, \dots, \underline{\mathbf{A}}_{n}^{S}\right)}{\underline{\mathbf{A}}_{i}^{RM} \Rightarrow \underline{\mathbf{A}}_{i}^{S}, \text{ for } 1 \leq i \leq n}$$

$$\mathbf{RR}\text{-}\mathbf{ACA5}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \left(\underline{\mathbf{B}}_{1}^{RM}, \underline{\mathbf{pred}}_{1}^{RM}\right); \dots; \left(\underline{\mathbf{B}}_{n-1}^{RM}, \underline{\mathbf{pred}}_{n-1}^{RM}\right); \underline{\mathbf{B}}_{n}^{RM} \Rightarrow}{\underline{\mathbf{P}}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \left(\underline{\mathbf{B}}_{1}^{S}, \underline{\mathbf{pred}}_{1}^{S}\right); \dots; \left(\underline{\mathbf{B}}_{n-1}^{S}, \underline{\mathbf{pred}}_{n-1}^{S}\right); \underline{\mathbf{B}}_{n}^{S}} \xrightarrow{\underline{\mathbf{B}}_{n}^{RM}} \\ \frac{\underline{\mathbf{B}}_{i}^{RM} \Rightarrow \underline{\mathbf{B}}_{i}^{S}}{\underline{\mathbf{pred}}_{i}^{RM}} \Rightarrow \underline{\mathbf{pred}}_{i}^{S}, \text{ for } 1 \leq i \leq n$$

$$\mathbf{RR}\text{-}\mathbf{ACA6}: \frac{\mathsf{P}_{RM|D}\left[\underline{\mathsf{C}}^{D}\right] \bullet \underline{p}^{D} \to \mathsf{RM}\left[\underline{\mathsf{C}}^{RM}\right] \bullet \underline{p}^{RM} \Rightarrow \mathsf{P}_{\underline{S}|D}\left[\underline{\mathsf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathsf{sca}}, \underline{\mathsf{G}}^{S}}{\mathsf{P}_{\underline{S}|RM}\left[\underline{\mathsf{C}}^{RM}\right] \bullet \underline{p}^{RM}} \to \underline{\mathsf{sca}}, \underline{\mathsf{G}}^{S}}$$

$$\mathbf{RR}\text{-}\mathbf{ACA7}: \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}^{D}}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\gamma}\left(\underline{\mathbf{A}^{RM}}\right) \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}^{D}}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\gamma}\left(\underline{\mathbf{A}^{S}}\right)}{\underline{\mathbf{A}^{RM}} \Rightarrow \underline{\mathbf{A}^{S}}}$$

$$\mathbf{RR} \cdot \mathbf{ACA8} : \frac{\mathbf{P}_{RM|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\gamma}\left(\underline{\mathbf{A}}^{RM}, \underline{\mathbf{pred}}^{RM}\right) \Rightarrow \mathbf{P}_{\underline{S}|D}\left[\underline{\mathbf{C}}^{D}\right] \bullet \underline{p}^{D} \to \underline{\mathbf{sca}}, \underline{\gamma}\left(\underline{\mathbf{A}}^{S}, \underline{\mathbf{pred}}^{S}\right)}{\underline{\mathbf{A}}^{RM} \Rightarrow \underline{\mathbf{A}}^{S}}, \underline{\mathbf{pred}}^{RM} \Rightarrow \underline{\mathbf{pred}}^{S}}$$

#### A.3. rewritten-rules to rewrite matching function signatures

Consider the following notation to describe the RR-MFs:

- $\mathbf{Y}$  is a variable that can be instantiated with a component pattern expression of forms:  $\underline{\mathbf{S}}[\underline{\mathbf{C}}^S] \bullet \underline{p}^S$  or  $\underline{\varphi}(\underline{\mathbf{S}}[\underline{\mathbf{C}}^S] \bullet \underline{p}^S)$ ; being that the letter S in  $\mathbf{Y}^S$  means that all elements into that expression belong to schema S.
- **S**<sub>1</sub> and **S**<sub>2</sub> are variables that can be instantiated with distinct schema names belonging to  $\mathcal{L}$ .
- C<sup>D</sup><sub>i</sub>, for 1 ≤ i ≤ n, are variable that can be instantiated with any class/relation of the schema D; mutatis mutandis to C<sup>RM</sup><sub>i</sub> and C<sup>S</sup><sub>i</sub>.
- $p_i^{RM}$ , for  $1 \le i \le n$ , are variables that can be instantiated with any property of a class/relation of the schema **RM**, mutatis mutandis to  $p_i^S$  and  $p_i^D$ .
- $\tau$  and  $\tau_i$  are variable that can be instantiated with any data type belonging to  $\mathcal{T}$ , being that the letter S in  $\tau^S$  (and  $\tau_i^S$ ) means that the type was described in the schema **S**, mutatis mutandis to  $\tau^{RM}$ ,  $\tau_i^{RM}$ ,  $\tau^D$  and  $\tau_i^D$ .
- **E** is a variable that can be instantiated with a class/relation with a selection condition (notation: C(**pred**)) or without a selection condition (notation: C), being that the letter D in **E**<sup>D</sup> means that all elements into that expression belong to schema **D**.

The RR-MFs are formed by seven rules defined as follows:

$$\mathbf{RR}\mathbf{-MF1}: \frac{\mathbf{match}: \left( \left( \mathrm{RM}\left[ \underline{\mathbf{C}^{RM}} \right], \underline{\tau^{RM}} \right) \times \left( \mathrm{D}\left[ \underline{\mathbf{E}^{D}} \right], \underline{\tau^{D}} \right) \right) \to \mathrm{Boolean} \Rightarrow}{\mathrm{P}_{\underline{S}|RM}\left[ \underline{\mathbf{C}^{RM}} \right] \to \underline{\mathbf{S}}\left[ \underline{\mathbf{C}^{S}} \right], \\ \underline{\tau^{S}} = type(\underline{\mathbf{C}^{S}})$$

$$\begin{split} \mathbf{RR}\text{-}\mathbf{MF2}: & \frac{\mathbf{match}:\left(\left(\mathrm{RM}\left[\underline{\mathbf{C}_{1}^{RM}}\right],\underline{\tau_{1}^{RM}}\right)\times\left(\mathrm{RM}\left[\underline{\mathbf{C}_{2}^{RM}}\right],\underline{\tau_{2}^{RM}}\right)\right) \rightarrow \mathrm{Boolean} \Rightarrow \\ \mathbf{RR}\text{-}\mathbf{MF2}: & \frac{\mathbf{match}:\left(\left(\underline{\mathbf{S}_{1}}\left[\underline{\mathbf{C}_{1}^{S1}}\right],\underline{\tau_{1}^{S1}}\right)\times\left(\underline{\mathbf{S}_{2}}[\underline{\mathbf{C}_{2}^{S2}}],\underline{\tau_{2}^{S2}}\right)\right) \rightarrow \mathrm{Boolean}}{\mathbf{P}_{\underline{S1}|RM}\left[\underline{\mathbf{C}_{1}^{RM}}\right] \rightarrow \underline{\mathbf{S}_{1}}\left[\underline{\mathbf{C}_{1}^{S1}}\right], \\ & \mathbf{P}_{\underline{S2}|RM}\left[\underline{\mathbf{C}_{2}^{RM}}\right] \rightarrow \underline{\mathbf{S}_{2}}\left[\underline{\mathbf{C}_{2}^{S2}}\right], \\ & \frac{\tau_{1}^{S1}}{\mathbf{1}} = \mathbf{type}(\underline{\mathbf{C}_{1}^{S1}}), \\ & \frac{\tau_{2}^{S2}}{\mathbf{1}} = \mathbf{type}(\underline{\mathbf{C}_{2}^{S2}}) \end{split}$$

**RR-MF3**: 
$$\frac{\text{match}: \left( \left( D\left[\underline{\mathbf{E}}_{\underline{1}}^{D}\right], \underline{\tau}_{\underline{1}}^{D}\right) \times \left( D\left[\underline{\mathbf{E}}_{\underline{2}}^{D}\right], \underline{\tau}_{\underline{2}}^{D}\right) \right) \rightarrow \text{Boolean} \Rightarrow}{\text{match}: \left( \left( D\left[\underline{\mathbf{E}}_{\underline{1}}^{D}\right], \underline{\tau}_{\underline{1}}^{D}\right) \times \left( D\left[\underline{\mathbf{E}}_{\underline{2}}^{D}\right], \underline{\tau}_{\underline{2}}^{D}\right) \right) \rightarrow \text{Boolean}}$$