

An approach for semantic integration of heterogeneous data sources

Giuseppe Fusco¹, **Lerina Aversano**^{Corresp. 1}

¹ Department of Engineering, University of Sannio, Benevento, BN, Italia

Corresponding Author: Lerina Aversano
Email address: aversano@unisannio.it

Data Integration is the problem of combining data distributed among heterogeneous information sources, and providing the user with a unified view of these data. Information sources can be structured, semi-structured or unstructured, such as database, HTML pages, files, and so on. Actually, the problem of dealing with data provided by heterogeneous source is very relevant in real applications. In this paper it is proposed an approach and the related software system for supporting the semantic integration of heterogeneous sources.

An Approach for Semantic Integration of Heterogeneous Data Sources

Giuseppe Fusco¹ and Lerina Aversano²

¹Department of engineering

²University of Sannio

Corresponding author:

Lerina Aversano¹

Email address: aversano@unisannio.it

ABSTRACT

Data Integration is the problem of combining data distributed among heterogeneous information sources, and providing the user with a unified view of these data. Information sources can be structured, semi-structured or unstructured, such as database, HTML pages, files, and so on. Actually, the problem of dealing with data provided by heterogeneous source is very relevant in real applications. In this paper it is proposed an approach and the related software system for supporting the semantic integration of heterogeneous sources.

1 INTRODUCTION

The large availability of data within the enterprise context and even in any inter-enterprise context, arise the problem of managing information sources that do not use the same technology, or, do not have the same data representation, or that have not been designed according to the same approach. The issues emerging to face up heterogeneity problems are the following:

- The autonomous and dynamic nature of information sources: each source manages its data independently of each other.
- The data access: Each source has its own data access mechanisms. In a process of integrating heterogeneous sources, one must also consider the problem of how to retrieve data using a global and general access mechanism.
- Data reconciliation: data from different sources need to be interpreted, that is, transformed into a common representation. Therefore, they must be converted, reconciled and combined.

Over the years, several data integration solutions have been proposed:

- Distributed databases can be considered the first attempt to integrate databases. Data, instead of being stored on a single machine, is stored on different machines. Compared to the centralized case, database schemas are made more complicated by the need to physically distribute data over multiple machines. Distributed databases require the complete integration of existing systems into a single homogeneous database. This is difficult to achieve both technical (prohibitive conversion costs) and organizational difficulties (existing DBMSs belong to different organizations).
- Federated databases have been proposed to address these limits. They are a set of multiple independent sources each of which can exchange information with the others. A connection is established for each pair of sources, and such architecture is particularly suitable when communications in the system occur predominantly between pairs of sources.

The solution often adopted consists of the cooperative information systems (Figure 1), in which there are two types of components: mediator and wrapper. The mediator coordinates the data flow between

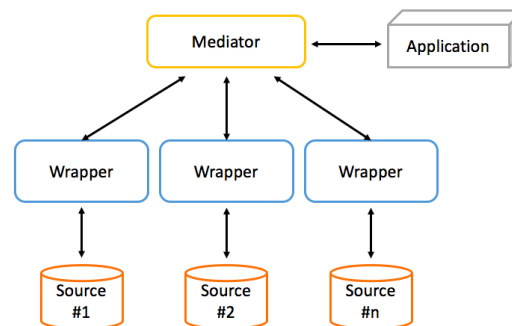


Figure 1. Architecture of a generic mediation system.

local sources and user applications. The mediator doesn't store data, since it only stores a virtual and global view of real data (or global schema) and the mappings between the global and local views. In this way, applications will run queries over the virtual view. It will then be the mediator to build queries for individual sources of information. Instead, wrappers are software components that interact directly with their respective local sources as follows:

- to translate the conceptual schema of the local source into a global language;
- to submit queries to local sources;
- to retrieve results by sending them to the mediator, which will provide the user with a unified result.

This approach allows to exploit semantic relationships between sources to provide users with a uniform interface (called mediated schema or global schema or global view) of sources, freeing them from manually managing each source.

This paper proposes a framework, Data Integration Framework (DIF), to support the analysis, acquisition and processing of data from heterogeneous sources. DIF allows the integration of heterogeneous data sources, enabling the use of a global and virtual view of the sources obtained from the analysis of the conceptual schema of acquired sources. The purpose of this view is to provide a unified view of real data, so that applications and users who use data will have the perception of accessing a single data source rather than multiple sources. In this context, the work faced aspects of acquisition, integration and processing of heterogeneous data sources.

The paper is organized as follows. Section 2 presents the state of the art, problems that characterize data integration and proposed solutions. Section 3 presents in detail the approach and architecture of the software system developed to support the integration of heterogeneous sources. Section 4 presents the DIF tool design and the main algorithms implemented. Section 5 presents a case study in order to show the results of the proposed solution. Finally, Section 6 concludes this paper by submitting concluding remarks and mentioning some research issues related to data integration that are not addressed in this paper.

2 RELATED WORK

Data integration systems using the mediation approach are characterized by an architecture (Figure 1) based on a global schema and a set of source schemas. The sources contain real data, while the global scheme provides a unified, integrated and reconciled view of local sources. The main components are: the mediator, which coordinates data flow between local sources and user applications, and wrappers, which directly interact with their respective local sources. Designing a data integration system is a complex task, which involves dealing with different aspects:

- Heterogeneity of sources: sources adopt different models and data storage systems. This poses problems in defining the schema/global view. The purpose is to provide a view with an appropriate level of abstraction for all data in the sources.

- Mappings between global schema and local sources: in literature, in order to model the correspondences between the schemes, different approaches [16] have been proposed as *global-as-view* and *local-as-view*.

With *global-as-view* (GaV) approach the global schema is expressed in terms of views over data sources. With *local-as-view* (LaV) approach the global schema is specified independently of the sources, and the relationships between global schema and sources are established by defining each source as a view built over the global schema. Differences between the two approaches are discussed in [16]. In order to overcome the limits of GaV and LaV approaches, techniques that combine the benefits of these approaches have also been proposed, mitigating their disadvantages in order to provide an alternative to data integration that is more flexible and scalable. The most interesting techniques are *GLaV* (Global and Local as View) [13][2] and *BGLaV* (BYU Global Local as View) [26].

- Schema matching activity: once the mapping approach is defined, it is necessary to define the methods and techniques to be used to generate mappings between the global and the local schema. The set of mappings is called *alignment*. A matching process [23] (Figure 2) defines an alignment (A') for each pair of schemas (o_1, o_2), making use of input parameters p if necessary (for example, thresholds, weights), a previously generated input alignment (A) and additional external resources r .

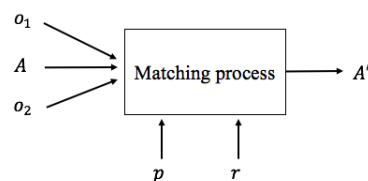


Figure 2. Matching process.

- Schema merging activity: generate the global schema based on mappings defined in the schema matching activity. A merging process (Figure 3) consists of integrating several existing schemes (o_1, o_2, \dots, o_n) into a single global schema (O) based on the correspondences generated by the schema matching process A , any input parameters p and external resources r . Different techniques

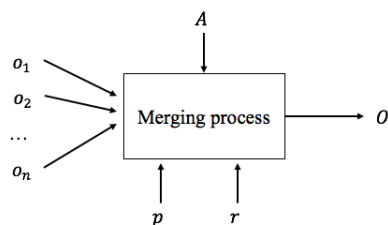


Figure 3. Merging process.

and methodologies about schema merging have been proposed in the literature [15] [8] [7] [14].

- Data storage: compared to managed data there are two approaches, called materialization and virtualization. With materialization, data is also present in the global schema. On the opposite, in the virtualization approach, data that resides in sources is only available when query processing activity is executed.
- Query processing activity: how to process a query posed over the global schema, that is how to express it in terms of a set of queries that can be processed by the sources acquired. In the LaV approach the proposed solutions consist of *query rewriting* (or *view-based query rewriting*) and *query answering* (or *view-based query answering*). In the GAV approach *query unfolding* techniques are used. The differences between the two approaches are discussed in [16].

- Data reconciliation: data from different sources need to be interpreted, that is, transformed into a common representation. Therefore, they must be converted, reconciled and combined.

To overcome problems due to semantic heterogeneity, it is useful to use *ontologies* [25]. Depending on how ontologies are used, data integration systems can adopt different approaches, such as *single ontology* (adopted in SIMS [2][1]), *multiple ontology* (adopted in OBSERVER [18][17]) and *hybrid* (adopted in KRAFT [20][21] and COIN [11])

2.1 Data Integration Systems

Several systems, methodologies and approaches have been proposed in literature to support data integration from heterogeneous sources, also based on ontologies [5]. Below is a brief description of the most interesting systems: TSIMMIS, GARLIC and MOMIS.

2.1.1 TSIMMIS

TSIMMIS (The Stanford-IBM Manager of Multiple Information Sources) [6] is based on an architecture that exposes a wrapper hierarchy (called Translators) and mediators.

TSIMMIS approach is global-as-view. Wrappers convert data to a common data model called OEM (Object Exchange Model) and mediators combine and integrate them. The global scheme consists of a set of OEM objects exported by wrappers to mediators. Mediators are specified using a language called Mediator Specification Language (MSL). Queries are expressed in MSL or in a specific language called LOREL (Lightweight Object Repository Language), an object-oriented extension of SQL. Each query is processed by a module, the Mediator Specification Interpreter (MSI).

It should be emphasized that TSIMMIS does not implement a real integration, as each mediator performs integration independently of each other. It means that does not exist the concept of a unified global scheme. The result of a query could be seen inconsistent and completely different from other mediators. This form of integration is called query-based.

2.1.2 GARLIC

GARLIC integration system is based on an architecture with Data Repositories at lowest level, which represent the data sources. Above each data repository we find a wrapper (called Repository Wrapper), which is responsible for communication between a data repository and the rest of the system. In addition, each wrapper ensures the transformation of the local schema of a source into a unified schema and transforming user queries into queries executable by data source.

The global schema has an object-oriented data model, managed by the Query Services and Runtime System components, and stored in the Metadata Repository, based on the ODMG standard. ODMG objects are exported by wrappers using Garlic Data Language (GDL), based on the ODL (Object Definition Language) standard.

Unlike the TSIMMIS system, there is no mediator concept in GARLIC, and the integration of ODMG objects from different sources is performed by wrappers.

2.1.3 MOMIS

MOMIS (Mediator Environment for Multiple Information Sources) [19][3] is a data integration system that manages structured and semistructured data sources. MOMIS is based on I^3 architecture [12], consisting of several wrappers and a mediator.

The integration methodology starts with an extraction activity where user uses a wrapper that transforms the structure of a data source into a ODL_3 (Object Definition Language) model based on descriptive logic. The integration process generates an integrated view of data sources using global-as-view approach, building the global schema incrementally. At the end of the MOMIS integration process, starting when the query is posed by the user over the global schema, the mediator generates a OQL_3 query and sends it to wrappers, which translate it into a query executable from the corresponding data source.

2.2 Discussions

In the previous subsections, the approaches used in mappings definition between the global schema and local ones have been described. Table 1 summarizes what has been said.

Based on the comparison approaches in Table 1 it is possible to observe that:

Table 1. Comparison between GaV and LaV.

	GaV	LaV
<i>Mapping</i>	Global schema expressed in terms of views over data sources	Data sources expressed in terms of views over global schema
<i>Query processing</i>	Query unfolding	Query rewriting/Query answering
<i>Global schema quality</i>	Exact or Sound	Complete or Exact
<i>Management effort</i>	High: data source changes affect the global schema and other sources	Low: data source changes only impact the global schema

- The LaV approach involves a priori presence of a global schema, which must be well-built in terms of concepts and relationships between them. If it is not well-built, or the integrated schemas differ greatly from each other, the global schema must be continually modified, also taking into account the previously integrated data sources. If not, the changes affect only the global schema.
- With GaV approach, the global schema is incrementally built: it is modified every time a new data source is integrated, adding and/or modifying concepts and relationships based on current and previously integrated data sources. Conversely, in the LaV case, changes don't impact previous integrated data sources (if the overall schema is well-built).
- In the context of semantic integration, the hybrid approach is surely the best solution but it reduces the reuse of local ontologies, since they have to refer to a common vocabulary.
- The LaV approach offers greater scalability when the number of integrated data sources increases, but when that number is relatively small, and the global schema is not well-built, the GaV approach increases the quality of global schema.

Based on these observations, and considering possible future reuse of local ontologies, it is possible to combine the presented approaches differently in order to support different cases and to present a data integration approach in order to provide different solutions as needed. The proposed approach, called *DIF* is based on these observations and seeks to combine the GaV and LaV approaches, exploiting ontologies to reach the goals.

In addition, a further issue is related to the interoperability between data integration systems, described in subsection 2.1, and systems outside their application context. TSIMMIS, GARLIS and MOMIS use their own description language for both local and global schemas, and queries. However, if a generic external application wants to communicate with one of the systems presented, it should know the specific query language and/or the specific language used to describe the schemas. The problem of translation between languages is widened if we consider interoperability with the Web. For this reason, in literature, were defined ontologies and languages to describe them and construct queries over them. Considering a wider scenario in data integration, the developed system, called Data Integration Framework (DIF), exploits the mechanism of ontologies.

Table 2. Comparison between data integration systems.

	TSIMMIS	GARLIS	MOMIS
<i>Integration approach</i>	GaV	GaV	GaV
<i>Integration type</i>	Query-based	Schema-based	Schema-based
<i>Data model</i>	OEM (Object Exchange Model)	GDL (Garlic Data Language)	ODL ₃
<i>Query language</i>	MSL (Mediator Specification Language) or LOREL (Lightweight Object Repository Language)	Object-oriented SQL-like	ODL ₃
<i>Query processing</i>	Query unfolding	Query unfolding	Query unfolding
<i>Structured data sources</i>	Yes	Yes	Yes
<i>Semistructured data sources</i>	Yes	Yes	Incomplete software
<i>Unstructured data sources</i>	Yes	Yes	No

3 DATA INTEGRATION FRAMEWORK

The proposed approach, called *DIF* (Data Integration Framework), is a generalization of both GaV and LaV approaches, based on data virtualization, and provides the possibility to define a mapping in one of the following ways:

- A correspondence between a view expressed in terms of the global schema and a view expressed in terms of the local schema (as in the GaV approach);
- A correspondence between a view expressed in terms of the local schema and a view expressed in terms of the global schema (as in the LaV approach);

In addition, to overcome problems due to semantic heterogeneity and to support interoperability with external systems, ontologies are used as a conceptual schema to represent both data sources to be integrated and the global schema, and therefore:

- Each mapping is defined as a correspondence between elements of ontologies: concepts (or classes), datatype properties, and object properties;
- Since the data virtualization approach is also used to define local ontologies, the construction of an ontology to represent a local source is guided by additional mappings, called *source-mappings*, defined as correspondences between elements of local ontology and elements that characterize the source itself (for example, for a relational source a mappings will be defined as a correspondence between an ontology concept and the table that represents it).

In the proposed solution, the query rewriting is used to reformulate a query posed over the global ontology into a set of queries posed over the local ontologies. This is due to the choice of using ontologies also to describe data sources to be integrated. In this way, though, the mediation process is not completed yet, since local ontologies don't contain real data. To complete the mediation process, a second query translation task is required to reformulate a query posed over the local ontology into a set of queries posed over the corresponding source.

Definition 3.1 (Data Integration Framework). *The data integration framework DIF is a 5-uple (O_g, O_l, A, MT, SML) where:*

- O_g is the global ontology, expressed in a LO_g logic.
- O_l is the local ontology, expressed in a LO_s logic.
- A (Alignment) is a set of mappings M_1, M_2, \dots, M_n between ontologies O_g and O_l . Each mapping M_i is a 5-uple (id, e_s, e_t, n, rel) where:
 - id is the unique mapping identifier;
 - e_s and e_t , respectively, are the elements of the source ontology O_s and target O_t . In the case of a GaV mapping type, O_s represents the local ontology and O_t the global one, vice versa in the case of a LaV mapping type;
 - n is a measure of confidence (typically within a range $[0, 1]$) that indicates the similarity between e_s and e_t ;
 - rel is a relationship between e_s and e_t (for example, equivalence, subsumption, disjunction).
- MT (Mapping Table) is a table whose rows represent an element e_g of the global ontology O_g and columns represent elements $e_{l1}, e_{l2}, \dots, e_{ln}$ of the local ontology O_l that are mapped to e_g .
- SML (Source Mapping List) is a set of mappings SM_1, SM_2, \dots, SM_n between the local ontology O_l and the correspondent data source S_i . Each mapping SM_i , called *source-mapping*, is a triple (id, src_k, dst_h) where:
 - id is the unique mapping identifier;
 - src_k is a source element of the local ontology O_l .

230 – dst_h is a destination element of the local data source S_i (for example, a table of a relational source).

231 The framework must be able to handle both the integration process and the mediation process
232 (Figure 4), making activities as automated as possible.

233 The integration process is divided into the following activities:

- 234 • *Source Wrapping*: for each source you want to integrate, you build an ontology to describe it. In
235 addition, source-mappings are defined between the ontology and the data source, which will be
236 subsequently used during the mediation process.
- 237 • *Schema Matching*: for each local ontology, mappings are generated between it and global ontology.
238 The matching activity generates mappings between a source ontology and a target one. Therefore,
239 considering as target ontology the local one, it is possible to generate LaV mappings. Conversely,
240 the followed approach will be GaV. Mappings are eventually validated or modified by the integrator
241 designer. If the number of data sources to be integrated is 1, global and local ontologies are the
242 same.
- 243 • *Schema Merging*: each local ontology, taking into account the set of mappings defined in the
244 previous activity, is integrated into the global ontology and the mapping table is updated. At this
245 stage, global ontology is built incrementally.

246 The mediation process, however, following a query submission, is divided into the following phases:

- 247 • *Global Query Processing*: a query posed over the global ontology is reformulated, through rewriting,
248 into a set of queries posed over local ontologies, using the mapping table generated at the end of
249 the integration process;
- 250 • *Local Query Processing*: each local query is reformulated into a set of queries over the correspond-
251 ing data source, using source-mappings generated in the source wrapping activity. This set of
252 queries, once executed, allows you to retrieve the real data.
- 253 • *Data Reconciliation*: extracted data from the previous activity is reconciled and combined before
254 being presented to the user.

255 Local and global ontologies are expressed in OWL-DL¹, whose basic elements are classes c , object
256 properties op and datatype properties dp . Instances i are not considered mapping because the data
257 management approach is virtualization rather than materialization.

258 4 DIF SUPPORTING TOOL

259 The tool, designed and developed to support the DIF framework, presents the typical architecture of
260 integration systems based on the mediation approach (Figure 5(a)), providing two main components:
261 *mediator* and *wrapper*.

262 According to Definition 3.1 and the description of the activities to be performed during integration
263 and mediation processes, the architecture (Figure 5(b)) is composed by *Acquisition*, *Integration* and
264 *Mediation* subsystems.

265 4.1 Source wrapping

266 Data sources that the framework allows to integrate are structured and semi-structured (in particular,
267 spreadsheet, JSON, and XML data sources). The source wrapping activity is performed by a class that
268 implements the *IWrapper* interface. The output of this activity, for each data source S , is a pair (O, SML)
269 composed by the local ontology O , which describes the data source S , and the associated source mapping
270 list SML .

¹<https://www.w3.org/TR/owl-features/>

271 **4.1.1 Relational data sources integration**

272 The system allows the integration of relational data sources via JDBC connection² and supported databases
273 are: MySQL, PostgreSQL, H2, DB2, Microsoft SQL Server, Oracle, Telid and MonetDB.

274 Relational data sources are connected to the framework by defining R2RML³ mappings. Each
275 R2RML mapping therefore represents a source-mapping, according with Definition 3.1. Local ontology
276 is generated by identifying conditions associated to the database tables [10] and, through identified
277 conditions, associating each database element (table, column) to the corresponding ontology (class,
278 datatype property, object property).

279 In addition to R2RML mapping, you can use a more compact notation, called axiom mapping,
280 consisting of three elements:

- 281 • MappingID: mapping identifier;
- 282 • Source: a SQL query posed over the database;
- 283 • Target: RDF triple containing variables that refer to the names of the columns mentioned in the
284 source query.

285 Each source source mapping $SM(id, src_k, dst_h)$ (Definition 3.1, contained in the source mapping list
286 SML , contains an OWL resource (or local schema element) src_k and an R2RML mapping (or an axiom
287 mapping) dst_h .

288 **4.1.2 Spreadsheet data sources integration**

289 Spreadsheet data sources are integrated with a new approach that seeks to examine the internal structure
290 of the tables in order to extract an ontology that reflects as much as possible the data source. The approach
291 is divided into several phases:

- 292 1. *Source Scanning*: the spreadsheet file is scanned in order to locate tables. At the end of the scan, a
293 text string that describes the table structure is produced.
- 294 2. *Parsing*: the text string is parsed in order to generate the ontology elements, the relationships
295 between them, and the physical location of cells within the spreadsheet data source. The output of
296 this step is a list of schema attribute tables.
- 297 3. *Analysis*: an analysis of the list of attribute tables built to the previous step is performed to aggregate
298 attributes with the same name in different concepts.
- 299 4. *Restructuring*: the generated ontology is refined in order to improve its quality.

300 Each source-mapping $SM(id, src_k, dst_h)$ contained the source mapping list SML (Definition 3.1)
301 contains an OWL resource (or local schema element) src_k and a data structure to track cells within the
302 spreadsheet data source dst_h .

303 **4.1.3 XML data sources integration**

304 XML data sources integration is based on its XSD schema [9]. If the XML schema does not exist, it is
305 generated. Possible mappings are:

- 306 • Class mapping: XML nodes mapped as classes are complex types and element-group declarations
307 (if they contain complex types). The XML schema supports two inheritance mechanisms (extensions
308 and restriction), which are also mapped as classes.
- 309 • Property mapping: if an XML node has a simple type, or is an attribute-group declaration, or is an
310 element-group declaration without additional associated complex types, it is mapped as properties,
311 and its domain is the class corresponding to the parent node. Attributes are treated as simple types.
312 In [9], instead, all element-groups and attributes-groups are mapped as classes.
- 313 • Relation mapping: if an element is a complex type and contains another element, whose type is
314 also complex, a relationship is created between the respective classes.

²<http://www.oracle.com/technetwork/java/javase/tech/index-jsp-136101.html>

³<https://www.w3.org/TR/r2rml/>

The algorithm, in the ontology generation step, receives the XSG graph of the XML schema (XML Schema Graph) input. Starting from the root node, a deep visit is performed, generating an XPath expression for each visited node. Each source-mapping $SM(id, src_k, dst_h)$ contained the source mapping list SML (Definition 3.1) contains an OWL resource (or local schema element) src_k and the XPath expression dst_h .

4.2 Schema matching

The goal of schema matching activity is to generate a set of mapping between local and global ontologies, which will then be validated by the user. The adopted approach generates mappings between classes, considering both semantic and contextual characteristics. Before to execute schema matching, a semantic annotation activity of ontologies is performed, whose output is a set of annotations AN , one for each e_i element of the schema, where each annotation AN_i is a triple $(tok_i, POS_i, sense_i)$ consisting of:

- tok_i : the token associated with the element e_i ;
- POS_i : the lexical category of the token tok_i ;
- $sense_i$: the meaning associated with tok_i token for the lexical category POS_i , obtained as the output of the disambiguation process.

In the semantic matching task, a semantic-based matcher is applied to all pairs (C_{Gi}, C_{Lj}) , where C_{Gi} is the i -th class of the global schema, while C_{Lj} is the j -th class of the local schema. The semantic matcher, for each pair (C_G, C_L) , generates the following information:

- *SemanticRel*: the type of semantic relation $(\equiv, \sqsubseteq, \sqsupseteq, idk)$;
- *SemanticSim*: is a coefficient $\in [0, 1]$ that specifies the reliability of the generated semantic relation..

Given n and m the number of local and global schema classes, respectively, the output of the semantic matching activity is:

$$(C_{Gi}, C_{Lj}) \Rightarrow \begin{matrix} i \in [1, m] \\ j \in [1, n] \end{matrix} \left\{ \begin{matrix} \textit{SemanticRel} \\ \textit{SemanticSim} \end{matrix} \right\}$$

The contextual matching activity generates mappings between classes taking into account how they are modeled, that is considering their properties. First, you must determine the equivalent properties between the two schemes. This is done by applying to all pairs (P_G, P_L) , where P_G and P_L are the properties of the global and local schema respectively, the syntax-based or the semantic-based matcher. The syntax-based matcher, by analyzing the syntactic structure of words, returns for each pair (P_G, P_L) a coefficient in a range $[0, 1]$. If the latter is greater than or equal to the β threshold value, a mapping is generated between P_G and P_L . The semantic-based matcher, instead, using WordNet to extract property sense, generates a mapping between P_G and P_L if there is an equivalence relation for at least one token pair. The semantic-based matcher is useful if synonyms are used to represent the same property and/or to discover mappings 1:n. Once the mappings have been discovered, it is possible to calculate, for all pairs (P_G, P_L) , the degree of contextual similarity, defined as *ContextualSim*, by applying the Jaccard measure:

$$\textit{ContextualSim}(C_G, C_L) = \frac{|P(C_G) \cap P(C_L)|}{|P(C_G) \cup P(C_L)|}$$

where $P(C_G)$ and $P(C_L)$ are the set of properties of the classes C_G and C_L respectively. The cardinality of the intersection of the two sets is equal to the number of existing mappings between the properties of the two classes. Given n and m the number of local and global schema classes, respectively, the output of the contextual matching activity is a set of pair as:

$$(C_{Gi}, C_{Lj}) \Rightarrow \begin{matrix} i \in [1, m] \\ j \in [1, n] \end{matrix} \left\{ \begin{matrix} \textit{ContextualSim} \\ M_{P(C_{Gi}), P(C_{Lj})} \end{matrix} \right\}$$

where $M_{P(C_{Gi}),P(C_{Lj})}$ is the set of mappings between the properties of classes C_{Gi} and C_{Lj} :

$$M_{P(C_{Gi}),P(C_{Lj})} = \left\{ \begin{array}{ccc} P_1(C_{Gi}) & \longleftrightarrow & P_1(C_{Lj}) \\ P_2(C_{Gi}) & \longleftrightarrow & P_2(C_{Lj}) \\ \dots & & \dots \\ P_k(C_{Gi}) & \longleftrightarrow & P_z(C_{Lj}) \end{array} \right\}$$

where k is the number of properties of the class C_{Gi} and z is the number of properties of the class C_{Lj} .
To determine which mappings can be returned to the user, a selection step is performed. The principle is that if there is a semantic relation *SemanticRel* and the degree of contextual similarity *ContextualSim* is greater than or equal to a threshold value, the corresponding mappings can be returned. By lowering these values, more weight is given to semantic characteristics rather than contextual ones. Given a semantic relation *Rel* and a threshold value α , the algorithm selects 1:1, 1:n, n:1 and n:m mappings between all pairs of classes (C_{Gi}, C_{Lj}) if there is a semantic relation equal to *Rel* and if $ContextualSim \geq \alpha$. If more mappings 1:n, n:1 and n:m for the same pair of classes (C_{Gi}, C_{Li}) have a threshold value greater than α , is returned the mapping with the largest number of classes. The output is the set of selected mappings.

The output of schema matching activity, according to the Definition 3.1, is an alignment A consisting of a set of mappings between the local and global schema classes, obtained after the selection step:

$$A \Rightarrow M(\{C_G\}_k, \{C_L\}_z) \Rightarrow (C_{Gi}, C_{Lj}) \Rightarrow \begin{array}{l} i \in [1, k] \\ j \in [1, z] \end{array} \left\{ \begin{array}{l} SemanticRel \\ Similarity \\ M_{P(C_{Gi}),P(C_{Lj})} \end{array} \right\}$$

Mappings can be 1:1, 1:n, n:1 and n:m.

4.3 Schema merging

The goal of schema merging activity, starting from user validated mappings, is the fusion between the local and global schema, generating a new virtual view. Schema merging activity is divided into two steps:

- In the first step, changes in the global schema are generated;
- In the second step, based on the proposed changes, the fusion of schemas is performed.

In the first step, given an input alignment A (the mappings list), the global schema G and the local one L , the new global schema T is initially created, which is initially equal to G , and the empty mapping list ML that will contain the mappings between L and T elements. Merging is performed by applying merge operators to each input mapping. Next, the local schema classes and relations, not included in the global schema, are added to T . The new resulting schema is modified by deleting redundant relationships and performing refactoring operations. The framework has an internal data structure to track changes to the new global schema T .

In the second step, given the changes produced and after deciding whether to validate or not, the real schema merging is performed.

Output of schema merging activity, besides to the new global schema T , according to the definition 3.1, also consists of the mapping table MT , whose rows represent a e_G element of the global schema G and columns represent the $e_{L1}, e_{L2}, \dots, e_{Ln}$ elements of the local schemas L that are mapped to e_G . Since it is possible to define complex mappings n:m, the mapping table will be a table whose rows represent an E_{Gi} expression of an element of the global schema G and the columns represent the expressions E_{Ljk} of the j -th element of the k -th local schema. A generic $MT[i, j]$ element of the mapping table represents, in fact, a mapping $M(id, E_{Gi}, E_{Ljk}, n, rel)$ between the expression of an element i -th of the global schema and an expression of an element j -th of the k -th local schema.

4.3.1 Mapping table

In the mapping table, according with the definition 3.1, rows represent elements of the global schema, and columns represent elements of the local schemes. Elements are generic OWL expressions and Table 3 shows the possible mappings in the mapping table:

	Global schema g	Local schemas $1, 2, \dots, n$
<i>CE mapping</i>	CE_g	$\bigcup_{i \in [1, n]} CE_i$
<i>OPE mapping</i>	OPE_g	$\bigcup_{i \in [1, n]} OPE_i$
<i>DPE mapping</i>	DPE_g	$\bigcup_{i \in [1, n]} DPE_i$

Table 3. Mapping table.

The framework, however, allows mappings to be defined in a generic way, without explicit reference to a global or local schema. For this reason, the framework must be configured by setting a parameter $dir = \{global, local\}$ indicating the direction of the mappings in such a way as to support queries reformulation. If not specified, it is assumed that in the rows there are expressions referring to the global schema and in the columns the expressions referring to the local schemes. When it is necessary to insert a mapping in the opposite direction, it is inverted.

The mapping table is represented using the EDOAL⁴ language. For example, we consider the following mapping:

$(Hospital, Infirmary) \Rightarrow$

$$\left\{ \begin{array}{c} \equiv \\ 0.4 \\ Name \longleftrightarrow Name \end{array} \right\}$$

For the mapping $M(Hospital, Infirmary)$, assuming to assign to the first schema the prefix $src\#$ and to the second schema the prefix $trg\#$, the mapping representation for the property $Name$ will be as follows:

```

380 <map>
381 <Cell>
382   <entity1>
383     <edoal:Property rdf:about="src#Name"/>
384   </entity1>
385   <entity2>
386     <edoal:Property>
387       <edoal:or rdf:parseType="Collection">
388         <edoal:Property>
389           <edoal:and rdf:parseType="Collection">
390             <edoal:Property rdf:about="trg#Name"/>
391             <edoal:PropertyDomainRestriction>
392               <edoal:class>
393                 <edoal:Class rdf:about="trg#Infirmary"/>
394               </edoal:class>
395             </edoal:PropertyDomainRestriction>
396           </edoal:and>
397         </edoal:Property>
398       </edoal:or>
399     </edoal:Property>
400   </entity2>
401   <relation>=</relation>
402   <measure rdf:datatype='http://www.w3.org/2001/XMLSchema#float'>1.0</measure>
403 </Cell>
404 </map>

```

4.4 Query processing

The framework allows query execution by defining a query, posed over the global schema, through SPARQL⁵.

The query rewriting process [24] exploits correspondences 1:n between global and local schema elements, expressed in descriptive logic, and applies a set of transformation rules to such correspondences.

Inputs of query rewriting process are a SPARQL query and a mapping table MT (in EDOAL format) and generates a set of queries, also expressed in SPARQL. Subsequently generated queries are transmitted to the acquisition subsystem for their evaluation, that is, to perform the local query processing task.

4.4.1 Global query processing

Query rewriting process is performed by rewriting the graph pattern of a SPARQL query, applying the transformation rules to each triple pattern in it. Since a triple pattern can refer to data (for example,

⁴<http://alignapi.gforge.inria.fr/edoal.html>

⁵<https://www.w3.org/TR/rdf-sparql-query/>

instance relationships) or schema (class and/or property relationships), or both, a pattern subdivision is performed based on the type. A triple pattern is a triple (*subject*, *predicate*, *object*), which can be classified as:

- DTP (Data Triple Pattern): if it is related to information concerning data and not the schema;
- STP (Schema Triple Pattern): if it is related to information concerning data and not the schema.

The reformulation process (Algorithm 1) applies the three-step transformation rules. In the first step, the triple is rewritten by considering the specified mappings for the *predicate* part. In the second step are considered mappings for the *object* part, and finally for the *subject* part. SPARQL variables, constants, and RDF/RDFS/OWL properties, which may appear in the subject, predicate, and object part of a triple, are not rewritten. As a result, they will also appear in the rewritten query.

Algorithm 1 SPARQL rewriting

Input: SPARQL query Q_{in} , mapping table MT

Output: SPARQL query Q_{out}

- 1: $GP_{in} \leftarrow$ graph pattern of Q_{in}
 - 2: $GP_{out} \leftarrow GP_{in}$ after replacing IRIs in *FILTER*, using 1:1 mappings MT
 - 3: $GP_{out} \leftarrow$ TRIPLE PATTERN REWRITING(GP_{out} , MT , *predicate*)
 - 4: $GP_{out} \leftarrow$ TRIPLE PATTERN REWRITING(GP_{out} , MT , *object*)
 - 5: $GP_{out} \leftarrow$ TRIPLE PATTERN REWRITING(GP_{out} , MT , *subject*)
 - 6: $Q_{out} \leftarrow$ new query containing GP_{out}
-

Transformation rules [24] are described by a set of functions of the type:

$$D_y^x(t, \mu) \rightarrow TR \quad (1)$$

$$S_y^x(t, \mu) \rightarrow TR \quad (2)$$

where t is a DTP (in 1) or STP (in 2), μ is the mapping between e_s (source schema entity) and e_t (target schema entity) for the subject, predicate or object part of t , $x \in \{s, p, o\}$ denotes the part of the triple used by the function, $y \in \{c, op, dp, *\}$ represents the type of x (a class, relation, property or any, respectively) and TR represents the transformation rule. A mapping μ is a generic element $MT[i, j] = M(id, E_{Gi}, E_{Ljk}, n, rel)$ of the mapping table MT . Although the mapping table allows managing 1:1, 1:n, n:1 and n:m mappings, the query reformulation process does not consider n:1 and n:m mappings. Functions 1 and 2 are used to rewrite each triple of the input graph pattern. Output of global query processing is a set of queries, posed over the local ontologies, still expressed in SPARQL.

4.4.2 Local query processing

Local query processing is the second activity of the mediation process. Each reformulated query is still expressed in SPARQL and a second reformulation step is required for those data sources that use a language other than SPARQL to retrieve data.

Relational sources that the framework allows to integrate use SQL to express a query. To perform query reformulation, a SPARQL engine is used, which uses query rewriting techniques and inference mechanisms: Quest [22].

Query processing for XML data source is supported by a framework, integrated in the system, that allows query reformulation from SPARQL to XQuery⁶: SPARQL2XQuery [4].

5 CASE STUDY

In order to validate the proposed framework, a case study was conducted using three heterogeneous data sources, designed in different contexts, related to the health domain applications. The first two sources integrated are relational data sources, whose entity-relationship diagrams are shown, respectively, in Figure 6(a) and 7(a), while the third source is semi-structured, specifically it is a spreadsheet.

⁶<https://www.w3.org/TR/xquery-30/>

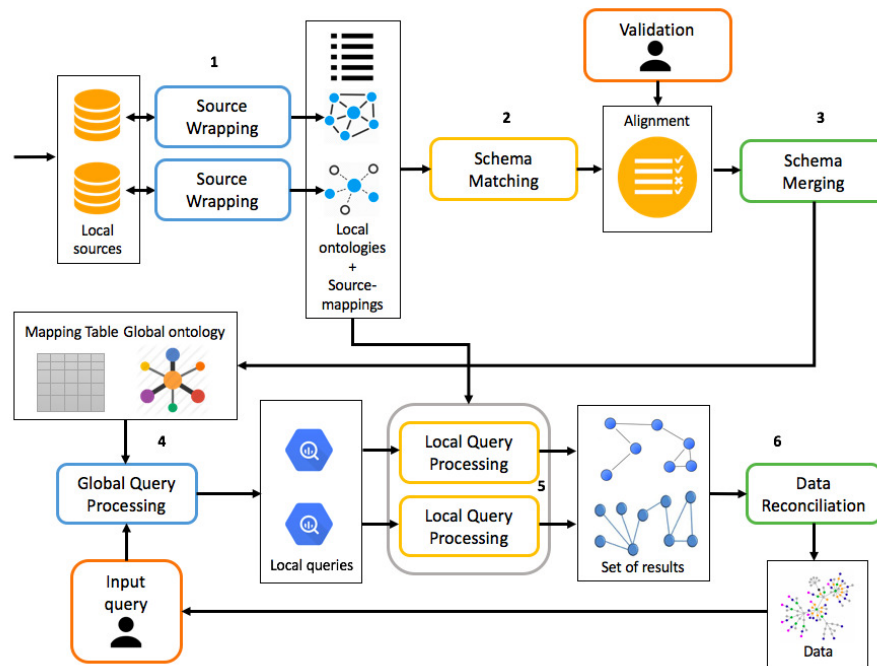


Figure 4. Overview of integration and mediation processes.

As initial step the first source is acquired, and the its scheme becomes the new virtual view (Figure 6(b)). In this case the only steps that must be performed are those of source wrapping and annotation of the schema. The extraction of semantic information, through the schema annotation activity, is necessary as this information will be used to generate the mapping with the source schemes that will be acquired later.

Then, the second source is acquired, and the virtual view is consequently updated. Therefore, the source wrapping and schema annotation activities are performed. The obtained schema is depicted in Figure 7(b).

Subsequently, the schema matching activity is performed. To this aim, the following thresholds setting is adopted:

$$\begin{aligned}\alpha_{\equiv} &= 0.22 \\ \alpha_{\sqsubseteq/\sqsupseteq} &= 0.3 \\ \alpha_{idk} &= 0.8 \\ \beta &= 0.8\end{aligned}$$

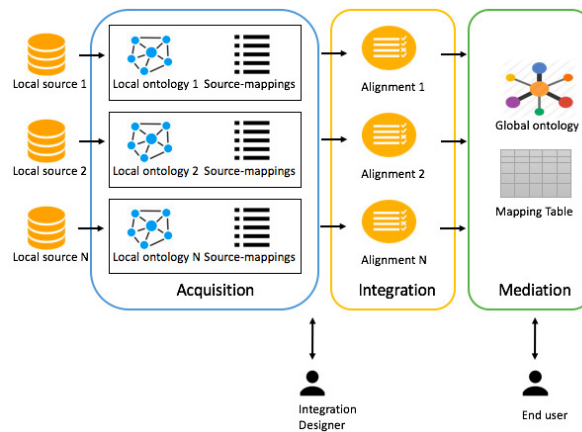
Once the schema matching is completed the mappings are obtained. Some examples are in the following:

$$M(Hospital, Hospital) \Rightarrow$$

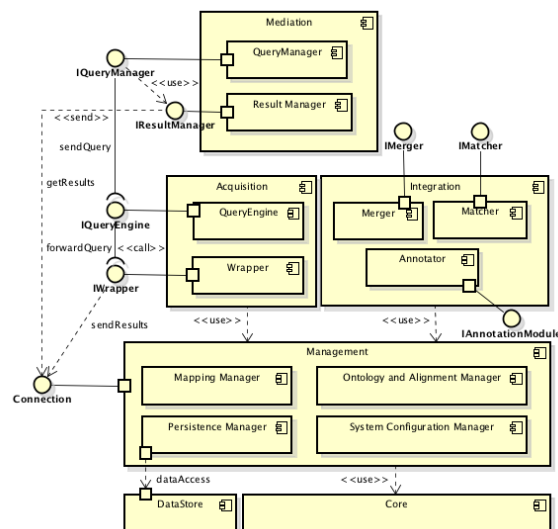
$$\left\{ \begin{array}{lcl} \equiv & & \\ 1 & & \\ Address & \leftrightarrow & Address \quad 1 \\ Name & \leftrightarrow & Name \quad 1 \\ Code & \leftrightarrow & Hosp.Code \quad 0.99 \end{array} \right\}$$

$$M(Operating_Room, Surgery) \Rightarrow$$

$$\left\{ \begin{array}{lcl} \sqsubseteq & & \\ 1 & & \\ Code & \leftrightarrow & Surgery.Code \quad 0.99 \\ NumTables & \leftrightarrow & NumberTables \quad 0.84 \end{array} \right\}$$

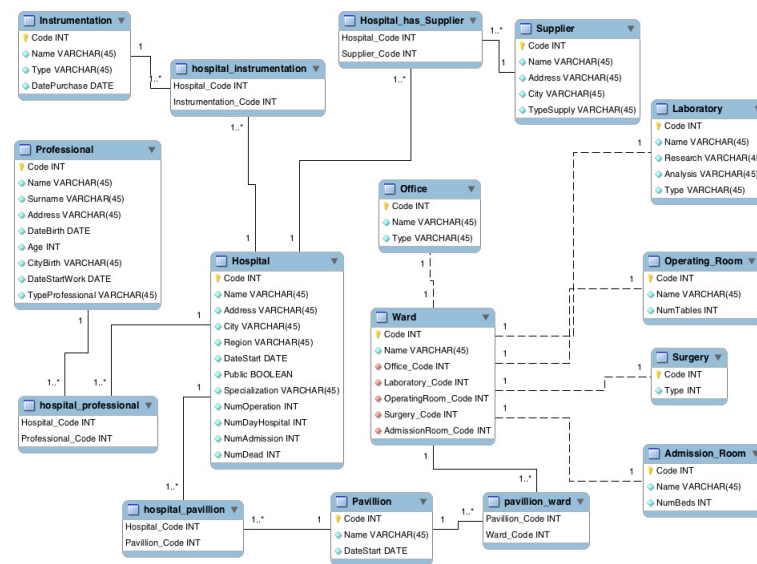


(a) Overview.

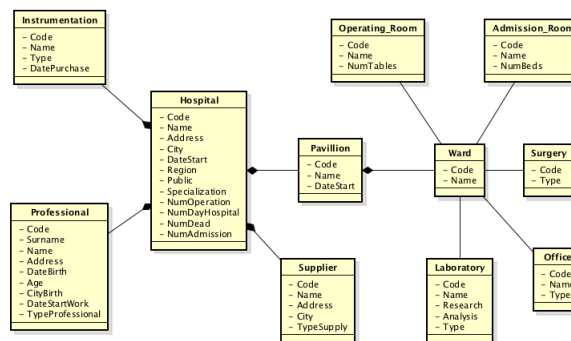


(b) UML diagram.

Figure 5. Framework architecture.

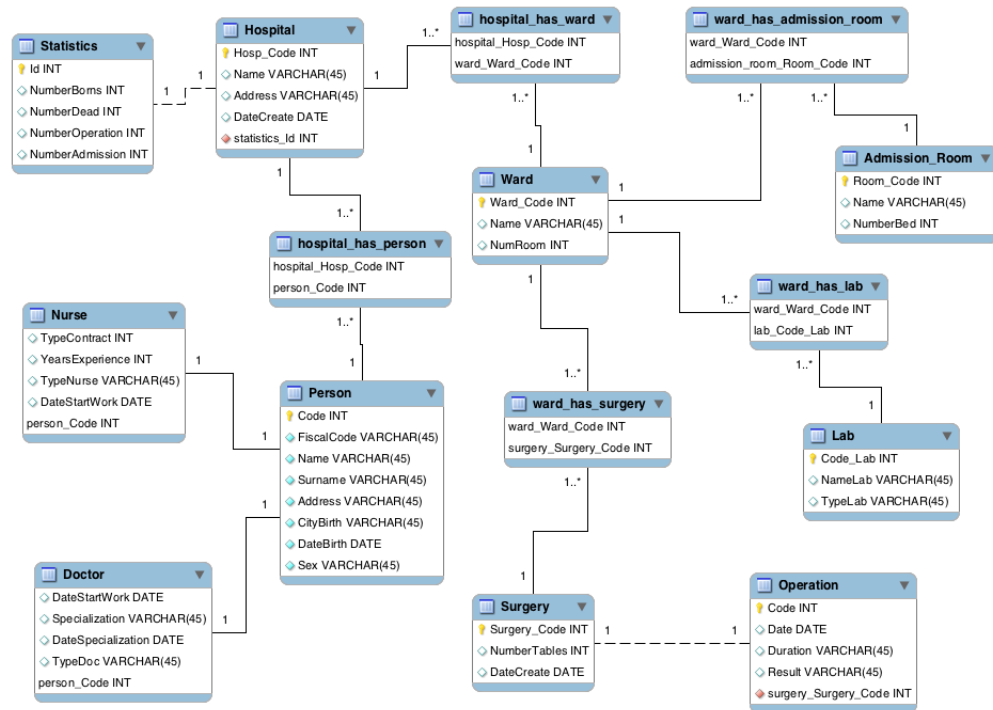


(a) ER diagram

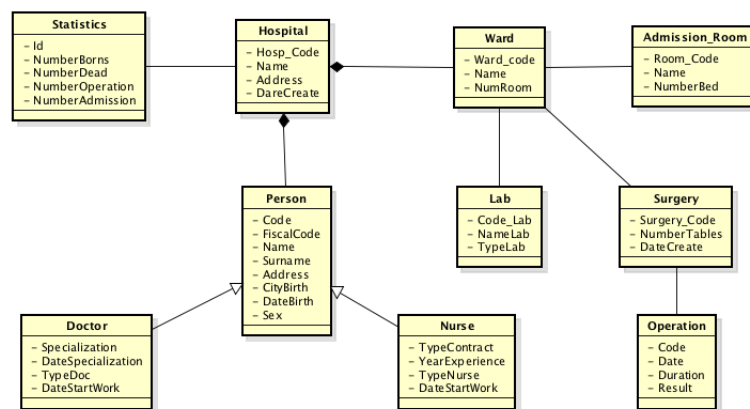


(b) Model of the Local view

Figure 6. First data source



(a) ER diagram



(b) Model of the Local view

Figure 7. Second data source

$$M(\text{Supplier}, \text{Person}) \Rightarrow$$

$$\left\{ \begin{array}{cc} \sqsubseteq & \\ 1 & \\ \text{Name} & \leftrightarrow \text{Name} \quad 1 \\ \text{Address} & \leftrightarrow \text{Address} \quad 1 \\ \text{Code} & \leftrightarrow \text{Code} \quad 1 \\ \text{City} & \leftrightarrow \text{CityBirth} \quad 0.99 \end{array} \right\}$$

$$M(\text{Office}, \text{Lab}) \Rightarrow$$

$$\left\{ \begin{array}{cc} idk & \\ 1 & \\ \text{Type} & \leftrightarrow \text{TypeLab} \quad 0.99 \\ \text{Name} & \leftrightarrow \text{NameLab} \quad 0.99 \\ \text{Code} & \leftrightarrow \text{CodeLab} \quad 0.99 \end{array} \right\}$$

458 During the validation of the mappings, the following changes have been performed:

- 459 • Delete the mapping $M(\text{Office}, \text{Lab})$;
- 460 • The semantic relationship of the mapping $M(\text{OperatingRoom}, \text{Surgery})$ becomes (\equiv) , The seman-
- 461 tic relationship of the mapping becomes, in that the SURGERY class, in the first scheme, refers to a
- 462 room in which a doctor can be consulted, while in the second scheme to an operating room;
- 463 • In mapping $M(\text{Supplier}, \text{Person})$ delete the correspondence between properties *City* and *CityBirth*.

464 The $M(\text{Office}, \text{Lab})$ mapping is returned because the two classes match all properties and, as a

465 result, the contextual similarity measure is 1. This mapping must be deleted otherwise during the schema

466 merging activity a wrong association relationship will be created between the two classes. The α_{idk}

467 threshold was chosen at 0.8 to highlight this observation. If association relationships have no reason to be

468 created, the schema matching activity should be performed with a high value for the α_{idk} threshold.

469 The threshold value $\alpha_{\sqsubseteq/\sqsupset}$ was chosen equal to 0.3 because, if it was lower, the mapping $M(\text{Ward}, \text{Person})$

470 would be added a semantic relationship of hyponymy, but this mapping is wrong. In reality, in the set

471 of mappings should also appear the mapping 1:n $M(\text{Hospital}, \{\text{Hospital}, \text{Statistics}\})$ but is not returned

472 because of the threshold α_{idk} high. To take into account the representation of the hospital concept through

473 the HOSPITAL and STATISTICS classes, there are therefore two alternatives:

- 474 • Keep the threshold value of α_{idk} high and insert the mapping
- 475 $M(\text{Hospital}, \{\text{Hospital}, \text{Statistics}\})$ manually;
- 476 • Lower the threshold and eliminate the other mappings in which there is a *idk* relationship, except
- 477 the mapping $M(\text{Hospital}, \{\text{Hospital}, \text{Statistics}\})$.

However, the following mapping is also obtained:

$$M(\text{Hospital}, \{\text{Hospital}, \text{Statistics}\}) \Rightarrow$$

$$(\text{Hospital}, \text{Hospital}) \Rightarrow$$

$$\left\{ \begin{array}{cc} \equiv & \\ 1 & \\ \text{Address} & \leftrightarrow \text{Address} \quad 1 \\ \text{Name} & \leftrightarrow \text{Name} \quad 1 \\ \text{Code} & \leftrightarrow \text{Hosp_Code} \quad 0.99 \end{array} \right\}$$

$$(\text{Hospital}, \text{Statistics}) \Rightarrow$$

$$\left\{ \begin{array}{cc} idk & \\ 0.21 & \\ \text{NumAdmission} & \leftrightarrow \text{NumberAdmission} \quad 0.88 \\ \text{NumDead} & \leftrightarrow \text{NumberDead} \quad 0.80 \\ \text{NumOperation} & \leftrightarrow \text{NumberOperation} \quad 0.88 \end{array} \right\}$$

478 The global view is initially the same as that of Figure 6(b). At this point the schema merging activity is

479 performed. The new global view is shown in Figure 8

480 As an example, considering the mapping $M(\text{Hospital}, \{\text{Hospital}, \text{Statistics}\})$. Assuming the pre-

481 fix *merged#* to the global view, the prefix *hospital_1#* to the first local view and *hospital_2#* to the

482 second local view, once the schema merging activity is performed, the representation of the mapping

483 $M(\text{Hospital}, \{\text{Hospital}, \text{Statistics}\})$, will be as follows:

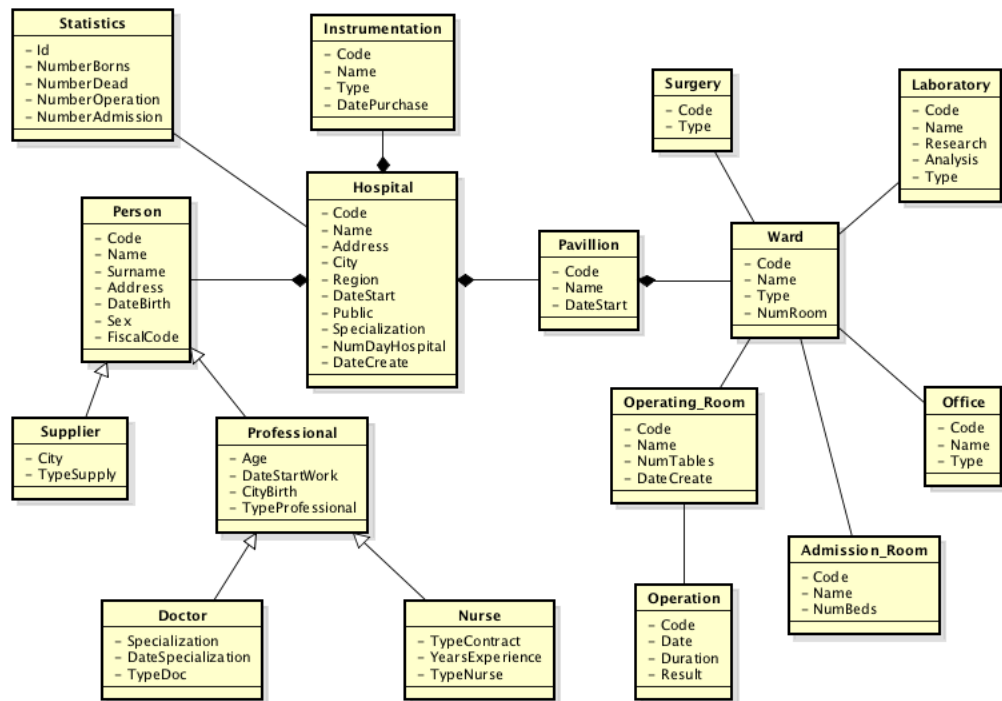


Figure 8. Global view after the second source integration.

```

484 <map>
485   <entity1>
486     <edoal:Class rdf:about="merged#Hospital"/>
487   </entity1>
488   <entity2>
489     <edoal:Class>
490       <edoal:or rdf:parseType="Collection">
491         <edoal:Class rdf:about="hospital_1#Hospital"/>
492         <edoal:Class rdf:about="hospital_2#Hospital"/>
493       </edoal:or>
494     </edoal:Class>
495   </entity2>
496   <relation>=</relation>
497 </map>
498 <map>
499   <entity1>
500     <edoal:Class rdf:about="merged#Statistics"/>
501   </entity1>
502   <entity2>
503     <edoal:Class>
504       <edoal:or rdf:parseType="Collection">
505         <edoal:Class rdf:about="hospital_2#Statistics"/>
506       </edoal:or>
507     </edoal:Class>
508   </entity2>
509   <relation>=</relation>
510 </map>

```

511 In a similar way the correspondences for the other elements of the schemes are defined.

512 The third source acquired is a composed of different spreadsheets.

After source wrapping and schema annotation activities are performed, the schema matching activity is performed using the following threshold values:

$$\begin{aligned}
 \alpha_{\equiv} &= 0.2 \\
 \alpha_{\sqsubseteq/\sqsupseteq} &= 0.3 \\
 \alpha_{idk} &= 0.95 \\
 \beta &= 0.1
 \end{aligned}$$

513 The threshold value of the contextual similarity β is equal to 0.1 because, although there are classes
 514 designed with different attributes, they represent the same concept of the real world and for which,

515 therefore, the mappings must be returned. This situation is managed by lowering the value of β but not
 516 those of $\alpha_{\sqsubseteq/\sqsupseteq}$ and α_{\equiv} . The high value of α_{idk} is meant to filter all *idk* mappings, since they are not
 517 correct.

An example of returned mappings, to be validated by the user, are the following:

$M(Hospital, Rehabilitation.Hospital) \Rightarrow$

$$\left\{ \begin{array}{lcl} & \equiv & \\ & 0.5 & \\ Code & \leftrightarrow & Hospital.Code \quad 0.99 \\ Name & \leftrightarrow & HospitalName \quad 0.99 \\ Address & \leftrightarrow & Street.Address \quad 0.99 \\ City & \leftrightarrow & City \quad 1 \end{array} \right\}$$

$M(Person, Contact.Person) \Rightarrow$

$$\left\{ \begin{array}{lcl} & \sqsupseteq & \\ & 0.5 & \\ Address & \leftrightarrow & Email.Address \quad 0.99 \\ Name & \leftrightarrow & Name \quad 1 \end{array} \right\}$$

$M(Person, Administrator) \Rightarrow$

$$\left\{ \begin{array}{lcl} & \sqsupseteq & \\ & 1 & \\ Address & \leftrightarrow & Email.Address \quad 0.99 \\ Name & \leftrightarrow & Name \quad 1 \end{array} \right\}$$

518 During the validation the changes to perform are the following:

- 519 • Insert a *idk* mapping between the HOSPITAL and
 520 SPECIALIZED_HOSPITAL_CHARACTERISTICS classes.
- 521 • Delete the correspondences between the *Address* and *Email.Address* properties in the mappings
 522 $M(Person, Administrator)$ and $M(Person, Contact.Person)$;
- 523 • Insert a mapping \equiv between the OPERATING_ROOM and SURGERY classes, as they both refer to
 524 an operating room.

525 Once the schema matching activity has been completed, the next step is the schema merging activity.
 526 The new global view, in which all attributes are not shown, is shown in Figure 9.

527 Also in this case, as for the previous data sources, a part of the mapping table representation is
 528 provided. We assign the prefix *merged#* to the global view and prefixes *hospital_1#*, *hospital_2#* and
 529 *hospital_3#*, respectively, to the first, second and third local views. In particular, the extract of the
 530 mapping table refers to the mapping of the *Code* property, related to the HOSPITAL class (the mappings
 531 for the other classes are omitted).

```

532 <map>
533 <Cell>
534 <entity1>
535 <edoal:Property rdf:about="merged#Code"/>
536 </entity1>
537 <entity2>
538 <edoal:Property>
539   <edoal:or rdf:parseType="Collection">
540     <edoal:Property>
541       <edoal:or rdf:parseType="Collection">
542         <edoal:Property>
543           <edoal:and rdf:parseType="Collection">
544             <edoal:Property rdf:about="hospital_2#Hosp_Code"/>
545             <edoal:PropertyDomainRestriction>
546               <edoal:class>
547                 <edoal:Class rdf:about="hospital_2#Hospital"/>
548               </edoal:class>
549             </edoal:PropertyDomainRestriction>
550           </edoal:and>
551         </edoal:Property>
552       </edoal:or>
553     </edoal:Property>
554   <edoal:Property>
555     <edoal:and rdf:parseType="Collection">
556       <edoal:Property rdf:about="hospital_3#Hospital_Code"/>
557       <edoal:PropertyDomainRestriction>

```

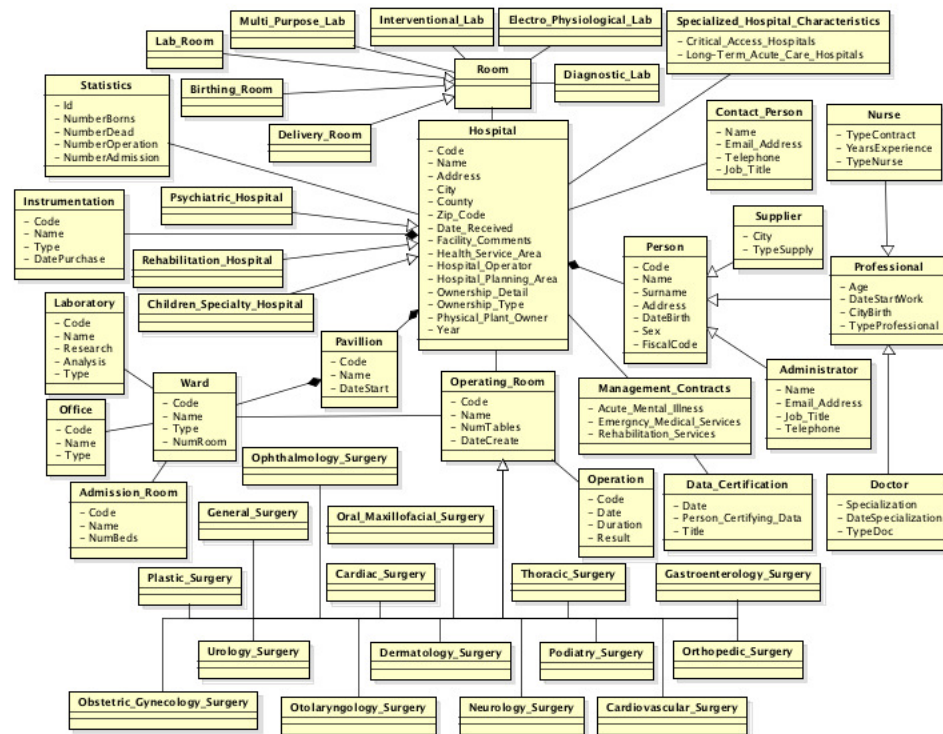


Figure 9. New global view after third data source is integrated.

```

558     <edoal:class>
559         <edoal:Class rdf:about="hospital_3#Hospital"/>
560     </edoal:class>
561 </edoal:PropertyDomainRestriction>
562 </edoal:and>
563 </edoal:Property>
564 <edoal:Property>
565     <edoal:or rdf:parseType="Collection">
566         <edoal:Property>
567             <edoal:and rdf:parseType="Collection">
568                 <edoal:Property rdf:about="hospital_1#Code"/>
569                 <edoal:PropertyDomainRestriction>
570                     <edoal:class>
571                         <edoal:Class rdf:about="hospital_1#Hospital"/>
572                     </edoal:class>
573                 </edoal:PropertyDomainRestriction>
574             </edoal:and>
575         </edoal:Property>
576     </edoal:or>
577 </edoal:Property>
578 </edoal:or>
579 </edoal:Property>
580 </entity2>
581 <relation>=</relation>
582 </Cell>
583 </map>

```

As an example, it is possible to consider the following query: "Return all instances of the HOSPITAL class, with the corresponding names":

```

586 SELECT ?x ?y
587 WHERE {?x rdf:type merged:Hospital. ?x merged:Name ?y}

```

In this case, in the global view are merged *hospital_1#*, *hospital_2#* and *hospital_3#*, respectively, from the first, second and third local views.

The mediation subsystem translates the above query in three queries, one for each of the integrated data sources. The reformulated queries are the followings:

```

592 SELECT ?x ?y
593 WHERE
594 { ?x rdf:type hospital_1#Hospital.

```

```

595         ?x hospital_1#Name ?y ;
596         rdf:type hospital_1#Hospital
597     }

```

```

598 SELECT ?x ?y
599 WHERE
600 { ?x rdf:type hospital_2#Hospital.
601   ?x hospital_2#Name ?y ;
602   rdf:type hospital_2#Hospital
603 }

```

```

604 SELECT ?x ?y
605 WHERE
606 { ?x rdf:type hospital_3#Hospital.
607   ?x hospital_3#HospitalName ?y ;
608   rdf:type hospital_3#Hospital
609 }

```

610 Considering the following query: "Return all instances of the PERSON class, with the corresponding
611 names":

```

612 SELECT ?x ?name
613 WHERE { ?x rdf:type merged:Person. ?x merged:Name ?name}

```

614 The mediation subsystem translates the above query into three queries, one for each of the integrated
615 data sources. The reformulated queries are the followings:

```

616 SELECT ?x ?name
617 WHERE
618 { ?x rdf:type hospital_1#Professional
619   { ?x hospital_1#Name ?name ;
620     rdf:type hospital_1#Supplier}
621   UNION
622   { ?x hospital_1#Name ?name ;
623     rdf:type hospital_1#Professional}
624 }

```

```

625 SELECT ?x ?name
626 WHERE
627 { ?x rdf:type hospital_2#Person.
628   ?x hospital_2#Name ?name ;
629   rdf:type hospital_2#Person
630 }

```

```

631 SELECT ?x ?name
632 WHERE
633 { ?x rdf:type hospital_3#Administrator.
634   ?x hospital_3#Name ?name ;
635   rdf:type hospital_3#Administrator
636 }

```

637 In Figure 10(a) and Figure 10(b), the outputs, in tabular form, of the queries, are shown.

638 6 CONCLUSIONS

639 The purpose of this paper is to allow unified access to heterogeneous and independent source data, offering
640 a data integration approach that addresses all the issues discussed. The architecture adopted is that of
641 mediation systems, which create a virtual view of the real data and allow to external applications to access
642 data through that view in a transparent manner. Transparency is guaranteed by translating queries posed
643 over the virtual view into queries that are directly executable from local sources.

644 The proposed approach allows unified access to heterogeneous sources through the following activities:

- 645 • Source wrapping: the initial activity is the construction of an ontology for each source you want to
646 integrate, whose structure is subsequently refined by using information extraction techniques to
647 improve the quality of ontology.
- 648 • Schema matching: ontologies are then put in a matching process in order to automatically search
649 mappings between the elements of the structures, using both syntax-based and semantic-based
650 techniques. Mappings are identified by combining both semantic and contextual characteristics of
651 each element. These mappings are then validated and, if necessary, modified by the user.
- 652 • Schema merging: based on the generated mappings, a global ontology is created which is the virtual
653 view of the system.

x	y
hospital_1#HospitalCode=125	Methodist Hospital of Chicago
hospital_1#HospitalCode=141	Blessing Hospital at 11th Street
hospital_1#HospitalCode=315	Advocate Christ Medical Center
hospital_1#HospitalCode=331	Clay County Hospital
hospital_1#HospitalCode=364	CGH Medical Center
hospital_1#HospitalCode=414	Community Memorial Hospital
hospital_1#HospitalCode=455	Crawford Memorial Hospital
hospital_1#HospitalCode=471	Decatur Memorial Hospital

(a) Result of the query "Return all instances of the HOSPITAL class, with the corresponding names".

x	name
hospital_2#PersonCode=2	Jim Renneker, RN
hospital_2#PersonCode=3	Heather Kattenbraker
hospital_2#PersonCode=4	Robert J. Stolba
hospital_2#PersonCode=5	Randall Rushing
hospital_2#PersonCode=6	Sherri Hitchcock
hospital_2#PersonCode=7	Mike McManus
hospital_2#PersonCode=8	Robert Ellison

(b) Result of the query "Return all instances of the PERSON class, with the corresponding names".

Figure 10. Tabular form of the output of the queries.

- Query reformulation: at this stage, a query posed over the virtual view is reformulated into a set of queries directly executable by local sources. The reformulation task is performed automatically, generating an execution plan of the reformulated queries, with the possibility for the user to modify each single query.

Overall, the approach is semi-automatic, but compared to existing systems, the user's effort is minimized as he only intervenes in the matching configuration activity, by setting the threshold values for the mappings generation, and mappings validation. Both simple (1:1) and complex mappings (1:n, n:1 and n:m) are generated.

The outlined approach is supported by a specially designed and developed software system. The system provides a first level of abstraction of the activities and components involved in their execution and a second level of component specialization. Although the design of the system is aimed at covering all aspects of data integration described so far, implementation has some limitations. In particular, the acquisition of unstructured sources is not yet contemplated in development and the data reconciliation process requires the development of appropriate components. Except for such activities, integration and mediation processes are fully supported by the system.

Research activities that will be carried out in the future will have the goal of overcoming the limitations shown and consolidating, at the same time, the part of the system developed so far. In particular, accurate experimentation is required for validating the proposed approach, for ensuring high quality of mappings and local and global views, for optimizing the mediation process.

REFERENCES

- [1] Yigal Arens, Chun-Nan Hsu, and Craig A Knoblock. "Query processing in the sims information mediator". In: *Advanced Planning Technology* 32 (1996), pp. 78–93.
- [2] Yigal Arens et al. "Retrieving and integrating data from multiple information sources". In: *International Journal of Intelligent and Cooperative Information Systems* 2.02 (1993), pp. 127–158.
- [3] Domenico Beneventano and Sonia Bergamaschi. "The MOMIS methodology for integrating heterogeneous data sources". In: *Building the Information Society*. Springer, 2004, pp. 19–24.
- [4] Nikos Bikakis et al. "Querying xml data with sparql". In: *International Conference on Database and Expert Systems Applications*. Springer. 2009, pp. 372–381.

- 683 [5] Diego Calvanese, Domenico Lembo, and Maurizio Lenzerini. “Survey on methods for query
684 rewriting and query answering using views”. In: *Integrazione, Warehousing e Mining di sorgenti*
685 *eterogenee* 25 (2001).
- 686 [6] Sudarshan Chawathe et al. “The TSIMMIS project: Integration of heterogenous information
687 sources”. In: (1994).
- 688 [7] Laura Chiticariu, Phokion G Kolaitis, and Lucian Popa. “Interactive generation of integrated
689 schemas”. In: *Proceedings of the 2008 ACM SIGMOD international conference on Management of*
690 *data*. ACM. 2008, pp. 833–846.
- 691 [8] Joseph Fong et al. “Schema integration for object-relational databases with data verification”. In:
692 *Proceedings of the 2000 international computer symposium workshop on software engineering and*
693 *database systems, Taiwan*. 2000, pp. 185–192.
- 694 [9] Raji Ghawi and Nadine Cullot. “Building ontologies from XML data sources”. In: *Database*
695 *and Expert Systems Application, 2009. DEXA’09. 20th International Workshop on*. IEEE. 2009,
696 pp. 480–484.
- 697 [10] Raji Ghawi and Nadine Cullot. “Database-to-ontology mapping generation for semantic interoper-
698 ability”. In: *VDBL’07 conference, VLDB Endowment ACM*. 2007, pp. 1–8.
- 699 [11] Cheng Hian Goh et al. “Context interchange: New features and formalisms for the intelligent
700 integration of information”. In: *ACM Transactions on Information Systems (TOIS)* 17.3 (1999),
701 pp. 270–293.
- 702 [12] R Hull, R King, et al. “Arpa i3 reference architecture, 1995”. In: *Reperibile presso: [http://www.](http://www.isse.gmu.edu/I3_Arch/index.html)*
703 *isse.gmu.edu/I3_Arch/index.html* (1995).
- 704 [13] Yannis Katsis and Yannis Papakonstantinou. “View-based data integration”. In: *Encyclopedia of*
705 *Database Systems*. Springer, 2009, pp. 3332–3339.
- 706 [14] Hyunjang Kong, Myungwon Hwang, and Pankoo Kim. “A new methodology for merging the
707 heterogeneous domain ontologies based on the wordnet”. In: *Next Generation Web Services*
708 *Practices, 2005. NWeSP 2005. International Conference on*. IEEE. 2005, 6–pp.
- 709 [15] MongLi Lee and TokWang Ling. “A methodology for structural conflict resolution in the integration
710 of entity-relationship schemas”. In: *Knowledge and Information Systems* 5.2 (2003), pp. 225–247.
- 711 [16] Maurizio Lenzerini. “Data Integration: A Theoretical Perspective”. In: *Proceedings of the Twenty-*
712 *first ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. PODS ’02.
713 Madison, Wisconsin: ACM, 2002, pp. 233–246. ISBN: 1-58113-507-6. DOI: 10.1145/543613.
714 543644. URL: <http://doi.acm.org/10.1145/543613.543644>.
- 715 [17] Eduardo Mena et al. “Managing multiple information sources through ontologies: relationship
716 between vocabulary heterogeneity and loss of information”. In: (1996).
- 717 [18] Eduardo Mena et al. “OBSERVER: An approach for query processing in global information
718 systems based on interoperation across pre-existing ontologies”. In: *Cooperative Information*
719 *Systems, 1996. Proceedings., First IFCIS International Conference on*. IEEE. 1996, pp. 14–25.
- 720 [19] M Orsini et al. “Query Management in Data Integration Systems: the MOMIS approach”. PhD
721 thesis. PhD thesis, International Doctorate School in Information, Communication Technologies of
722 the University of Modena, and Reggio Emilia, 2009.
- 723 [20] Alun Preece, Kit Hui, and Peter Gray. “KRAFT: Supporting virtual organisations through knowl-
724 edge fusion”. In: *Artificial Intelligence for Electronic Commerce: Papers from the AAAI-99 Work-*
725 *shop*. 1999, pp. 33–38.
- 726 [21] Alun Preece et al. “The KRAFT architecture for knowledge fusion and transformation”. In:
727 *Knowledge-Based Systems* 13.2 (2000), pp. 113–120.
- 728 [22] Mariano Rodriguez-Muro, Josef Hardi, and Diego Calvanese. “Quest: efficient SPARQL-to-SQL
729 for RDF and OWL”. In: *Proceedings of the 2012th International Conference on Posters &*
730 *Demonstrations Track-Volume 914*. CEUR-WS. org. 2012, pp. 53–56.
- 731 [23] Pavel Shvaiko and Jérôme Euzenat. “A survey of schema-based matching approaches”. In: *Journal*
732 *on data semantics IV*. Springer, 2005, pp. 146–171.

- 733 [24] Élodie Thiéblin et al. “Rewriting SELECT SPARQL queries from 1: n complex correspondences”.
734 In: *Ontology Matching* (2016), p. 49.
- 735 [25] Holger Wache et al. “Ontology-based integration of information-a survey of existing approaches”.
736 In: *IJCAI-01 workshop: ontologies and information sharing*. Vol. 2001. Citeseer. 2001, pp. 108–
737 117.
- 738 [26] Li Xu and David W Embley. “Combining the Best of Global-as-View and Local-as-View for Data
739 Integration.” In: *ISTA*. Vol. 48. 2004, pp. 123–136.