

# Extensional Knowledge for Semantic Query Optimization in a Mediator Based System <sup>\*</sup>

D. Beneventano<sup>1,2</sup>, S. Bergamaschi<sup>1,2</sup>, and F. Mandreoli<sup>1</sup>

<sup>1</sup> Dipartimento di Scienze dell'Ingegneria - Università di Modena e Reggio Emilia  
DSI - Via Vignolese 905, 41100 Modena

{domenico.beneventano,sonia.bergamaschi,mandreoli.federica}@unimo.it

<sup>2</sup> CSITE-CNR Bologna V.le Risorgimento 2, 40136 Bologna

**Abstract.** Query processing in global information systems integrating multiple heterogeneous sources is a challenging issue in relation to the effective extraction of *information* available on-line. In this paper we propose intelligent, tool-supported techniques for querying global information systems integrating both structured and semistructured data sources. The techniques have been developed in the environment of a data integration, wrapper/mediator based system, MOMIS, and try to achieve the goal of *optimized query reformulation* w.r.t local sources. The developed techniques rely on the availability of *integration knowledge* whose semantics is expressed in terms of description logics. Integration knowledge includes local source schemata, a virtual mediated schema and its *mapping descriptions*, that is semantic mappings w.r.t. the underlying sources both at the *intensional* and *extensional* level. Mapping descriptions, obtained as a result of the semi-automatic integration process of multiple heterogeneous sources developed for the MOMIS system, include, unlike previous data integration proposals, *extensional intra/interschema knowledge*. Extensional knowledge is exploited to perform semantic query optimization in a mediator based system as it allows to devise an *optimized query reformulation* method. The techniques are under development in the MOMIS system but can be applied, in general, to data integration systems including extensional intra/interschema knowledge in *mapping descriptions*.

## 1 Introduction

The purpose of data integration is to provide a uniform interface to multiple heterogeneous sources. Applications range from *searching information on the net* to *providing an uniform consistent view* of data associated with the various legacy systems of an enterprise. Query processing in such global information systems environment is a challenging issue which has been faced in many previous works in both the AI and Database Community [2, 8, 14, 17]. A data integration system, based on conventional wrapper/mediator architectures, allows the user to pose a query and receive a unified answer without the need of: *locating* the sources

---

<sup>\*</sup> This work was partially supported by MIUR co-funded project D2I.

relevant to the query, *interacting* with each source in isolation and *combining* the data coming from the different sources. Data integration systems usually follow this architecture: each data source provides a schema and a mediated (global) *virtual* schema of all the sources is obtained manually or semi-automatically, for a particular integration application. The mediated schema has a set of *mapping descriptions* (called source descriptions in [14]) that specify the semantic mapping between the mediated schema and the source schemas. Unlike previous mentioned approaches, in the MOMIS system [3, 4], *mapping descriptions* obtained as a result of the semi-automatic integration process include *extensional intra/interschema knowledge* which represents fundamental knowledge for a correct and complete schema integration [20]. The mediator uses mapping descriptions to reformulate a user query into queries over the source schemata. From a theoretical point of view, solving a user (mediated) query, i.e. giving a single unified answer w.r.t. multiple sources, implies to face two main problems: *query reformulation/optimization* [11, 12, 14, 15, 18] and *object fusion* [16, 21].

In this paper we deal with the query reformulation/optimization problem and propose a theoretical framework providing intelligent, tool-supported techniques to query global information systems integrating both structured and semistructured data sources. The framework exploits all the available integration knowledge, and, in particular, extensional inter/intra schema knowledge. Therefore, we do not address the object fusion problem, which concerns the grouping of information (from the same or different sources) about the same real-world entity. Anyway, we believe that it is a relevant topic now and will become more important in the future as integration systems cope with more and more information that has not been nicely structured and partitioned in advance.

The proposed techniques are under development in the MOMIS environment where they rely on the availability of integration knowledge, whose semantics is given in terms of description logics. Integration knowledge is expressed in terms of: local source schemata, a virtual mediated schema and its *mapping descriptions*, *extensional intra/interschema knowledge*.

Extensional knowledge is exploited to perform semantic query optimization in a mediator based system as it allows to devise an *optimized query reformulation* method. In particular, starting from the method developed in [20], we exploit the “base extension” approach and the description logics inference techniques in order to face the reformulation/optimization problem of a mediated query.

The outline of the paper is the following. Section 2 summarizes the MOMIS approach to schema integration and introduces the running example. Section 3 introduces mapping descriptions and base extensions. Finally, Section 4 presents the MOMIS Query Manager implementing the proposed theoretical framework by means of examples.

## 2 Preliminaries

Like other integration projects [1, 19], MOMIS follows a “semantic approach” to information integration based on the conceptual schema, or metadata, of the

information sources, and on the  $I^3$  architecture [13]; for a detailed description of the MOMIS system see [8, 3] available at <http://www.dbgroup.unimo.it/Momis>.

**The ODL<sub>I<sup>3</sup></sub> language** The system supports an object-oriented language, called ODL<sub>I<sup>3</sup></sub>, which is very close to the ODL language [13, 10] and allows a semantically rich representation of source schemas and object patterns associated with information sources to be integrated. ODL<sub>I<sup>3</sup></sub> is a source independent language used for information extraction to describe heterogeneous schemas of structured and semistructured data sources in a common way. ODL<sub>I<sup>3</sup></sub> introduces three main extensions with respect to ODL [13]: intentional relationships, that are *terminological relationships* expressing intra and inter-schema knowledge for the source schemas, extensional relationships and integrity constraint rules. The intentional relationships of synonymy (SYN), hypernymy (BT) and positive association (NT) between two classes  $C_1$  and  $C_2$  may be “strengthened” by establishing that they are also *extensional* relationships [9]. Consequently, the following extensional relationships can be defined in ODL<sub>I<sup>3</sup></sub>:

- $C_1$  SYN<sub>ext</sub>  $C_2$ : the instances of  $C_1$  are the same of  $C_2$ .
- $C_1$  BT<sub>ext</sub>  $C_2$ : the instances of  $C_1$  are a superset of the instances of  $C_2$ .
- $C_1$  NT<sub>ext</sub>  $C_2$ : the instances of  $C_1$  are a subset of the instances of  $C_2$ .
- $C_1$  DISJ<sub>ext</sub>  $C_2$ : the instances of  $C_1$  are disjoint from the instances of  $C_2$ .

In contrast with [20] we do not introduce an *overlap relationship* as we assume a default overlap relationships among two classes if no extensional relationship is specified. Moreover, extensional relationships “constrain” the structure of the two classes  $C_1$  and  $C_2$ . If an extensional relationship  $C_1$  NT<sub>ext</sub>  $C_2$  is issued, we have that:

- strict inheritance between  $C_1$  and  $C_2$  is enforced for the common attributes;
- both  $C_1$  and  $C_2$  may have further attributes as we adopt usual description logics semantics (i.e. open world semantics).

Extensional relationships can be partially automatically extracted and partially explicitly declared by the integration designer.

The integrity constraint rules are introduced in ODL<sub>I<sup>3</sup></sub> in order to express, in a declarative way, *if then* integrity constraint rules at both intra- and inter-sources level. The antecedent and the consequent of an integrity constraint rule are expressions formulated by the *intersection* operator (**and**), the *union* operator (**or**) and the *negation* operator (**not**).

ODL<sub>I<sup>3</sup></sub> descriptions are translated into OLC<sub>D</sub> (*Object Language with Completions allowing Descriptive cycles*) descriptions [6] in order to perform Description Logics inferences that will be useful for semantic integration. Indeed, OLC<sub>D</sub> supports some interesting reasoning techniques: *subsumption* and *equivalence* detection between types (i.e. “is-a” relationships implied by type descriptions), and *inconsistency* type detection.

In particular, MOMIS translates into OLC<sub>D</sub> ODL<sub>I<sup>3</sup></sub> class descriptions, integrity constraint rules and extensional relationships, giving rise to a schema

called *inter-sources schema*, formally defined as follows. Let  $\mathbf{L}$  be a set of *local class names* (denoted by  $L_1, L_2, \dots$ ) and let  $\mathbf{LA}$  be a set of *local attributes* (denoted by  $la_1, la_2, \dots$ ).  $L_A$  is a total function  $L_A: \mathbf{L} \rightarrow 2^{\mathbf{LA}}$  which associates local class names with attributes. The *inter-sources schema*  $\sigma: \mathbf{L} \rightarrow \mathbf{S}(\mathbf{LA}, \mathbf{L})$  associates local class names to their descriptions; an instance  $\mathcal{I}$  of  $\sigma$  defines the instances of the local classes: given a local class  $L$ ,  $\mathcal{I}(L)$  denotes the set of its real-world.

**Intelligent schema integration** The MOMIS approach to intelligent schema integration is supported by a tool, SI-Designer [3] and is articulated in the following phases:

1. *Generation of a Common Thesaurus* - The Common Thesaurus, consisting of intra and inter-schema relationships expressed in  $ODL_{I^3}$ , is generated.
2. *Affinity analysis of  $ODL_{I^3}$  classes* - Relationships in the Common Thesaurus are used to evaluate the level of intra and inter source *affinity* between classes. The concept of affinity is introduced to formalize the kind of relationships that can occur between classes from the integration point of view. The affinity of two classes is established by means of affinity coefficients based on class names, class structures and their relationships in the Common Thesaurus.
3. *Clustering  $ODL_{I^3}$  classes* - Classes with affinity in different sources are grouped together in clusters using hierarchical clustering techniques. The goal is to identify the classes that have to be integrated since describing the same or semantically related information.
4. *Generation of the mediated schema*- A *global class* is defined for each cluster, which is representative of all cluster's classes and is characterized by the union of their attributes. The *Global Schema* for the analyzed sources is composed of all the global classes derived from clusters, and is the basis for posing queries against the sources.

**The Query Manager** The user application interacts with MOMIS to query the Global Schema by using the  $OQL_{I^3}$  language. This phase is performed by the QM that generates the  $OQL_{I^3}$  queries for wrappers. Using Mapping Descriptions and DLs techniques, the QM generates in an automatic way the reformulation/optimization of the generic  $OQL_{I^3}$  query into different sub-queries, one for each involved local source.

To achieve the mediated query result, the QM has to assemble each local sub-query result into a unified data set. This process involves the solution of redundancy and reconciliation problems, due to the incomplete and overlapping information available on the local sources, i.e. *Object Fusion*.

As a mediator is not the *owner* of the data stored in the local classes but it only provides a virtual view, this means that the mediator has to *recognize* instances of the sources to be fused in an object. This recognition is a difficult task as each source may have its own techniques to identify objects, like keys

### UNIVERSITY source (*UNI*)

```
Research_Staff(name,e_mail,dept_code,s_code)
School_Member(name,school,year,e_mail)
Department(dept_name,dept_code,budget)
Section(section_name,s_code,length,room_code)
Room(room_code,seats_number,notes)
```

### COMPUTER\_SCIENCE source (*CS*)

```
CS_Person(first_name,last_name)
Professor:CS_Person(belongs_to:Division,rank)
Student:CS_Person(year,takes:set(Course),rank,e_mail)
Division(description,address:Location)
Location(city,street,number,country)
Course(course_name,taught_by:Professor)
```

### TAX\_POSITION\_XML source (*TP*)

```
<!ELEMENT ListOfStudent (Student*)>
<!ELEMENT Student (name, s_code, school_name, e_mail, tax_fee)>
<!ELEMENT name (#PCDATA)>
...
```

Fig. 1. Three heterogeneous University Sources

for relational or OIDs for object sources, and, usually instances referring to the same real word object are identified with different keys or OIDs, depending on the source the object is stored. The idea of our approach is to find *semantically homogeneous attributes* for each instance of each local class, on the basis of the available integration knowledge (see [5]).

**Running example** We consider three sources with different data model (see Figure 1). The first source is a relational database, **University** ( $S_1$ ), containing data about the staff and the students of a given university. The relations are: **Research\_Staff**, **School\_Member**, **Department**, **Section** and **Room**. For a given professor (in **Research\_Staff**) his department (**dept\_code**) and his section (**s\_code**) are stored. In the relation **School\_Member** the information **name**, **year**, **e\_mail**, and **school** about students enrolled at the university are stored.

The second source **Computer\_Science** ( $S_2$ ) is an object-oriented database containing information about people belonging to the Computer Science department of the same university, and is an object-oriented database. There are six classes: **CS\_Person**, **Professor**, **Student**, **Division**, **Location** and **Course**. Information is quite similar to the first source: it stores data on professors and students, also giving the possibility to retrieve the division of a given professor. This division may be part of another department, being a logical specialization of **Department**. The class **Location** maintains the division address. With respect to students, we may know the courses they take and their enrollment year. A third source is also available, **Tax\_Position** ( $S_3$ ), derived from the Registry Office. It consists of an XML file, storing information about student's tax.fees.

The Common Thesaurus may contain the following extensional relationships:

1. UNI.School\_Member SYN<sub>Ext</sub> TP.Student
2. CS.Student NT<sub>Ext</sub>UNI.School\_Member
3. CS.Professor NT<sub>Ext</sub> UNI.Research\_Staff
4. CS.Professor DISJ<sub>Ext</sub> UNI.School\_Member
5. UNI.Research\_Staff DISJ<sub>Ext</sub> TP.Student
6. UNI.Research\_Staff DISJ<sub>Ext</sub> CS.Student
7. CS.Student NT<sub>Ext</sub> CS.CS\_Person
8. CS.Professor NT<sub>Ext</sub>CS.CS\_Person

Relationships from 1 to 6 are designer-supplied inter-schema relationships; 7 and 8 are intra-schema extensional relationships automatically extracted by the MOMIS system from the isa relationships of the COMPUTER\_SCIENCE source.

An example of integrity constraint rule at intra-source level is the following:

```
rule Rule1 forall X in CS.Professor:  
  X.rank = 'full' then X.belongs_to.description='department';
```

At inter-sources level, both rules between local classes and global/local classes can be expressed. Examples of inter-sources integrity constraint rules are the following:

```
rule Rule2 forall X in UNI.School_Member:  
  X.school = 'cs' then X in TP.Student;
```

```
rule Rule3 forall X in CS.Student:  
  X in CS.Student then (X in TP.Student)  
  and not (X in UNI.Research_Staff);
```

### 3 Mapping Descriptions and Base Extensions

#### 3.1 Global Class and Mapping Tables

Let us assume that the first three phases of the integration have been performed on our University integration example<sup>1</sup>, thus obtaining the mediated global schema with five clusters and a global class for each cluster. For each global class a persistent *mapping-table* storing all the mappings is generated; it is a table whose rows represent the set of the local classes belonging to the cluster and whose columns represent the global attributes. An element  $MT[L][ga]$  represents the set of attributes of the local class  $L$  which are mapped into the global attribute  $ga$ : the value of the  $ga$  attribute is a function of the values assumed by the set of attributes  $MT[L][ga]$ . Some simple and frequent cases of such function are the following (see Figure 2 as an example):

- *identity* : the  $ga$  value is equal to the  $la$  value; we denote this case as  $MT[L][ga] = la$ .

<sup>1</sup> For a detailed description of the mappings selection and of the tool SI-Designer which assist the designer in this integration phase see [3].

University_Person	name	dept	e_mail	section	school
UNI_Research_Staff	name	dept_code	e_mail	s_code	null
UNI_School_Member	name	null	e_mail	null	school
CS_CS_Person	first_name and last_name	null	null	null	"cs"
CS_Student	first_name and last_name	null	e_mail	null	"cs"
CS_Professor	first_name and last_name	null	null	null	"cs"
TP_Student	name	null	e_mail	null	school_name

year	belong_to	takes	rank	s_code	tax_fee
null	null	null	"professor"	null	null
year	null	null	"student"	null	null
null	null	null	null	null	null
year	null	takes	rank	null	null
null	"belong_to"	null	rank	null	null
null	null	null	"student"	s_code	tax_fee

Fig. 2. Mapping Table of the Global Class University\_Person

- *concatenation* : the  $ga$  value is obtained as a concatenation of the values assumed by a set of local attributes  $la_i$  of the local class  $L$ ; we denote this case as  $MT[L][ga] = la_1 \text{ and } \dots \text{ and } la_n$  (see  $MT[CS.Student][name]$ ).

When the  $ga$  has no correspondence with any attribute of the local class  $L$ , the possible choices are:

- *constant* :  $ga$  assumes into the local class  $L$  a constant value set by the designer; we denote this case by  $MT[L][ga] = const$  (see the **rank** attribute).
- *undefined* :  $ga$  is set undefined into the local class  $L$ ; we denote this case by  $MT[L][ga] = null$ .

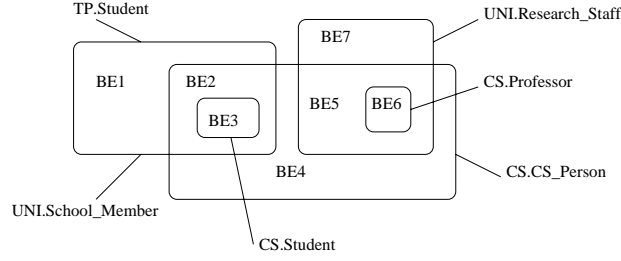
Formally, let  $\sigma$  be an inter-sources schema,  $\sigma: \mathbf{L} \rightarrow \mathbf{S}(\mathbf{LA}, \mathbf{L})$ . In order to define a global class  $G$ , we have to consider a subset of  $\mathbf{L}$ , which are the local classes of the cluster related to  $G$ ; this subset is denoted with  $\mathbf{L}_G$  and with  $\mathbf{LA}_G$  we denote the attributes of  $\mathbf{L}_G$ .

**Definition 1 (Global Class).** Given a set  $\mathbf{L}$  of local class names and a set  $\mathbf{GA}$  of global attributes (denoted by  $ga_1, ga_2, \dots$ ), a global class  $G$  is a tuple  $G = (\mathbf{L}_G, \mathbf{GA}, MT)$  where  $MT$ , called mapping table, is a total function  $MT: \mathbf{L}_G \times \mathbf{GA} \rightarrow 2^{\mathbf{LA}_G} \cup \{const\} \cup \{null\}$ .

### 3.2 Base Extensions

Intuitively, given a global class, a Base Extension gathers up all objects of some local classes such that the set of Base Extensions satisfies all the extensional relationships defined over the set of local classes and allows a partitioning of the set of the sources objects.

**Definition 2 (Base Extension set).** Let  $G = (\mathbf{L}_G, \mathbf{GA}, MT)$  be a global class and  $\mathcal{I}$  be an instance of the inter-sources schema  $\sigma: \mathbf{L} \rightarrow \mathbf{S}(\mathbf{LA}, \mathbf{L})$ . A set of base



**Fig. 3.** Base extension Set for the global class University\_Person

CL	BE	BE1	BE2	BE3	BE4	BE5	BE6	BE7
UNI.School_Member		1	1	1	0	0	0	0
UNI.Research_Staff		0	0	0	0	1	1	1
CS.CS_Person		0	1	1	1	1	1	0
CS.Student		0	0	1	0	0	0	0
CS.Professor		0	0	0	0	0	1	0
TP.Student		1	1	1	0	0	0	0

**Fig. 4.** Tabular representation of the Base Extension Set of University\_Person

extensions of  $G$  on  $\mathcal{I}$  is a pair  $(\mathbf{B}, F)$ , where  $\mathbf{B}$  is a set of base extension names (denoted by  $B_1, B_2, \dots$ ),  $\mathcal{I}(B_i) = \bigcap_{L \in F(B)} \mathcal{I}(L)$ , and  $F$  is a total function  $F : \mathbf{B} \rightarrow 2^{L_G}$  such that:  $\bigcup_{B \in \mathbf{B}} F(B) = L_G$ , and the set  $\left\{ \mathcal{I}(B) - \bigcup_{L \in (L_G - F(B))} \mathcal{I}(L) \mid B \in \mathbf{B} \right\}$  is a partition of  $\bigcup_{L \in L_G} \mathcal{I}(L)$ .

At present, we adopt the algorithm of [20] to determine a base extension set<sup>2</sup>. Figure 3 shows the Base Extension Set for University\_Person. A Base extension set of a Global Class  $G$  is represented by a table. Table rows represent the local classes of the global class and table columns represent the base extensions. The presence of a 1 in the table cell  $(L, B)$  means  $L \in F(B)$  (e.g. see Figure 4).

The attributes of a Base Extension  $B$  are the global attributes which are mapped, by a not null mapping, into a local class of  $B$ . Formally:

**Definition 3 (Base Extension Attributes).** Let  $G = (L_G, \mathbf{GA}, MT)$  be a global class and  $(\mathbf{B}, F)$  be the set of base extensions of  $G$ , then the attributes of a base extension  $B \in \mathbf{B}$  are defined as:

$$A(B) = \{ga \in \mathbf{GA} \mid \exists L \in F(B), MT[L][ga] \neq null\}.$$

The Base Extension Attributes of our example are shown in Figure 5.

<sup>2</sup> More than one base extension set can be obtained on the basis of the above definition; the discussion about the quality of the selected set is out of the scope of this paper.



$$\begin{aligned}
A(BE1) &= \{\text{name, year, school, rank, e\_mail, s\_code, tax\_fee}\} \\
A(BE2) &= \{\text{name, year, school, rank, e\_mail, s\_code, tax\_fee}\} \\
A(BE3) &= \{\text{name, year, school, rank, e\_mail, s\_code, tax\_fee, takes}\} \\
A(BE4) &= \{\text{name, school}\} \\
A(BE5) &= \{\text{name, rank, dept, e\_mail, section, school}\} \\
A(BE6) &= \{\text{name, rank, dept, e\_mail, section, belong\_to, school}\} \\
A(BE7) &= \{\text{name, rank, dept, e\_mail, section}\}
\end{aligned}$$

**Fig. 5.** Base Extension Attributes

**Definition 4 (Domination).** *Given a global class  $G = (\mathbf{L}_G, \mathbf{GA}, MT)$ , the set of base extensions  $(\mathbf{B}, F)$  of  $G$ , and  $B_1, B_2 \in \mathbf{B}$  we say that  $B_1$  dominates  $B_2$  w.r.t. the set of global attributes  $X \subseteq \mathbf{GA}$  iff  $X \subseteq A(B_1) \cap A(B_2) \wedge F(B_1) \subset F(B_2)$ .*

## 4 The MOMIS Query Manager

The MOMIS Query Manager functional architecture is shown in figure 6. As shown in figure, the optimized query reformulation is obtained on the basis of the computed Base Extension set and by using the semantic query optimization techniques previously developed by the authors [7]. In particular, base extensions allows a query to be rewritten with respect to the local sources, whereas, semantic query optimization techniques are used to transform the query on the basis of the intra- and inter-sources integrity constraint rules.

We will show the effectiveness of our optimization method by means of the following (mediated) query:

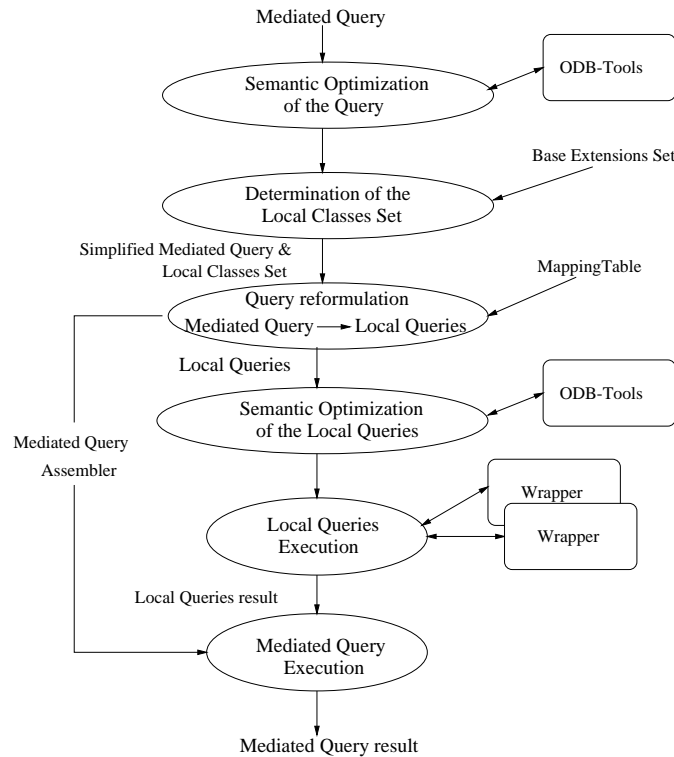
```

Q: select e_mail
   from University_Person
   where school = 'cs'
   and (s_code = 'a1x' or year = '2001' or tax_fee < 200)

```

Processing the above query, without considering extensional knowledge, would individuate all the local classes for which at least one of the mediated query attributes has a not null mapping with it. The candidate local classes are thus all the local classes of `University_Person`. The query has thus to be reformulated on the basis of the above classes.

By considering extensional knowledge, i.e. integrity constraint rules and extensional relationships, and its exploitation by means of base extensions, description logics, types and mapping table, the access to some local class can be avoided, obtaining an effective query optimization, independent from any specific cost model. Let us follow the processing phases of Figure 6 in order to show this result.



**Fig. 6.** Query processing phases

**1. Semantic optimization of the mediated query**

1. the mediated query is transformed by considering the integrity rules at both intra- and inter-sources level

**2. Determination of the Local Classes set**

1. the subset of Base Extensions including the global attributes of the query is computed
2. the local classes included into the Base Extension subset computed at point 1. are determined

**3. Query Reformulation w.r.t local sources**

1. query reformulation into local queries on the basis of the classes identified in the previous phase

**4. Semantic optimization of the local queries**

1. The local queries are transformed by considering the integrity constraint rules at both intra- and inter-sources level

**5. Local query execution**

1. local queries are sent to the wrappers to be translated and executed by the local sources

**6. Mediated query execution**

1. for each Base Extension object fusion of the local queries results is performed
2. the union of the results of the queries identified in the previous phase is performed

**1. Semantic optimization of the mediated query** This optimization can be performed if some integrity constraint rule at inter-sources level are available. In this case, the query can be transformed by considering the integrity rules. For example, let us consider the following query:

```
Q1: select name
      from University_Person X
      where X.rank = 'graduate'
      and X.tax_fee < 200
```

If the following integrity rules is available:

```
rule Rule4 forall X in University_Person:
  X.rank = 'graduate' then (X.tax_fee <100)
```

then the factor  $(X.\text{tax\_fee} < 200)$  can be eliminated from the query Q1. Notice that to perform semantic optimization of the mediated query we need to use the mapping table. In fact, the mediated query refers to a global class (global attributes) and intra-sources integrity constraint rules refers to local classes (local attributes).

**2. Determination of the Local Class Set** In order to determine the local classes for query reformulation we first add to the original query a non-null condition on the select clause attributes; for the query Q we have:

```
Q: select e_mail
      from University_Person
      where school = 'cs'
      and (s_code = 'a1x' or year = '2001' or tax_fee < 200)
      and e_mail is not null
```

then we consider the *Disjunctive Normal Form - DNF* of the query condition:  $DNF = F1 \text{ or } F2 \text{ or } F3$  where:

```
F1 = (school='cs') and (s_code='a1x') and (e_mail is not null)
F2 = (school='cs') and (tax_fee<200)and (e_mail is not null)
F3 = (school='cs') and (year='2001')and (e_mail is not null)
```

For each factor  $F$  of  $DNF$  we define the set:  $BE(F) = \{B \mid \forall ga \text{ of } F, ga \in A(B)\}$ , i.e.,  $B \in BE(F)$  iff  $A(B)$  contains all the global attributes of the factor  $F$ . In our example:

- $BE(F1) = BE(F3) = \{BE1, BE2, BE3\}$
- $BE(F2) = \emptyset$

Intuitively, a factor  $F$  of  $DNF$  such that  $BE(F) = \emptyset$  can be eliminated as the value of  $F$  is always false. In our example, we obtain a simplified  $DNF = F1$  or  $F3$ .

On the basis of this simplification, we obtain the following result: we do not have to access the local classes `UNI.Research_Staff` and `CS.Professor` as the only predicate related to their attributes (`tax_fee<200`) has been eliminated.

In the presence of more complex queries and a large set of extensional relationships the optimization results that can be obtained by using base extensions can be effective. Furthermore, in some cases, as:

```
Q: select email from University_Person
   where school = 'cs' and section = 'info1'
```

we have  $BE(F) = \emptyset$ , that is an empty answer (no access to the sources).

Local classes are determined by considering the union of all the local classes included in a subset of the previously evaluated base extensions obtained by eliminating the *dominated* ones. Indeed, the concept of domination as introduced in Def. 4 can be applied w.r.t. the global attributes of each factor  $F$  for which we compute  $BE(F)$ . With reference to our example, we have  $BE(F1) = BE(F3) = \{BE1\}$ , as  $BE1$  dominates both  $BE2$  and  $BE3$ ; as a consequence the identified local classes are:  $\{TP.Student, UNI.School_Member\}$ . Notice that the classes `CS.Student` and `CS.CS_Person` are excluded from query execution too: they are useless as their objects are also instances of other local classes (as, for instance, stated by `CS.Student NTExt UNI.School_Member`) and their contributions to the query, that is attributes `school` and `e_mail`, are already provided by the identified local classes.

As to summarize, the result of this phase are: a simplified Mediated Query, a set of base extensions and, then, a Set of Local Classes.

**3. Query Reformulation** Query reformulation is based on the base extensions selected in the previous phase. First, for each pair of local classes belonging to the same base extension the related join attributes are considered. In our example, we have only a base extension,  $BE1$ , then the local classes are `TP.Student` and `UNI.School_Member` and the join attribute is `name` for both the classes. Then, we consider the simplified  $DNF$  obtained in the previous phase and, for each factor  $F$ , we build a *local query* for each local class of  $BE(F)$ .

```
QL: select <select-list> from L
     where <condition>
```

in two steps:

1. `<condition>` is the conjunction of all predicates of the factor  $F$  which can be *solved* in the local class  $L$  (at least one predicate since  $L \in BE(F)$ ); these predicates are rewritten w.r.t. the attributes of local class  $L$  on the basis of the mapping table. In our example, for factor  $F1$  and the class `TP.Student` we have `<condition> = (school_name='cs')and(s_code='a1x')and(e_mail`

is not null) whereas for the class `UNI.School.Member` we have `<condition> = (school='cs') and (e_mail is not null)`. Note that, the same predicate may be mapped into semantic homogeneous attributes of more than one class. From a theoretical point of view, such a multiple mapping introduces no problem as semantic homogeneous attributes have been individuated in the integration activity; on the other hand, from an optimized execution point of view only some classes supporting such a predicate could be selected according to cost (either financial or computational) criteria.

2. `<select-list>` is obtained by adding to the query select clause all the join attributes.

In our example, we have two local queries:

```
QF1L1: select e_mail, name      QF1L2: select e_mail, name
        from TP.Student          from UNI.School_Member
        where (school_name='cs')  where (school = 'cs')
        and (s_code='aix')        and (e_mail is not null)
        and (e_mail is not null)
```

**4. Semantic optimization of the local queries** The local queries are transformed, first, by considering the inter-sources integrity constraint rules and, then by considering the intra-sources ones.

By considering the integrity constraint rules at inter-sources level, the subsumption relationships between local queries are computed: if a local query  $QL_1$  is subsumed by a local query  $QL_2$  and the *select* attributes of  $QL_1$  are a subset of the *select* attributes of the query  $QL_2$  thus it is not necessary to execute the local query  $QL_1$ . For example, by applying rule `Rule2`, we obtain that the local query `QF1L2` is subsumed by the local query `QF1L1`. In this case we avoid the `QF1L2` execution as its select clause attributes coincide with those of `QF1L1`.

For each source a local semantic query optimization is performed by considering, for each local query, the intra source integrity constraint rules. Note that this kind of optimization is independent from the physical organization of data of a source as our techniques are independent from any specific cost model.

**5. Local Query Execution** Local queries are sent to the wrappers to be translated and executed by the local sources.

**6. Mediated Query Execution** The first step of the Mediated Query Execution is the object fusion of the local query results belonging to the same base extension; this is implemented by a query, called *object fusion query*, which is performed for each factor and each base extension previously individuated.

In our example, for factor `F1` we have only the base extension `BE1` and the object fusion of `QF1L1` and `QF1L2` is based on the direct-join on the join attribute `name`, thus we obtain the following object fusion query:

```
QF1BE1: select e_mail
         from QF1L1,QF1L2
         where QF1L1.name=QF1L2.name
```

of course, all the object fusion queries have the same `<select-list>`.

The second and last step of the Mediated Query Execution is the union of the object fusion queries. In our example, the outcome is the union between QF1BE1 and QF3BE1, that is the analogue of QF1BE1 for factor F3 and *BE1*.

## 5 Conclusions

In this paper we proposed tool-supported techniques to query global information systems. The techniques are under development in the environment of a data integration, wrapper/mediator based system, MOMIS, and try to achieve the goal of: *optimized query reformulation* w.r.t local sources.

The techniques rely on the availability of *integration knowledge*, whose semantics is given in terms of description logics. Integration knowledge includes: local sources schemata, a virtual mediated schema and its *mapping descriptions*, that is semantic mappings w.r.t. the underlying sources both at the intensional and *extensional* level, *extensional intra/inter-schema knowledge*. Extensional knowledge is exploited to perform semantic query optimization in a mediator based system as it allows to devise an optimized query reformulation method. The method is based on *base extensions* and description logics inference techniques. In particular, base extensions permit to transform a query on the basis of the extensional relationships between classes, whereas, semantic query optimization techniques are used to transform the query on the basis of the intra- and inter-sources integrity constraint rules.

## References

1. Y. Arens, C.Y. Chee, C. Hsu, and C. A. Knoblock. Retrieving and integrating data from multiple information sources. *International Journal of Intelligent and Cooperative Information Systems*, 2(2):127–158, 1993.
2. Y. Arens, C. A. Knoblock, and W. Shem. Query reformulation for dynamic information integration. *Journal of Intelligent Information Systems (JIIS)*, 6(1):99–130, 1996.
3. I. Benetti, D. Beneventano, S. Bergamaschi, A. Corni, F. Guerra, and G. Malvezzi. Si-designer: a tool for intelligent integration of information. *Int. Conference on System Sciences (HICSS2001)*, 2001.
4. D. Beneventano, S. Bergamaschi, S. Castano, A. Corni, R. Guidetti, G. Malvezzi, M. Melchiori, and M. Vincini. Information integration: The momis project demonstration. In *VLDB 2000, Proc. of 26th International Conference on Very Large Data Bases, 2000, Egypt*, 2000.
5. D. Beneventano, S. Bergamaschi, F. Guerra, M. Vincini, and M. Vincini. Exploiting extensional knowledge for a mediator based query manager. In *Convegno Nazionale su Sistemi Evoluti per Basi di Dati - SEBD2001, Venezia*, 2001.

6. D. Beneventano, S. Bergamaschi, S. Lodi, and C. Sartori. Consistency checking in complex object database schemata with integrity constraints. *IEEE Transactions on Knowledge and Data Engineering*, 10:576–598, July/August 1998.
7. D. Beneventano, S. Bergamaschi, C. Sartori, and M. Vincini. ODB-QOPTIMIZER: A tool for semantic query optimization in oodb. In *Int. Conference on Data Engineering - ICDE97*, 1997. <http://sparc20.dsi.unimo.it>.
8. S. Bergamaschi, S. Castano, D. Beneventano, and M. Vincini. Semantic integration of heterogenous information sources. *Journal of Data and Knowledge Engineering*, 36(3):215–249, 2001.
9. T. Catarci and M. Lenzerini. Representing and using interschema knowledge in cooperative information systems. *Journal of Intelligent and Cooperative Information Systems*, 2(4):375–398, 1993.
10. R. G. G. Cattell, editor. *The Object Database Standard: ODMG93*. Morgan Kaufmann Publishers, San Mateo, CA, 1994.
11. O. M. Duschka and M. R. Genesereth. Answering recursive queries using views. In *Proc. of the Sixteenth ACM SIGMOD Symposium on Principles of Database Systems*, 1997.
12. S. Gnanaprakasam E. Lambrecht, S. Kambhampati. Optimizing recursive information-gathering plans. In *Proceedings of the International Joint Conference on Artificial Intelligence, IJCAI 99*, 1999.
13. R. Hull and R. King et al. Arpa i<sup>3</sup> reference architecture, 1995. Available at [http://www.isse.gmu.edu/I3\\_Arch/index.html](http://www.isse.gmu.edu/I3_Arch/index.html).
14. Z. G. Ives, D. Florescu, M. Friedman, A. Y. Levy, and D. S. Weld. An adaptive query execution system for data integration. In *Proc. ACM SIGMOD Int. Conf. on Management of Data, USA*, 1999.
15. C.A. Knoblock J.L. Ambite. Flexible and scalable query planning in distributed and heterogeneous environments. In *Proc. of the 4th Int. Conf. on Artificial Intelligence Planning Systems. AAAI*, 1998.
16. Y. Papakonstantinou, S. Abiteboul, and H. Garcia-Molina. Object fusion in mediator systems. In *VLDB Int. Conf.*, Bombay, India, September 1996.
17. Y. Papakonstantinou, A. Gupta, and L. M. Haas. Capabilities-based query rewriting in mediator systems. *Distributed and Parallel Databases*, 6(1):73–110, 1998.
18. R. Pottinger and A. Y. Levy. A scalable algorithm for answering queries using views. In *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, pages 484–495. Morgan Kaufmann, 2000.
19. M.T. Roth and P. Scharz. Don't scrap it, wrap it! a wrapper architecture for legacy data sources. In *Proc. of the 23rd Int. Conf. on Very Large Databases*, Athens, Greece, 1997.
20. I. Schmitt and C. Türker. An Incremental Approach to Schema Integration by Refining Extensional Relationships. In G. Gardarin, J. French, N. Pissinou, K. Makki, and L. Bougamin, editors, *Proc. of the 7th ACM CIKM Int. Conf. on Information and Knowledge Management, November 3–7, 1998, Bethesda, Maryland, USA*, pages 322–330, New York, 1998. ACM Press.
21. R. Yerneni, Y. Papakonstantinou, S. Abiteboul, and H. Garcia-Molina. Fusion queries over internet databases. In *Advances in Database Technology - EDBT'98, 6th International Conference on Extending Database Technology*, volume 1377 of *Lecture Notes in Computer Science*, 1998.