

# A Semantic Matching Approach for Making Heterogeneous Sources Interoperable

Michel Schneider, Damien Thevenet

LIMOS, Complexe des Cézeaux, 63173 AUBIERE Cedex  
{michel.schneider@isima.fr, thevenetdamien@free.fr}

**Abstract.** Approaches to make multiple sources interoperable were essentially investigated when one are able to resolve a priori the heterogeneity problems. This requires that a global schema must be elaborated or that mappings between local schemas must be established before any request can be posed. The object of this paper is to study to what extend a mediation approach can be envisaged when none of these features are a priori available. Our solution consists in matching a query with each of the local schema. Such a solution is particularly suitable when sources are liable to evolve all the time. We are investigating this solution by considering the interoperability of heterogeneous XML sources. Local schemas are represented in the OWL language. Queries are formulated using an XQUERY-like language. Matching of names is solved by using an ontology of the domain. We have developed a prototype and conducted a number of experiments to evaluate the capacity of the approach.

**Keywords :** Interoperability, Heterogeneity, Mediation, Matching.

## 1 Introduction

The interoperability of multiple heterogeneous sources represents an important challenge considering the proliferation of numerous information sources both in private networks (intranet) and in public networks (internet). Heterogeneity is the consequence of the autonomy: sources are designed, implemented and used independently. Heterogeneity can appear for different reasons: different types of data, different representations of data, different management software packages. The interoperability consists in allowing the simultaneous manipulation of these sources so as to link the data which they contain. It is necessary to make different sources interoperable in numerous domains such as electronic business, environment, economy, medicine, genomics.

Interoperability problems occur in very different ways depending on whether sources are structured (data bases), semi-structured (HTML or XML pages), non-structured (any

file). The access interfaces also influences the possibilities of interoperability. For example two data bases can be difficult to make interoperable when they are only accessible through specific web interfaces.

One interoperability approach which has been studied for several years is based on mediation [17], [4]. A mediator analyzes the query of a user, breaks it down into sub-queries for the various sources and re-assembles the results of sub-queries to present them in a homogeneous way. The majority of mediation systems operate in a closed world where one knows a priori the sources to make interoperable. There are several advantages to this. First it is possible to build an integrated schema which constitutes a reference frame for the users to formulate their queries. Then it is possible to supply the mediator with various information which are necessary for the interoperability and particularly to resolve heterogeneity problems. The different kinds of heterogeneity to be resolved are now clearly identified: heterogeneity of concepts or intentional semantic heterogeneity; heterogeneity of data structures or structural semantic heterogeneity; heterogeneity of values or extensional semantic heterogeneity. Different solutions have been studied and experimented on to solve these problems. For example we can cite the work of [6] and [8]. From these initial investigations, very numerous works intervened to propose automatic approaches of integration of schemas. An approach was particularly investigated: the mapping of schemas. It led to the elaboration of several systems such as SEMINT, LSD, SKAT, DIKE, COMA, GLUE, CUPID. One will find analyses and comparisons of such systems in [14] or [5] or [12]. The practical aspects of the application of such systems are discussed in [1]. The role of ontologies was also investigated. In [2] and [11], the interest of ontologies for the semantic interoperability is underlined. Several approaches of integration of information based on ontologies were suggested. One will find a synthesis of it in [16]. It is necessary also to quote the work of [9] suggesting a logical frame for the integration of data. In all these works, the objective is to build a global schema which integrates all the local schemas.

When one operates in an evolutionary world where sources can evolve all the time, the elaboration of a global schema is a difficult task. It would be necessary to be able to reconstruct the integrated schema each time a new source is considered or each time an actual source makes a number of changes. To overcome these drawbacks an another approach has also been investigated : the query based approach. In this approach no global schema is required. The integration problems are solved during querying. Three main systems can be classified in this category: Information Manifold system, InfoSleuth system and Singapore system. Information Manifold [10] uses sources capabilities to determine how a query can be solved. InfoSleuth [13] is an agent-based system using ontologies for performing information gathering and analysis tasks. Singapore system [3] proposes an object language to formulate exact or fuzzy queries.

In this paper we suggest an approach which can also be considered as query-oriented. It is based on a semantic matching between the user query and each source schema. The user formulates its query by using its knowledge of the domain. Only the sources whose schemas match with the query are considered. The user query is rewritten for each of these sources according to its information capacity. These sources are then interrogated. Results are formatted and integrated. This approach offers several advantages. The rewriting process is simpler. A new source can be inserted at any time. The only obligation is to provide an adequate representation of this source.

The paper is organised as follows. In section 2 we give an overall presentation of our approach. Section 3 is devoted to the query language and section 4 to the OWL

representation of sources. In section 5 we explain the main features of our matching algorithm. Section 6 is relative to the rewriting of a query. Section 7 is devoted to some experiments with a prototype of the system. Section 8 presents a number of conclusions and perspectives.

## 2 Presentation of the approach

The approach which we propose does not use a global schema. The user thus formulates his query by using his implicit knowledge of the domain or by making an explicit reference to an ontology of the domain.

The matcher is the central element of the system. It receives the user query, and has the task to determine if this query can be applied to a data source (figure 1). To achieve this processing, it possesses a representation of each data source in a common formalism (we propose OWL to support this formalism, cf section 3). It must search for a correspondence between the query and each source by taking into account the terms and the structure of the query. Intuitively, so that a source can answer a query, the terms of the query must correspond to those of the source and the structure of the query must correspond to that of the source.

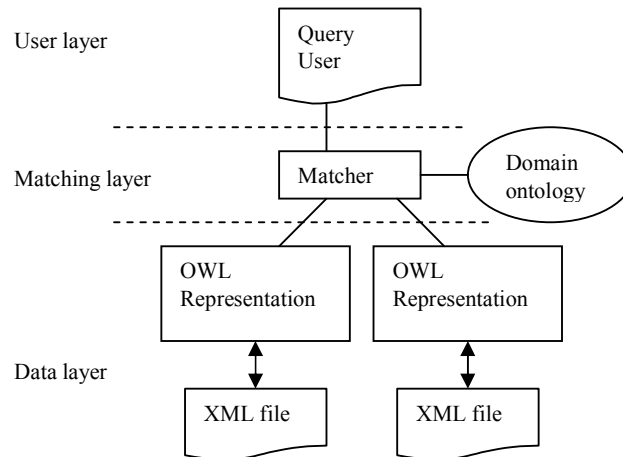


Figure 1: The architecture of our system

We propose a query language based on a simplified version of XQUERY. The structure of a query is thus defined through the various paths which appear in the clauses FOR, LET, WHERE. A correspondence is established with a source if each of these paths has a correspondent in the OWL representation of the source. More exactly, let  $E_1, E_2, \dots, E_k$  be a path. There is correspondence if one can find in the OWL representation classes  $C_1, C_2, \dots, C_k$  such that  $C_j$  is a synonym or hyponym of  $E_j$  for  $j \in [1, k]$  and such that every couple of classes  $C_i, C_{i+1}$  for  $i \in [1, k-1]$  is connected by a composition of properties in the OWL representation. In other words, it is necessary to find a subset of the OWL

representation which is subsumed by the path. The notions of synonym and hyponym are defined through the ontology of the domain.

For example consider the following query specified with our simplified XQUERY language :

```
Q : for $a in supplier, $b in customer
    where $b/name = "Ronald" and $a/region = $b/region
    return <b> { $b } </b>
```

It looks for customers with name "Ronald" and living in the same region as a supplier. The user supposes that a customer has a name and that both a customer and a supplier have a region.

Consider the two semi-structured sources of figure 2.

It is straightforward to infer that the query matches with the first source since the supplier element and the customer element both have a son element the name of which is region. The matcher must then check that this son element occurs only once. A matching for the second source cannot be inferred so easily. First the matcher must discover that buyer is an hyponym of customer. Then it must scan the hierarchy upward in order to establish that a supplier and a buyer are both connected to a unique region. So this second source is also a candidate for a rewriting of query Q.

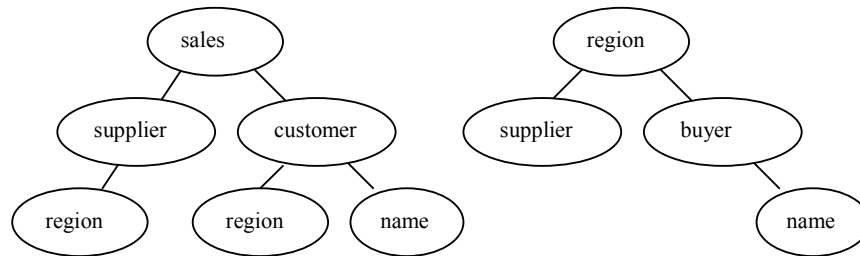


Figure 2 : Two semi-structured sources

We are now able to summarise the working of our system.

In the first phase the system initializes the connection with the ontology and gets back the names of the classes and the properties in the OWL representation. The system is then ready to handle queries.

In the second phase, when the system receives a query, it first interrogates the ontology to retrieve the synonyms and the hyponyms of the terms of the request. It then initiates the operation of matching for each of the paths of the query. Several rewriting possibilities can be proposed. To avoid inconvenient rewritings, we consider only the hyponyms of levels 1 to 3.

The third phase corresponds to the execution of one or several of these rewritings on the sources concerned. It may be necessary to transform the rewritings. This operation is performed with a wrapper associated to each source.

### 3 Query language

Our query language is a simplified version of XQuery. The user will have the possibility of using the FLWR construct of XQUERY with limitations indicated below.

- Since the user does not know the documents which can provide an answer, names of documents are omitted. So the root of a path is the name of an element. The system will make searches in all the documents having the root names in their description.

- Since the user does not know the structure of the data sources, it is not possible for him to decide if a term corresponds to an element or to an attribute. However he has the possibility of using the symbol '@' to indicate that he wants to search for an attribute. The system will first look for a correspondence with an attribute, but if this is not possible, it will continue its search on elements. If the symbol '@' is not present, the search will be made at the same moment on elements and attributes.

- Also, it is impossible for the user to know whether two elements are directly connected or if there are one or several intermediate elements. So it is not possible to differentiate the descent of one level "/" and the descent of several levels "//". So the system will be responsible for testing the descent at several levels.

- The functions of XPATH are not implemented in the simplified version.

- No difference is made in the query user between lower case and upper case letters.

The system will make sure the exact writing of a term is retrieved for the rewriting of the request.

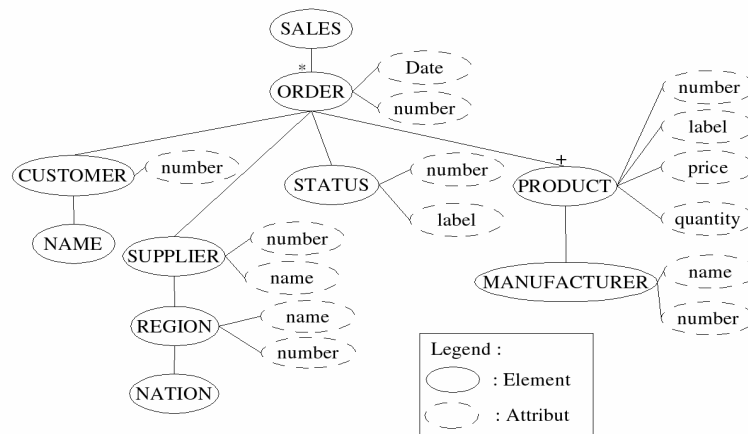


Figure 3: The semi-structured source A

### 4 OWL representation of the sources

To make the matching, the system could directly work on the DTD or the XML schema of the sources. We preferred to describe every source in OWL for various reasons. At first, as we want to operate with several types of sources, the matching will be implemented

with a single kind of description and will be so independent from the types of the sources. Then we foresee improved versions of the matcher, in particular by placing semantic annotations in the description of sources. The management of these annotations will be made more rationally with an OWL representation. Finally it is possible to take advantage of OWL's properties to implement the matching.

We elaborated an algorithm with which a DTD can be mapped into an OWL representation. This mapping is bijective: from the OWL representation, it is possible to regenerate the DTD. In the following we give the principles of the mapping.

The main idea is to represent every element of the DTD by an OWL class. Every father-son link between two elements is then represented by an OWL property. An attribute is also represented by a property. When a father element has only a single son element, the cardinality of this son is represented by creating a restriction on the property connecting the two elements. When the father element is a complex element, we add an intermediate class to be able to express correctly all the cardinalities.

```

<owl:Class rdf:ID="ORDER"><rdfs:subClassOf> <owl:Restriction>
<owl:onProperty rdf:resource="#ORDER.&complex1;"/>
<owl:cardinality rdf:datatype="&xsd;nonNegativeInteger"> 1 </owl:cardinality>
</owl:Restriction> </rdfs:subClassOf> </owl:Class>
<rdf:Property rdf:ID="ORDER.&complex1;"/>
<rdfs:domain rdf:resource="#ORDER"/> <rdfs:range rdf:resource="#&complex1;"/>
</rdf:Property>
<owl:Class rdf:ID="&complex1;"/>
<rdfs:subClassOf> <owl:Restriction>
<owl:OnProperty rdf:resource="#&complex1;.CUSTOMER"/> <owl:cardinality
rdf:datatype="&xsd;nonNegativeInteger"> 1 </owl:cardinality> </owl:Restriction> </rdfs:subClassOf>
.....
<rdfs:subClassOf> <owl:Restriction>
<owl:OnProperty rdf:resource="#&complex1;.PRODUCT"/> <owl:minCardinality
rdf:datatype="&xsd;nonNegativeInteger"> 1 </owl:minCardinality> </owl:Restriction> </rdfs:subClassOf>
</owl:Class>
<rdf:Property rdf:ID="&complex1;.CUSTOMER">
<rdfs:domain rdf:resource="#&complex1;"/>
<rdfs:range rdf:resource="#CUSTOMER"/>
</rdf:Property>

```

Figure 4: OWL representation for the ORDER element and its sons in source A

Agreements for the names of classes and properties are as following. The class representing an element will be named with the name of the element. For an intermediate class (associated to a complex element), the name of the class will contain the names of elements with their separator, quite in brackets. When this name is long, an entity can be used. A property between two classes will carry the two names separated by a point. For attributes, the symbol '@' is used to separate the name of the class and the name of the attribute.

As an example let us consider the element ORDER of the source A, the schema of which is shown in figure 3.

In the DTD, the definition of this element is:

```
<!ELEMENT ORDER(CUSTOMER, STATUS, SUPPLIER, PRODUCT+)>
```

In order to obtain its OWL representation, a class ORDER is created and also an intermediate class the name of which is (CUSTOMER, STATUS, SUPPLIER, PRODUCT+). For clearer understanding, the entity &complex1 is introduced to replace this name in the OWL file. Then a property connecting ORDER with the complex class is created, and the cardinality in the class ORDER is restricted. In the definition of the complex class the limitations of cardinalities are introduced for each of the elements. For CUSTOMER, STATUS and SUPPLIER, the cardinality is forced to be 1. Then properties are created to connect the complex class with each of the classes CUSTOMER, STATUS, SUPPLIER and PRODUCT (figure 4).

Using the same principles, it is possible to design an algorithm which maps a relational schema into a similar OWL representation. So our approach can be extended to deal also with relational sources.

## 5 Matching algorithm

We have to find a matching for each of the paths of the query.

To make the matching we represent a path as a tree having normal nodes and condition nodes. For example the path

```
order[customer/name="Pierre"] /product[price>15]
```

is represented by the tree in figure 5. In the tree, to every node is associated a simple path composed only of a succession of terms separated by the symbol “/”.

The matching of a path is then made through two main functions matchSimplePath(SP) and matchPath(P).

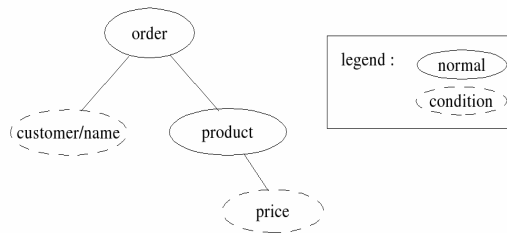


Figure 5: The tree corresponding to a path

Let  $SP_i$  be a simple path associated to node  $N_i$ . The function matchSimplePath( $SP_i$ ) looks in the OWL representation for the simple paths  $SP_{ij}$  which have a matching with  $SP_i$ . For example, let  $SP_i = E_1/E_2/E_3$ . A path  $SP_{i1}$  matches with  $SP_i$  if  $E_1, E_2, E_3$  have correspondents  $C_1, C_2, C_3$  in the OWL representation and if  $C_1$  is connected to  $C_2$  and  $C_2$  is connected to  $C_3$ .  $E_j$  corresponds to  $C_j$  if  $E_j$  or an  $E_j$ 's synonym or an  $E_j$ 's hyponym is identical to  $C_j$ .  $C_1$  is connected to  $C_2$  if  $C_1$  and  $C_2$  are connected either by a direct or inverse property or either by a composition of direct and inverse properties.

The set of paths  $SP_{ij} j \in [1, k]$  which have a matching with  $SP_1$  is associated to every node  $N_i(SP_i)$ . A node will thus be represented by  $N_i(SP_i, SP_{ij} j \in [1, k])$ .

The function  $matchPath(P)$  tests whether if a correct assembly of simple paths can be found which corresponds to the tree of  $P$ . Let  $N_i$  be a node of the tree and  $N_{i+1}$  one of its sons. One says that the assembly between the simple path  $SP_{im}$  of  $N_i$  and the simple path  $SP_{i+1,n}$  of  $N_{i+1}$  is correct if the last element of  $SP_{im}$  is connected with the first element of  $SP_{i+1,n}$ . The function  $matchPath(P)$  supplies all the possible correct assemblies for  $P$ . Each of these assemblies represents a path in the source which has a matching with  $P$ .

## 6 Rewritings of a query

To rewrite a query with regard to a source one looks for a rewriting of each of its paths. The rewriting of a path  $P_1$  then consists in replacing it in the query by one of the paths  $P_{ik}$  which matches with  $P_1$  and in inserting the navigation operators between the elements. When in  $P_{ik}$  one moves from a class  $C_1$  to a class  $C_2$  by a direct property, we only insert the descent operator  $//$  between the corresponding elements into  $P_1$ . If one moves from  $C_1$  to  $C_2$  by an inverse property, we insert the ascent operator.

At the end of this stage one can obtain several rewritings for a query.

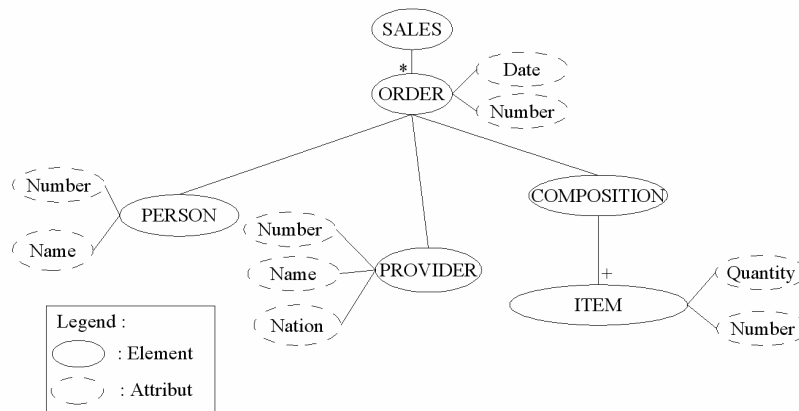


Figure 6: The semi-structured source B

## 7 Prototype and experiments

The prototype which we built implements the architecture presented in figure 1. We incorporated the tool SAXON-B [15] to access the OWL representations. We used the ontology WORDNET as the domain ontology. Since WORDNET is in fact a general ontology, we shall use sources for our experiments which do not contain highly



specialized terms. WORDNET is thus used only to provide synonyms and hyponyms of terms and to control that a succession of terms along a path is acceptable. Access to WORDNET is made through the JAVA API Java WordNet Library [7]. The body of the matcher is written in JAVA. We have implemented the two matching functions described in section 5.

Our experiments were conducted on source A already presented in figure 3 and on source B presented in figure 6.

We have submitted different queries to the prototype. We show the results obtained with the two sources A and B.

**Query 1 :** {order[customer/name="Pierre"]/product[price>15]}

Rewritings for source A:

1: {//ORDER[./CUSTOMER//NAME = "Pierre"] //PRODUCT[./@price>15]}

Rewritings for source B:

1:

For this query the matcher proposes a correct rewriting for source A. It does not use any synonym or hyponym. No rewriting is proposed for source B.

**Query 2 :** {for \$a in supplier where \$a/nation="FRANCE" return \$a}

Rewritings for source A:

1: {for \$a in //SUPPLIER where \$a//NATION = "FRANCE" return \$a}

Rewritings for source B:

1: {for \$a in //PROVIDER where \$a/@Nation = "FRANCE" return \$a}

In source A, NATION is an element and in source B, it is an attribute. In both cases, the matcher provides the correct rewriting.

**Query 3 :**

{for \$a in supplier, \$b in manufacturer where \$a/name=\$b/name and \$a/nation = "FRANCE" return \$a}

Rewritings for source A:

1: {for \$a in //SUPPLIER, \$b in //MANUFACTURER where \$a/@name = \$b/@name and \$a//NATION = "FRANCE" return \$a}

Rewritings for source B:

1:

For source B, the matcher does not provide any rewriting. For source A, it proposes a unique pertinent rewriting.

A rewriting such as:

{for \$a in //SUPPLIER, \$b in //MANUFACTURER where \$a//REGION/@name = \$b/@name and \$a//NATION = "FRANCE" return \$a} is also provided by our matching algorithm. This rewriting comes from the fact that there exists another attribute "name" of element REGION which can be reached from SUPPLIER. This rewriting contains strictly rewriting 1 and is not pertinent for the user. It is filtered in an additional step by using the following rule: "if a rewriting path P1 is a sub-path of another rewriting path P2 with the same starting node, delete P2 from the set of solutions".

**Query 4 :** {for \$a in person, \$b in supplier where \$a/name=\$b/name return \$a}

Rewritings for source A:

1: {for \$a in //MANUFACTURER, \$b in //SUPPLIER where \$a/@name = \$b/@name return \$a}

2: {for \$a in //CUSTOMER, \$b in //SUPPLIER where \$a//NAME = \$b/@name return \$a}

3: {for \$a in //NAME, \$b in //SUPPLIER where \$a = \$b/@name return \$a}

Rewritings for source B:

1: {for \$a in //PERSON, \$b in //PROVIDER where \$a/@Name = \$b/@Name return \$a}

The matcher provides three rewritings for source A since SUPPLIER, MANUFACTURER, NAME are hyponyms of PERSON (at level 3). The rewritings 1 and 2 are both pertinent and have immediate interpretation for a user. Rewriting 3 can surprise a user. It comes from the ambiguity of using in the schema of source A an element such NAME which is a hyponym of PERSON. In fact this rewriting is redundant with rewriting 2: its path is a sub-path of rewriting 2 and both terminate at the same element. So, rewritings 2 and 3 give the same result when executed on the source. We can use another filtering rule based on sub-paths. In this case we eliminate the rewriting 2 and keep only the rewriting 3.

Rewritings such as:

{\$a in //MANUFACTURER, \$b in //SUPPLIER where \$a/@name = \$b//REGION/@name return \$a}

are provided by our matching algorithm. Like for query 3, there are filtered in the additional step.

Our matcher filters also rewritings such as: {for \$b in //SUPPLIER, \$a in //SUPPLIER where \$a/@name = \$b/@name return \$a} which are correct but which correspond to a truth assertion and do not give pertinent results.

If we use only the hyponyms of level 1 for this query, the matcher give no answer for source A. This example shows clearly the interest of hyponyms, but also the problems which they can pose when confronting it to ambiguous schemas.

For source B the matcher provides a unique rewriting which is pertinent.

## 8 Conclusion and perspectives

Through the results obtained, it appears that our system is able to find data from an intuition of the user, intuition expressed through an implicit vision of the domain compatible with the ontology.

The main limitation comes from the fact that the system can propose irrelevant rewritings for a query. Most of them can be eliminated by adequate filtering rules. One can also act on the exploration depth in the ontology.

The quality of the ontology is highly important to avoid irrelevant rewritings. Ontology WORDNET used for our experiments is too general and contributes to sending back too many solutions. General terms such as name, number, ... contribute to increase the number of irrelevant solutions. It would so be necessary to minimize their use in the sources schemas and to resort to more precise terms.

More elaborated solutions exist to deal with this problem. A solution which we are investigating at present consists in placing annotations in the OWL representation at the level of classes or properties. These annotations will be exploited by the matcher to take into account semantic features (sense of a term, meaning of a property). These annotations could be installed manually by the administrator of the source or automatically by the system by seeking the opinion of the users when several rewritings are possible. To help the matching one can ask the user to clarify his query if the system detects some ambiguities.

We think that these improvements could result in an efficient system.

The system can be extended to deal with other types of sources (relational, object).

The main advantage of our approach is its robustness with regard to the evolution of sources. When a new source is inserted, it is sufficient to elaborate its OWL representation so that it can be exploited by the system. When a source evolves, it is sufficient to reshape its OWL representation.

We are also engaged in another improvement of our prototype in order to allow the join of results coming from different sources. In that case a query is rewritten in several sub-queries, each sub-queries being relative to a different source. Our matching algorithm can be easily adapted for this more general situation. It is necessary to look for sub-paths in different sources and to impose a join condition between sub-paths (the terminal node of a sub-path must correspond with the start node of another sub-path).

Another important point concerns the performances of our approach. Searching in a source the compatible paths of a query is the critical stage. We investigate different solutions in order to permit the system to scale up.

Such a system can be very useful for different applications. Incorporated into an intranet system, it would allow a user to reach the data sources without knowing their schemas, by being based only on the domain ontology. In a P2P system, it could be installed on some peers or on the super-peers to facilitate access to data by their semantics. The only obligation for a peer would be to publish its data by using the OWL representation.

## References

1. Bernstein P. A., Melnik S., Petropoulos M., and Quix C.: Industrial-strength schema matching. SIGMOD Record, Vol. 33, No 4, pp 38-43, 2004.
2. Cui Z., Jones D., O'Brien P.: Issues in Ontology-based Information Integration. IJCAI, Seattle, 2001.
3. Domenig R., Dittrich K.R.: A Query based Approach for Integrating Heterogeneous Data Sources. CIKM 2000, pp. 453-460.

4. Garcia-Molina H., Papakonstantinou Y., Quass D., Rajaraman A., Sagiv Y., Ullman J., Vassalos V. and Widom J.: The Tsimmis approach to mediation: Data models and languages. *Journal of Intelligent Information Systems*, Vol. 8, No. 2, pp. 117-132, 1997.
5. Hai Do H., Melnik S., Rahm E.: Comparison of Schema Matching Evaluations. *Web, Web-Services, and Database Systems*. pp 221-237, 2002.
6. Hull R.: Managing semantic heterogeneity in databases: A theoretical perspective. *Proc. of the Symposium on Principles of Database Systems (PODS)*, Tucson, Arizona, pp. 51-61, 1997.
7. JWNL. Java WordNet Library. <http://sourceforge.net/projects/jwordnet>.
8. Kedad Z., Métais E.: Dealing with Semantic Heterogeneity During Data Integration. *Proc of the International Entity Relationship Conference*, pp. 325-339, 1999.
9. Lenzerini M.: Logical Foundations for Data Integration. *SOFSEM 2005*. pp 38-40, 2005.
10. Levy A.L., Rajaraman A., Ordille J.J.: Querying Heterogeneous Information Sources Using Source Descriptions. *VLDB 1996*, pp. 251-262
11. Missikoff M., Taglino F., 2004. An Ontology-based Platform for Semantic Interoperability. *Handbook on Ontologies*. pp 617-634, 2004.
12. Mohsenzadeh M., Shams F., Teshnehlab M.: Comparison of Schema Matching Systems. *WEC (2)*, pp 141-147, 2005.
13. Noding M.H., Fowler J., Perry B.: Active Information Gathering in InfoSleuth. *CODAS 1999*, pp. 15-26.
14. Rahm E., Bernstein P.A.: A survey of approaches to automatic schema matching. *VLDB Journal* 10(4), pp 334-350, 2001.
15. Saxon. SAXON: The XSLT and XQuery Processor. <http://saxon.sourceforge.net/>.
16. Wache H., Vogeles T., Visser U., Stuckenschmidt H., Schuster G., Neumann H. and Hubner S.: Ontology-based integration of information - a survey of existing approaches. In Stuckenschmidt, H., ed., *IJCAI-01 Workshop: Ontologies and Information Sharing*, pp 108-117, 2001.
17. Wiederhold, G.: Mediators in the architecture of future information systems. *IEEE Computer*, Vol. 25, No 3, pp.38-49, 1992.