

# Mediators over Taxonomy-based Information Sources

Yannis Tzitzikas<sup>1\*</sup>, Nicolas Spyratos<sup>2</sup>, Panos Constantopoulos<sup>1</sup>

<sup>1</sup> Department of Computer Science, University of Crete, Greece and Institute of Computer Science, ICS-FORTH e-mail: {tzitzik|panos}@csi.forth.gr

<sup>2</sup> Laboratoire de Recherche en Informatique, Universite de Paris-Sud, France e-mail: spyratos@lri.fr

The date of receipt and acceptance will be inserted by the editor

**Abstract** We propose a mediator model for providing integrated and unified access to multiple taxonomy-based sources. Each source comprises a taxonomy and a database that indexes objects under the terms of the taxonomy. A mediator comprises a taxonomy and a set of relations between mediator's and sources' terms, called articulations. By combining different modes of query evaluation at the sources and the mediator, and different types of query translation, a flexible, efficient scheme of mediator operation is obtained, which can accommodate various application needs and levels of answer quality. We adopt a simple conceptual modeling approach (taxonomies and inter-taxonomy mappings) and we illustrate its advantages in terms of ease of use, uniformity, scalability and efficiency. These characteristics make this proposal appropriate for a large-scale network of sources and mediators.

**Key words** Mediators – Taxonomies – Approximate Query Translation – Information Integration

## 1 Introduction

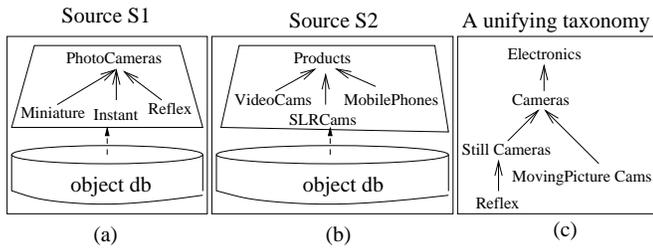
The need for integrated and unified access to multiple information sources has stimulated the research on *mediators* (initially proposed in [78]). Roughly, a mediator is a secondary information source aiming at providing a uniform interface to a number of underlying sources (which may be primary or secondary). Users submit queries to the mediator. Upon receiving a user query, the mediator queries the underlying sources. This involves selecting the sources to be queried and formulating the query to be sent to each source. These tasks are accomplished based on what the mediator “knows” about the underlying sources. Finally, the mediator appropriately combines the returned results and delivers the final answer to the user.

In this paper we consider information sources over a common domain consisting of a denumerable set of objects. For example, in the environment of the Web, the domain could be the set of all Web pages, specifically, the set of all pointers to Web pages. Each source has a *taxonomy*, i.e. a structured set of names, or *terms*, that are familiar to the users of the source. In particular, the taxonomies considered in this paper consist of a set of terms structured by a subsumption relation. In addition, each source maintains a database storing objects that are of interest to its users. Specifically, each object in the database of a source is indexed under one or more terms of the taxonomy of that source. In quest for objects of interest, a user can browse the source taxonomy until he reaches the desired terms, or he can query the source by submitting a boolean expression of terms. The source will then return the appropriate set of objects. In the environment of the Web, general purpose catalogues, such as Yahoo! or Open Directory<sup>1</sup>, domain-specific catalogues/gateways (e.g. for medicine, physics, tourism), as well as personal bookmarks of Web browsers can be considered as examples of such sources.

However, although several sources may carry information about the same domain, they usually employ different taxonomies, with terms that correspond to different natural languages, or different levels of granularity. For example, consider two sources  $S_1$  and  $S_2$  that both provide access to electronic products as shown in figures 1.(a) and 1.(b). Each source consists of a taxonomy plus a database that indexes objects under the terms of that taxonomy. However, the two sources provide *different* information about electronic products - as seen in the figures. Suppose now that we want to provide unified access to these two sources through a single taxonomy which is familiar to a specific group of users. An example of such a unifying taxonomy is shown in Figure 1.(c), and constitutes part of what we call a “mediator”.

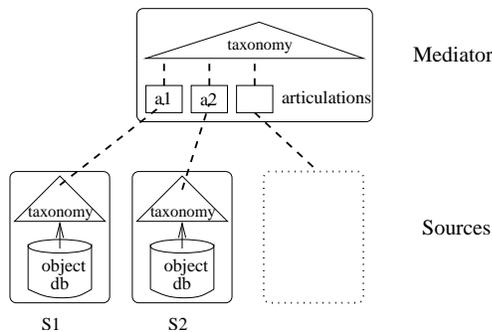
\* Current address: Istituto di Scienza e Tecnologie dell' Informazione, ISTI-CNR, Pisa, Italy

<sup>1</sup> <http://dmoz.org>



**Fig. 1** Two sources providing access to electronic products

A *mediator* is a secondary source that can bridge the heterogeneities that may exist between two or more sources in order to provide unified access to those sources. Specifically, a mediator has a taxonomy and structuring that reflects the needs of its potential users, but does *not* maintain a database of objects. Instead, the mediator maintains a number of *articulations* to the sources. An articulation to a source is a set of relationships between the terms of the mediator and the terms of that source. These relationships are defined by the designer of the mediator at design time and are stored at the mediator. Figure 2 shows the general architecture of a mediator.



**Fig. 2** The mediator architecture

Users formulate queries over the taxonomy of the mediator and it is the task of the mediator to choose the sources to be queried, and to formulate the query to be sent to each source. To this end, the mediator uses the articulations in order to translate queries over its own taxonomy to queries over the taxonomies of the articulated sources. Then it is again the task of the mediator to combine the results returned by the sources appropriately in order to produce the final answer.

An essential feature that distinguishes our work is that we adopt a simple conceptual modeling approach for both sources and mediators. This conceptual modeling approach has the following advantages: (a) it is very easy to create the conceptual model of a source or a mediator, and (b) the integration of information from multiple sources can be done very easily. Indeed, as we

shall see, the articulations offer a *uniform* and easy to use method to bridge naming, contextual and granularity heterogeneities between the conceptual models of the sources. Given this conceptual modeling approach, the mediator does not have to tackle complex structural differences between the sources (as it happens in mediators for relational databases).

Another essential feature that distinguishes our approach is that a source can provide two types of answer to a given query, namely, a *sure* answer or a *possible* answer. The first type of answer is appropriate for a user who does not want to retrieve objects which are not relevant to his information need, while the second for a user who does not want to miss objects which are relevant to his information need. Moreover, as exact translation of user queries is not always possible, a user query to the mediator admits two types of approximation, namely, *lower* or *upper* translation. What kind of translation will be used at the mediator level and what kind of answer will be requested at the source level is decided by the mediator designer at design time and/or the mediator user at query time. Therefore a prominent feature of our approach is that sources and mediators can operate in a variety of modes according to specific application needs. As a consequence, our mediators are quite flexible and can adapt to a variety of situations.

A main objective of this paper is to prescribe easy to use and formally sound methods for building mediators. In the context of the Web, our mediators can be used for providing unified access to multiple Web catalogues. An advantage of our approach is that a mediator can be constructed quite easily, therefore ordinary Web users can use it in order to define their own mediators. In this sense, this approach can be used for personalizing existing Web catalogues. Furthermore it can be used for building mediators over XFML [1] information bases (XFML aims at applying the faceted classification paradigm in the context of the Web).

The remainder of the paper is organized as follows: Section 2 describes the information sources and the query answering process at a single source. Section 3 defines the architecture of a mediator over a set of sources and the different modes in which a mediator can operate. Section 4 discusses query evaluation and Section 5 discusses enhancements of the query answering process. Section 6 discusses various extensions of our model. Section 7 discusses related work and, finally, Section 8 concludes the paper and discusses further research. All proofs are given in the Appendix.

## 2 The Sources

### Why taxonomies

Taxonomies is probably the oldest and most widely used conceptual modeling tool. Nevertheless, it is a pow-

erful tool still used is in Web directories (e.g. in Google and Yahoo!), Content Management (hierarchical structures are used to classify documents), Web Publishing (many authoring tools require to organize the contents of portals according to some hierarchical structure), Web Services (services are typically classified in a hierarchical form), Marketplaces (goods are classified in hierarchical catalogs), Personal File Systems, Personal Bookmarks for the Web, Libraries (e.g. Thesauri [40]) and in very large collections of objects (e.g. see [61]). Although more sophisticated conceptual models (including concepts, attributes, relations and axioms) have emerged and are recently employed even for meta-tagging in the Web [49, 75], almost all of them have a backbone consisting of a subsumption hierarchy, i.e. a taxonomy.

Furthermore, a taxonomy-based conceptual modeling approach has several advantages in large and open domains. In a very broad domain, such as the set of all Web pages, it is not easy to identify the classes of the domain because the domain is too wide and different users, or applications, conceptualize it differently, e.g. one class of the conceptual model according to one user may correspond to a value of an attribute of a class of the conceptual model according to another user. For example, Figure 3 shows two different conceptual models for the same domain. We consider only two objects of the domain, denoted by the natural numbers 1 and 2.

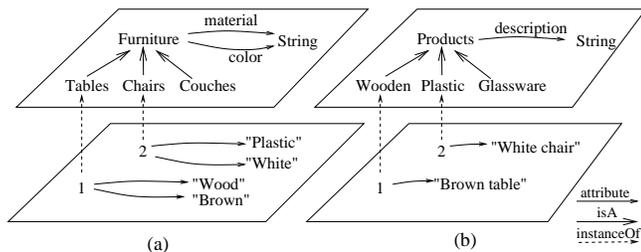


Fig. 3 Two different conceptual models for the same domain

The conceptual model of Figure 3.(a) is appropriate for building an information system for a furniture store, while the conceptual model of Figure 3.(b) is appropriate for building an information system for a department store. The classes of model (a), i.e. the classes **Tables**, **Chairs** and **Couches**, have been defined so as to distinguish the objects of the domain according to their *use*. On the other hand, the classes of model (b), i.e. the classes **Wooden**, **Plastic** and **Glassware**, have been defined so as to distinguish the objects of the domain according to their *material*. This kind of distinction is useful for a department store, as it determines (up to some degree) the placement of the objects in the various departments of the store. Figure 4 shows a conceptual model for the same domain which consists of terms and subsumption links only, i.e. a *taxonomy*. This conceptual modeling approach seems to be more application inde-

pendent. All criteria (characteristics) for distinguishing the objects are equally “honoured”.

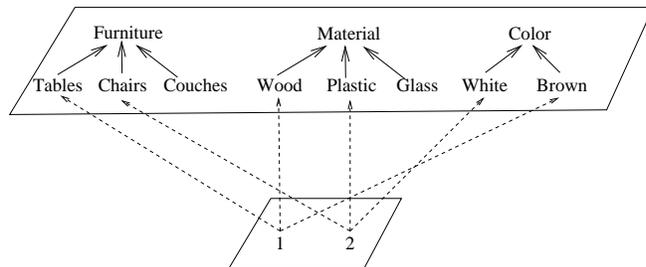


Fig. 4 A conceptual model that consists of terms and subsumption links only

A simple conceptual modeling approach where each conceptual model is a taxonomy, has three main advantages. The first is that it is very easy to create the conceptual model of a source or a mediator. Even ordinary Web users can design this kind of conceptual models. Besides, the queries submitted by ordinary users are mostly bags of words, not structured queries. Furthermore, the design can be done more systematically if done following a faceted approach (e.g. see [58, 57]). In addition, thanks to techniques that have emerged recently [68], taxonomies of compound terms can be also defined in a flexible and systematic manner.

The second, is that the simplicity and modeling uniformity of taxonomies allows integrating the contents of several sources without having to tackle complex structural differences. Indeed, as it will be seen in the subsequent sections inter-taxonomy mappings offer a *uniform* method to bridge *naming*, *contextual* and *granularity* heterogeneities between the taxonomies of the sources. Given this conceptual modeling approach, a mediator does not have to tackle complex structural differences between the sources, as happens with relational mediators (e.g. see [37, 47]) and Description Logics-based mediators (e.g. see [42, 13]). Moreover, it allows the integration of *schema* and *data* in a uniform manner. Another advantage of this conceptual modeling approach is that query evaluation in taxonomy-based sources and mediators can be done efficiently (polynomial time).

The third, is that this conceptual modeling approach makes the automatic construction of mappings possible [69]. The last is the major drawback of the current more expressive Web annotation languages.

Due to the above benefits (conceptual modeling simplicity, integration flexibility, query evaluation efficiently), taxonomies worth further investigation. The only assumption that we make is that the domain is a set of objects which we want to index and subsequently re-

trieve, without being interested in the relationships that may hold between the objects of the domain.

### Defining a Source

Let  $Obj$  denote the set of all objects of a domain common to several information sources. A typical example of such a domain is the set of all pointers to Web pages. We assume that each source has a *taxonomy*, defined as follows:

**Definition 1** A taxonomy is a pair  $(T, \preceq)$  where  $T$  is a *terminology*, i.e. a set of names, or *terms*, and  $\preceq$  is a *subsumption* relation over  $T$ , which is a reflexive and transitive relation over  $T$ .

If  $a$  and  $b$  are terms of  $T$ , we say that  $a$  is *subsumed* by  $b$  if  $a \preceq b$ ; we also say that  $b$  *subsumes*  $a$ ; for example,  $Databases \preceq Informatics$ ,  $Canaries \preceq Birds$ . We say that two terms  $a$  and  $b$  are *equivalent*, and write  $a \sim b$ , if both  $a \preceq b$  and  $b \preceq a$  hold, e.g.,  $Computer\ Science \sim Informatics$ . Note that the subsumption relation is a preorder over  $T$  and that  $\sim$  is an equivalence relation over the terms of  $T$ . Moreover  $\preceq$  is a partial order over the equivalence classes of terms.

We assume that, in addition to its taxonomy, each source has a stored *interpretation*  $I$  of its terminology, i.e. a function  $I : T \rightarrow 2^{Obj}$  that associates each term of  $T$  with a set of objects. Here, we use the symbol  $2^{Obj}$  to denote the powerset of  $Obj$ . Figure 5 shows an example of a source.

In this and subsequent figures the objects are represented by natural numbers and membership of objects to the interpretation of a term is indicated by a dotted arrow from the object to that term. For example, the objects 1 and 3 in Figure 5 are members of the interpretation of the term `JournalArticle` as these objects are connected to `JournalArticle` with dotted arrows. Moreover, as these are the only objects connected to `JournalArticle` with dotted arrows, they make up the interpretation of `JournalArticle`, i.e.

$$I(\text{JournalArticle}) = \{1, 3\}.$$

Subsumption of terms is indicated by a continuous-line arrow from the subsumed term to the subsuming term. For example, the term `RDB` in Figure 5 is subsumed by `DB` as there is a continuous-line arrow going from `RDB` to `DB`; this arrow indicates that  $RDB \preceq DB$ .

Note that we do not represent the entire subsumption relation but a subset of it sufficient to generate the entire relation. In particular, we do not represent the reflexive, nor the transitive arrows of the subsumption relation.

Equivalence of terms is indicated by a continuous non-oriented line connecting the terms that are equivalent. For example, the term `Databases` is equivalent with the term `DB` since these two terms are connected by a continuous non-oriented line. Note that equivalence captures the notion of synonymy, and that each equiva-

lence class simply contains alternative terms for naming a given set of objects.

For technical reasons that will become clear shortly, we assume that every terminology  $T$  contains two special terms, the *top term*, denoted by  $\top$ , and the *bottom term*, denoted by  $\perp$ . The top term subsumes every other term  $t$ , i.e.  $t \preceq \top$ . The bottom term is strictly subsumed by every other term  $t$  different than top and bottom, i.e.  $\perp \preceq \perp$ ,  $\perp \preceq \top$ , and  $\perp \prec t$ , for every  $t$  such that  $t \neq \top$  and  $t \neq \perp$ . Moreover we assume that every interpretation  $I$  of  $T$  satisfies the condition  $I(\perp) = \emptyset$ .

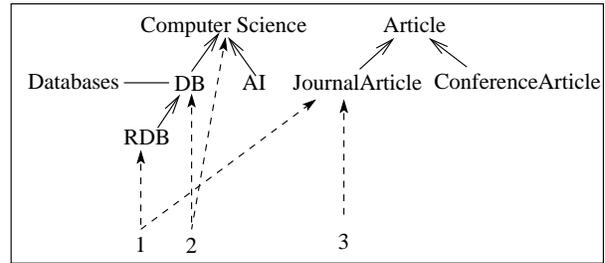


Fig. 5 Graphical representation of a source

### Querying a source

Each source responds to queries over its own terminology. A query is either a term or a combination of terms using the usual connectives  $\wedge$ ,  $\vee$ ,  $\neg$  and  $()$ . For technical reasons that will become clear shortly we shall also use the concept of *empty query* denoted by  $\epsilon$ . More formally, a query is defined as follows:

**Definition 2** Let  $T$  be a terminology. A *query* over  $T$  is any string derived by the following grammar, where  $t$  is a term of  $T$ :

$$q ::= t \mid q \wedge q' \mid q \vee q' \mid q \wedge \neg q' \mid (q) \mid \epsilon$$

Note that our use of negation corresponds to domain restricted negation.

In what follows, given a query  $q$  we define two answers of  $q$ , that we call the *sure* and *possible* answer. To this end, we need some preliminary definitions and notations.

The set of interpretations of a given terminology  $T$  can be ordered using pointwise set inclusion.

**Definition 3** Given two interpretations  $I, I'$  of  $T$ , we call  $I$  less than or equal to  $I'$ , and we write  $I \sqsubseteq I'$ , if  $I(t) \subseteq I'(t)$  for each term  $t \in T$ .

Note that  $\sqsubseteq$  is a partial order over interpretations.

A source answers queries based on the stored interpretation of its terminology. However, in order for query answering to make sense, the interpretation that a source

uses for answering queries must respect the structure of the source’s taxonomy (i.e. the relation  $\preceq$ ) in the following sense: if  $t \preceq t'$  then  $I(t) \subseteq I(t')$ . For example, consider a source whose taxonomy contains only three terms: `DB`, `AI` and `Computer Science`, where  $\text{DB} \preceq \text{Computer Science}$ , and  $\text{AI} \preceq \text{Computer Science}$ . Assume that in the stored interpretation  $I$  of the source we have:  $I(\text{DB}) \neq \emptyset$ ,  $I(\text{AI}) \neq \emptyset$  and  $I(\text{ComputerScience}) = \emptyset$ . Clearly,  $I$  does not respect the structure of the taxonomy, as  $\text{DB} \preceq \text{Computer Science}$ , and yet  $I(\text{DB}) \not\subseteq I(\text{ComputerScience})$ . However,  $I$  is acceptable as we can “augment” it to a new interpretation  $I'$  that *does* respect the structure of the taxonomy. The interpretation  $I'$  is defined as follows:  $I'(\text{DB}) = I(\text{DB})$ ,  $I'(\text{AI}) = I(\text{AI})$ ,  $I'(\text{ComputerScience}) = I(\text{ComputerScience}) \cup I(\text{DB}) \cup I(\text{AI})$ . An interpretation such as  $I'$  that respects the structure of a taxonomy is what we call a model of that taxonomy.

**Definition 4** An interpretation  $I$  is a *model* of a taxonomy  $(T, \preceq)$  if for all  $t, t'$  in  $T$ , if  $t \preceq t'$  then  $I(t) \subseteq I(t')$ .

For brevity hereafter we shall sometimes write  $T$  instead of  $(T, \preceq)$ , whenever no confusion is possible.

Now, as there may be several models of  $T$  in general, we assume that each source answers queries from one or more *designated* models induced by its stored interpretation. In this paper we will use two specific models for answering queries, the *sure model* and the *possible model*. In order to define these models formally we need to introduce the notions of *tail* and *head* of a term.

**Definition 5** Given a term  $t \in T$  we define

$$\text{tail}(t) = \{s \in T \mid s \preceq t\} \text{ and } \text{head}(t) = \{u \in T \mid t \preceq u\}$$

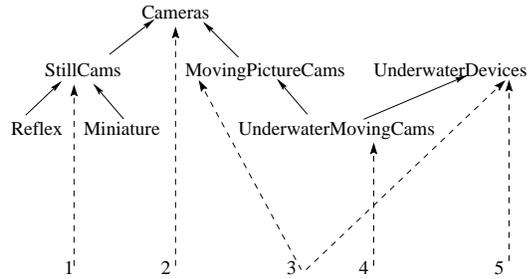
Note that  $t$ , and all terms that are equivalent to  $t$ , belong to both  $\text{tail}(t)$  and  $\text{head}(t)$ . Also note that  $\text{tail}(t)$  always contains the bottom term  $\perp$  and  $\text{head}(t)$  always contains the top term  $\top$ .

**Definition 6** Given an interpretation  $I$  of  $T$  we define the *sure model* of  $T$  generated by  $I$ , denoted  $I^-$ , as follows:

$$I^-(t) = \bigcup \{I(s) \mid s \in \text{tail}(t)\}$$

Intuitively the stored set  $I(t)$  consists of the objects that are known to be indexed under  $t$ . The set  $I^-(t)$  on the other hand consists of the objects known to be indexed under  $t$  plus the objects that are known to be indexed under terms subsumed by  $t$ . Therefore  $I^-(t)$  consists of all objects that are *surely* indexed under  $t$  with respect to  $I$  and  $\preceq$ . Figure 6 shows an example of a source and its sure model  $I^-$ .

**Proposition 1** If  $I$  is an interpretation of  $T$  then  $I^-$  is the unique minimal model of  $T$  which is greater than or equal to  $I$ .



(a)

Term	$I$	$I^-$	$I^+$
$\perp$	$\emptyset$	$\emptyset$	$\emptyset$
$\top$	$\emptyset$	$\{1,2,3,4,5\}$	$\{1,2,3,4,5\}$
Cameras	$\{2\}$	$\{1,2,3,4\}$	$\{1,2,3,4,5\}$
StillCams	$\{1\}$	$\{1\}$	$\{1,2,3,4\}$
Reflex	$\emptyset$	$\emptyset$	$\{1\}$
Miniature	$\emptyset$	$\emptyset$	$\{1\}$
MovingPictureCams	$\{3\}$	$\{3,4\}$	$\{1,2,3,4\}$
UnderwaterMovingCams	$\{4\}$	$\{4\}$	$\{3,4\}$
UnderwaterDevices	$\{3,5\}$	$\{3,4,5\}$	$\{1,2,3,4,5\}$

(b)

**Fig. 6** Graphical representation of a source

**Definition 7** Given a taxonomy  $T$  and interpretation  $I$  we define the *possible model* of  $T$  generated by  $I$ , denoted  $I^+$ , as follows:

$$I^+(t) = \bigcap \{I^-(u) \mid u \in \text{head}(t) \text{ and } u \not\preceq t\}$$

As it is clear from its definition, the set  $I^+(t)$  consists of the objects known to be indexed under each term strictly subsuming  $t$ . Therefore  $I^+(t)$  consists of all objects that are *possibly* indexed under  $t$  with respect to  $I$  and  $\preceq$ . An example of the possible model of a source is given in Figure 6. In this example we have:

$$I^+(\text{Reflex}) = \{1\}$$

$$I^+(\text{UnderwaterMovingCams}) = \{3,4\}$$

Note that the possible interpretations of the terms `Cameras` and `UnderwaterDevices` is the set of *all* stored objects. This is so because the head of each of these terms contains only the term itself and the top term  $\top$ , thus we have:

$$\begin{aligned} I^+(\text{Cameras}) &= I^+(\text{UnderwaterDevices}) = I^-(\top) = \\ &= \bigcup \{I(s) \mid s \preceq \top\} \end{aligned}$$

Note that since  $\text{head}(\top) = \{\top\}$ , the set  $\{u \in \text{head}(\top) \mid u \not\preceq \top\}$  is empty. This means that  $I^+(\top)$ , i.e.  $\bigcap \{I^-(u) \mid u \in \text{head}(\top) \text{ and } u \not\preceq \top\}$  is actually the intersection of an empty family of subsets of  $\text{Obj}$ . However, according to the Zermelo axioms of set theory<sup>2</sup> (see [10] for an overview), the intersection of an empty family of subsets of a universe equals to the universe. In our case, the universe is the set of all objects known to the source, i.e. the set  $I^-(\top)$ , thus we conclude that  $I^+(\top) = I^-(\top)$ .

<sup>2</sup> We do not mean here the Zermelo-Fraenkel axioms.

**Proposition 2** If  $I$  is an interpretation of  $T$  then  $I^+$  is a model of  $T$  and  $I \sqsubseteq I^- \sqsubseteq I^+$ .

It follows from the above proposition that for every term  $t$  we have  $I^-(t) \subseteq I^+(t)$ .

We view the stored interpretation  $I$  as the result of indexing. However, although we may assume that indexing is done correctly, certain objects may not have been indexed under all terms that could apply to them. For example, object 1 in Figure 6 is indexed under `StillCams` but not under `Cameras`, and object 3 is indexed under `MovingPictureCams` and `UnderwaterDevices` but not under `UnderwaterMovingCams`. Note that object 3 could in fact be an `UnderwaterMovingCamera` but it was not indexed under this term because either the indexer did not use this term, or the term `UnderwaterMovingCamera` was defined after the indexing of object 3 was performed.

By consequence, given a query that consists of a single term  $t$ , we may want to answer it in either of two ways: (a) by including in the answer only objects that are known to be indexed under  $t$ , or (b) by including in the answer objects that are possibly indexed under  $t$ . In the first case the answer is the set  $I^-(t)$ , while in the second it is the set  $I^+(t)$ .

*Remark.* If we consider that each term corresponds to a property or characteristic of the objects of the domain, then  $t \preceq t'$  means that if an object has the property  $t$  then it also has the property  $t'$ . In this view,  $I(t)$  consists of the objects that each have the set of properties  $head(t)$ ,  $I^-(t)$  consists of the objects that each have at least the set of properties  $head(t)$ , i.e. some of the objects in  $I^-(t)$  have one or more properties  $t'$  such that  $t' \preceq t$ , and finally,  $I^+(t)$  consists of the objects that each have at least the set of properties  $head(t) \setminus \{t\}$ .

◇

Referring to Def. 2 let us now define query answering for a general query  $q$ .

**Definition 8** Let  $q$  be a query over a terminology  $T$  and let  $I$  be an interpretation of  $T$ .

(a) The *sure answer* of  $q$ , denoted  $I^-(q)$ , is a set of objects defined as follows:

$$\begin{aligned} I^-(t) &= \bigcup \{I(s) \mid s \in tail(t)\} \\ I^-(q \wedge q') &= I^-(q) \cap I^-(q') \\ I^-(q \vee q') &= I^-(q) \cup I^-(q') \\ I^-(q \wedge \neg q') &= I^-(q) \setminus I^-(q') \\ I^-(\epsilon) &= \emptyset \end{aligned}$$

(b) The *possible answer* of  $q$ , denoted  $I^+(q)$ , is a set of objects defined as follows:

$$\begin{aligned} I^+(t) &= \bigcap \{I^-(u) \mid u \in head(t) \text{ and } u \not\preceq t\} \\ I^+(q \wedge q') &= I^+(q) \cap I^+(q') \\ I^+(q \vee q') &= I^+(q) \cup I^+(q') \\ I^+(q \wedge \neg q') &= I^+(q) \setminus I^-(q') \\ I^+(\epsilon) &= \emptyset \end{aligned}$$

It follows easily from the above definition that for every query  $q$  we have  $I^-(q) \subseteq I^+(q)$ . This means that the sure answer of a query  $q$  is always included in the possible answer of  $q$ .

Note that we interpret  $I^+(q \wedge \neg q')$  by  $I^+(q) \setminus I^-(q')$ , and *not* by  $I^+(q) \setminus I^+(q')$ . This is because if we had interpreted  $I^+(q \wedge \neg q')$  by  $I^+(q) \setminus I^+(q')$  then we could have found queries  $q$  for which  $I^-(q) \supset I^+(q)$ , contrary to intuition. For example, consider a terminology  $T$  with three terms  $a, b$  and  $c$  such that  $c \preceq b \preceq a$ , and an interpretation  $I$  such that  $I(c) = \emptyset$ ,  $I(b) = \{1\}$  and  $I(a) = \{2\}$ . Then for  $q = a \wedge \neg c$  we would have had:  $I^-(q) = I^-(a) \setminus I^-(c) = \{1, 2\}$  and  $I^+(q) = I^+(a) \setminus I^+(c) = \{2\}$ , i.e.  $I^-(q) \supset I^+(q)$ . However, with our definition we have  $I^+(a \wedge \neg c) = I^+(a) \setminus I^-(c) = \{1, 2\}$ , i.e. the relation  $I^- \sqsubseteq I^+$  is preserved.

User interaction with a source consists in submitting a query  $q$  plus the nature of the desired answer (sure or possible). The source then responds by computing  $I^-(q)$  or  $I^+(q)$  according to the user's desire. The possibility of providing two types of answer to a query can enhance the quality of user interaction with the source. For example, the user may submit a query and require a sure answer. If the sure answer is empty, this may mean either that no object has been indexed under the user's query, or that the objects have been indexed at a coarser level. So, if the sure answer turns out to be empty then the user can ask for the possible answer to his query. In the possible answer the user can see objects related to, but not necessarily indexed under his query. Another possibility is that the sure answer to the query is not empty but the user just likes to see more objects related to his query, but at a coarser level. In this case, again, the user can ask for a possible answer to his query.

A source can be implemented using any of a number of data models. For example, using the relational model [18], a source can be implemented as a database schema consisting of three tables, one for storing the terminology, one for storing the subsumption relation, and one for storing the interpretation  $I$ .

```
TERMINOLOGY(term-id:Int, term-name:Str)
SUBSUMPTION(term1:Int, term2:Int)
INTERPRETATION(term-id:Int, obj:Int)
```

Note that each term of the terminology is stored in the form of a pair  $\langle \text{term-id}, \text{term-name} \rangle$  where “term-id” is an internal identifier.

Concerning query evaluation at a source, there are basically two approaches. The first approach consists of computing and storing the models  $I^-$  and  $I^+$  and then using these stored models for computing answers to queries. This can be done using algorithms that follow easily from Definition 8. The advantage of this approach is that answers can be computed in a straightforward manner from the stored models. The disadvantage is in-

creased space requirements as well as increased maintenance costs for the stored models. Indeed, whenever the taxonomy or the interpretation  $I$  change,  $I^-$  and  $I^+$  must be updated appropriately. This requires an efficient method for handling updates since recomputing  $I^-$  and  $I^+$  from scratch would be inefficient.

The second approach consists of storing only the interpretation  $I$  and, whenever a query  $q$  is submitted, computing the appropriate answer,  $I^-(q)$  or  $I^+(q)$ , using  $I$ . The computation of  $I^-(q)$  can be done in a straightforward manner following Definition 8.(a). The computation of  $I^+(q)$  can be done again following Definition 8.(b) but requires the previous computation of  $I^-(t)$  for all terms  $t$  that subsume terms appearing in the query. The advantage of this approach is that we have no additional space requirements and no additional maintenance costs. The disadvantage is increased time cost for the computation of the answers.

The relative merits of the two approaches depend on the application at hand as well as on the frequency by which the taxonomy and/or the stored interpretation of the source are updated. In both approaches we need algorithms for computing the head and the tail of a term. However, if we compute the transitive closure of the subsumption relation by one of the existing algorithms (e.g. see [55]), then the algorithms for computing the head and tail of a term follow immediately from Definition 5. The complexity of evaluating the transitive closure of  $\preceq$  is polynomial. For instance, the time complexity of the Floyd-Warshall algorithm is cubic in the number of terms, and the space used is at most quadratic in the number of terms. If the entire subsumption relation  $\preceq$  is stored, i.e. if the transitive links are stored, then the computation of  $head(t)$  and  $tail(t)$  can be done in  $O(|\preceq|)$  time. If only the interpretation  $I$  is stored, then the computation of  $I^-(t)$  requires taking the union of at most  $|T|$  subsets of  $Obj$ . If  $U$  denotes the set of objects that are stored in the source<sup>3</sup> then the union of two subsets of  $Obj$  can be computed in  $O(|U|)$  time. Thus the computation of  $I^-(t)$  can be done in  $O(|T| * |U|)$  time<sup>4</sup>. If the sure model  $I^-$  is stored, then the computation of  $I^+(t)$  requires taking the intersection of at most  $|T|$  subsets of  $Obj$ . Thus the computation of  $I^+(t)$  can be done in  $O(|T| * |U|)$  time. If only the interpretation  $I$  is stored, then the computation of  $I^+(t)$  can be done as follows:

$$I^+(t) = \bigcap_{u>t} \left( \bigcup \{I(s) \mid s \preceq u\} \right)$$

This computation can be done in  $O(|T|^2 * |U|)$  time.

<sup>3</sup> Specifically,  $U = \{o \in Obj \mid \exists t \in T \text{ s.t. } o \in I(t)\}$ .

<sup>4</sup> Note that here we express the execution time with respect to two parameters: the size of the terminology and the number of the stored objects.

### 3 The Mediator

So far we have seen that an information source over an underlying set of objects  $Obj$  consists of:

- 1) a taxonomy  $(T, \preceq)$ , and
- 2) a stored interpretation  $I$  of  $T$ .

The terminology  $T$  contains terms that are familiar to the users of the source; the subsumption relation  $\preceq$  contains relationships between terms of  $T$ ; and the stored interpretation  $I$  associates each term  $t$  with the objects that are indexed under  $t$  (by the indexer).

Consider now a set of sources  $S_1, \dots, S_k$  over the *same* underlying set of objects  $Obj$ . In general, two different sources may have different terminologies either because the users of the two sources are familiar with different sets of terms, or because one source indexes objects at a different level of granularity than the other. The two sources may also have different subsumption relations as the relationships between any two given terms may be perceived differently in the two sources. Finally, two different sources may have different stored interpretations, for example some objects may have been indexed by one source but not by the other.

Clearly if one wants to combine or *integrate* information coming from different sources one has to cope with the above heterogeneities. One way of rendering all these heterogeneities transparent to users is through the use of *mediators* (initially proposed in [78]).

The problem of information integration has attracted considerable attention in the last few years, especially in the area of databases (see [32] for a comprehensive overview). The main idea is to have users access the information sources through a common schema that reflects their needs. Two main approaches seem to have emerged, namely the virtual view approach and the materialized view approach. In the first, only the common schema is stored (but no data), while in the second (which is also called the warehouse approach) both the common schema and data over that schema are stored. Our approach is similar in spirit to the virtual view approach.

In our approach, a mediator  $M$  has a taxonomy  $(T, \preceq)$  that reflects the needs of its potential users but has *no* stored objects. Instead, each term at the mediator is related directly or indirectly with terms in the underlying sources. More formally, a mediator is defined as follows:

**Definition 9** A mediator  $M$  over  $k$  sources  $S_1 = \langle (T_1, \preceq_1), I_1 \rangle, \dots, S_k = \langle (T_k, \preceq_k), I_k \rangle$  consists of:

- 1) a taxonomy  $(T, \preceq)$ , and
- 2) a set of *articulations*  $a_i$ , one for each source  $S_i$ ; each articulation  $a_i$  is a subsumption relation over  $T \cup T_i$ .

Roughly speaking, a mediator is just like a source but with an important difference: there is no interpretation stored at the mediator. What is stored at the mediator, instead, is the set of articulations  $a_i$ , one for each source  $S_i$ . For example, suppose that we want to

integrate two Web catalogues which provide access to pages about electronic products. In particular, consider the sources  $S_1$  and  $S_2$  shown in Figure 7 and assume that we want to provide access to these sources through a mediator  $M$  as shown in that figure. To achieve integration we enrich the mediator with articulations, i.e. with relationships that relate the terms of the mediator with the terms of the sources as shown in Figure 7. The articulations  $a_1$  and  $a_2$  shown in Figure 7 are the following sets of subsumption relationships:

$$a_1 = \{\text{PhotoCameras} \preceq \text{Cameras}, \\ \text{StillCameras} \preceq \text{PhotoCameras}, \\ \text{Miniature} \preceq \text{StillCameras}, \\ \text{Instant} \preceq \text{StillCameras}, \text{Reflex}_1 \preceq \text{StillCameras}, \\ \text{Reflex}_1 \preceq \text{Reflex}, \text{Reflex} \preceq \text{Reflex}_1\}$$

$$a_2 = \{\text{Products} \preceq \text{Electronics}, \text{SLRCams} \preceq \text{Reflex}, \\ \text{VideoCams} \preceq \text{MovingPictureCams}, \\ \text{MovingPictureCams} \preceq \text{VideoCams}\}$$

Note that  $a_1$  is a subsumption relation over  $T \cup T_1$  and  $a_2$  is a subsumption relation over  $T \cup T_2$ , as required by the definition of an articulation (Def. 9).

Figure 8 shows another example of a mediator over three sources. These three sources provide access to tourist information and the information is organized by location.

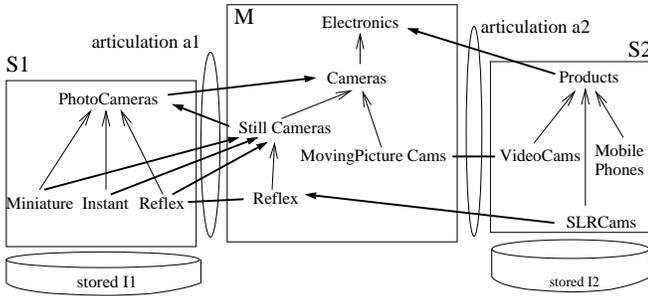


Fig. 7 A mediator over two catalogues of electronic products

Now, in the presence of several sources, one and the same term may appear in two or more sources. If the same term appears in two different sources then we consider the two appearances as two different terms. This is denoted here by subscripting each term of a source  $S_i$  by the subscript  $i$ , and can be implemented in practice by, say, prefixing each term by the name of the source in which the term appears. Take for example the term DB and suppose that it appears in sources  $S_i$  and  $S_j$ . Then, from the mediator's point of view there are two distinct terms: the term  $DB_i$  in source  $S_i$  and the term  $DB_j$  in source  $S_j$ . This is reasonable as the same term can have different interpretations (meanings) in different sources.

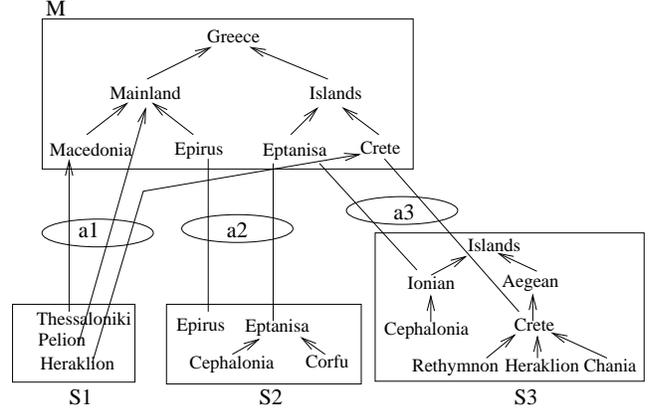


Fig. 8 A mediator over three catalogues of tourist information

Thus for every  $i \neq j$  we assume  $T_i \cap T_j = \emptyset$ ; and for every  $i$  we assume  $T \cap T_i = \emptyset$ . In this way we overcome the problems of homonyms. Under these assumptions, two terms are considered equivalent, e.g.  $DB_i \sim DB_j$ , only if they can be shown to be equivalent through the articulations  $a_i$  and  $a_j$ , e.g.  $DB_i$  and  $DB_j$  are equivalent if there is a term  $t$  in  $T$  such that  $t \sim_{a_i} DB_i$  and  $t \sim_{a_j} DB_j$ .

Integrating objects from several sources often requires *restoring the context* of these objects, i.e. adding information that is missing from the original representation of the objects which concerns the context of the objects. Consider for example a mediator which provides access to electronic products according to the *type* of the products and according to the *location* of the stores that sell these products. Suppose that the mediator has two underlying sources  $S_1$  and  $S_2$  as shown in Figure 9. Assume that  $S_1$  is the source of a store located in Heraklion, while  $S_2$  is the source of a store located in Paris. The context of the objects of each source, here the location of the store that sells each product, can be restored by adding appropriate relationships to the articulations. Specifically, for defining that all **PhotoCameras** of the source  $S_1$  are available through a store located in **Heraklion**, it suffices to put in the articulation  $a_1$  the relationship:

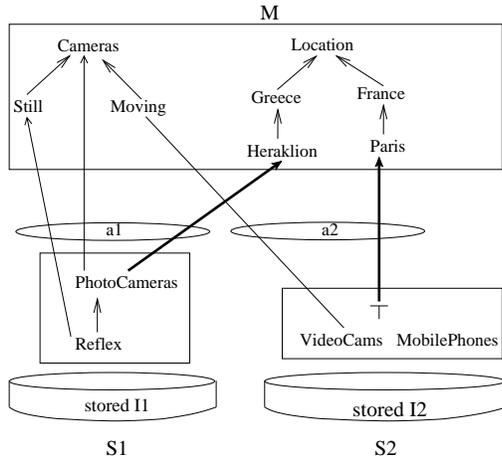
$$\text{PhotoCameras}_1 \preceq \text{Heraklion}$$

while for defining that all products of the source  $S_2$  are available through a store located in **Paris**, it suffices to put in the articulation  $a_2$  the following relationship:

$$\top_2 \preceq \text{Paris}$$

This example demonstrates how the articulations of the mediator can restore the context of the objects.

Turning now to query answering, we recall that the mediator receives queries over its own terminology  $T$ . Now, as the mediator has no stored interpretation of  $T$ , the only way to obtain one is by *querying* the underlying



**Fig. 9** Using articulations to restore the context of the objects of the sources

sources. However, as the mediator and the sources have different terminologies, for computing the interpretation of a term  $t \in T$ , the mediator sends to each source  $S_i$  a *translation* of  $t$ , i.e. a query that can be answered by the source, and then it takes the *union* of the answers returned by the sources. The definition of translations is based on the articulations of the mediator.

Thus we will actually define an interpretation  $I$  of the mediator terminology, based on the interpretations  $I_i$  stored at the sources, on the one hand, and on the articulations  $a_i$ ,  $i = 1, \dots, k$ , on the other. Conceptually, once the interpretation  $I$  of the mediator is defined, the mediator can answer queries just like any other source does, i.e. from its sure model  $I^-$  and from its possible model  $I^+$ .

In order to define the mediator interpretation  $I$  we proceed as follows: for every term  $t$  of the mediator terminology  $T$ :

1. first we define a translation  $t^i$  of  $t$  in  $a_i$ , in the form of a query to source  $S_i$ ,  $i = 1, \dots, k$ ;
2. then we evaluate the query  $t^i$  at source  $S_i$ ,  $i = 1, \dots, k$ , and
3. finally we define  $I(t)$  by taking the union of the answers to the queries  $t^i$  returned by the sources.

Now, there are two ways to translate  $t$  using the articulation  $a_i$ , that we shall call the *upper approximation* of  $t$  and the *lower approximation* of  $t$  in  $a_i$ . Roughly, the upper approximation of  $t$  in  $a_i$  is the conjunction of all terms of  $T_i$  that subsume  $t$  in  $a_i$ , and the lower approximation of  $t$  in  $a_i$  is the disjunction of all terms of  $T_i$  that  $t$  subsumes in  $a_i$ . In order to define these notions formally we need the notions of *tail* and *head* of a term *relative* to an articulation:

**Definition 10** Given a term  $t \in T$  and articulation  $a_i$  we define

$$\text{tail}_i(t) = \{s \in T_i \mid sa_i t\} \text{ and } \text{head}_i(t) = \{u \in T_i \mid ta_i u\}$$

**Definition 11** Let  $M = (T, \preceq, a_1, \dots, a_k)$  be a mediator over sources  $S_1, \dots, S_k$ . If  $t$  is a term of  $T$  then

- the *lower approximation* of  $t$  with respect to  $a_i$ , denoted  $t_l^i$ , is defined by

$$t_l^i = \bigvee \text{tail}_i(t)$$

- the *upper approximation* of  $t$  with respect to  $a_i$ , denoted  $t_u^i$ , is defined by

$$t_u^i = \begin{cases} \bigwedge \text{head}_i(t), & \text{if } \text{head}_i(t) \neq \emptyset \\ t_l^i, & \text{otherwise} \end{cases}$$

Note that if  $\text{head}_i(t) = \emptyset$  then we consider that  $t_u^i = t_l^i = \bigvee \text{tail}_i(t)$ . The reason behind this choice is that we want the interpretation obtained by using lower approximation to be less than or equal to ( $\sqsubseteq$ ) the interpretation obtained by using the upper approximation.

Here are some examples of approximations for the mediator shown in Figure 7:

$$\begin{aligned} \text{StillCameras}_l^1 &= \text{Miniature} \vee \text{Instant} \vee \text{Reflex} \\ \text{StillCameras}_u^1 &= \text{PhotoCameras} \\ \text{Reflex}_l^1 &= \text{Reflex} \\ \text{Reflex}_u^1 &= \text{Reflex} \wedge \text{PhotoCameras} \\ \text{Reflex}_l^2 &= \text{SLRCams} \\ \text{Cameras}_l^1 &= \text{PhotoCameras} \vee \text{Miniature} \vee \text{Instant} \vee \text{Reflex} \\ \text{Cameras}_u^1 &= \text{PhotoCameras} \vee \text{Miniature} \vee \text{Instant} \vee \text{Reflex} \\ \text{MovingPictureCams}_u^1 &= \text{MovingPictureCams}_l^1 = \epsilon \end{aligned}$$

Note that for a given term  $t \in T$  the evaluation of  $t_u^i$  requires the previous evaluation of  $\text{head}_i(t)$ , and the evaluation of  $t_l^i$  requires the previous evaluation of  $\text{tail}_i(t)$ . However, if we compute the transitive closure of  $a_i$  then the evaluation of  $\text{head}_i(t)$  and  $\text{tail}_i(t)$  is straightforward.

Now, the approximations  $t_u^i$  and  $t_l^i$  of  $t$  are actually queries to the source  $S_i$ , and as such each can have a sure answer and a possible answer (see Section 2). As a consequence, we can define at least four different interpretations  $I$  for the mediator. Assuming for simplicity that *all* sources respond in the same manner, i.e. either all give a sure answer or all give a possible answer, we can define exactly four interpretations for the mediator that we shall denote by  $I_{l-}$ ,  $I_{l+}$ ,  $I_{u-}$ ,  $I_{u+}$ . These interpretations are defined as follows:

- 1 Lower approximation of  $t$  at mediator and sure answer from sources:  

$$I_{l-}(t) = \bigcup_{i=1}^k I_i^-(t_l^i)$$
- 2 Lower approximation of  $t$  at mediator and possible answer from sources:  

$$I_{l+}(t) = \bigcup_{i=1}^k I_i^+(t_l^i)$$
- 3 Upper approximation of  $t$  at mediator and sure answer from sources:  

$$I_{u-}(t) = \bigcup_{i=1}^k I_i^-(t_u^i)$$
- 4 Upper approximation of  $t$  at mediator and possible answer from sources:  

$$I_{u+}(t) = \bigcup_{i=1}^k I_i^+(t_u^i)$$

So, the mediator can answer queries submitted by its users based on any of the above four interpretations. Moreover, for any of these four interpretations, the mediator can give either a sure answer or a possible answer - just like any source can (see Section 2). By consequence, we can distinguish eight possible modes in which a mediator can operate. Each mode essentially corresponds to a different answer model of the mediator. The operation modes of a mediator and the corresponding answer models are summarized in Table 1.

<i>oper. mode at the med.</i>	<i>term approx. at med.</i>	<i>query eval. at source</i>	<i>query eval. at med.</i>	<i>the answer model of the med.</i>
1	lower	sure	sure	$I_{l-}^-$
2	lower	possible	sure	$I_{l+}^-$
3	upper	sure	sure	$I_{u-}^-$
4	upper	possible	sure	$I_{u+}^-$
5	lower	sure	possible	$I_{l-}^+$
6	lower	possible	possible	$I_{l+}^+$
7	upper	sure	possible	$I_{u-}^+$
8	upper	possible	possible	$I_{u+}^+$

**Table 1** Modes in which a mediator can operate

Very roughly speaking, as we go down the table (from mode 1 to 8) the answer to the same user query is more likely to contain objects that are not “relevant” to the query. This is described more precisely in the following proposition.

**Proposition 3** The answer models of the mediator are ordered as follows:

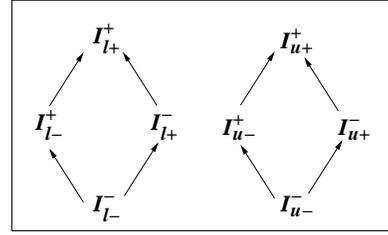
- (a)  $I_{l-}^- \sqsubseteq I_{l+}^-$
- (b)  $I_{u-}^- \sqsubseteq I_{u+}^-$
- (c)  $I_{l-}^+ \sqsubseteq I_{l+}^+$
- (d)  $I_{u-}^+ \sqsubseteq I_{u+}^+$
- (e)  $I_{l-}^- \sqsubseteq I_{l-}^+$
- (f)  $I_{l+}^- \sqsubseteq I_{l+}^+$
- (g)  $I_{u-}^- \sqsubseteq I_{u-}^+$
- (h)  $I_{u+}^- \sqsubseteq I_{u+}^+$

◇

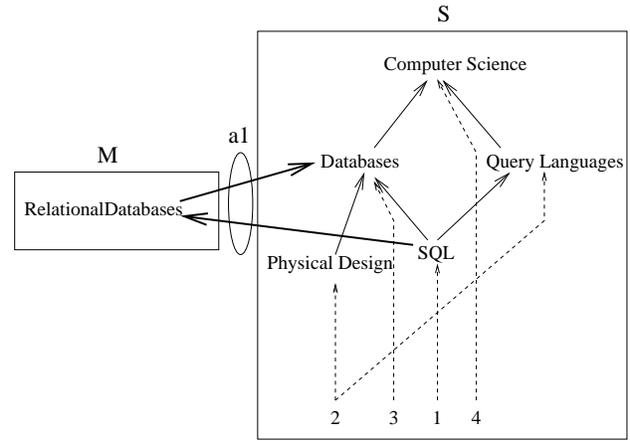
Figure 10 shows graphically the orderings of the above proposition. The nodes represent the answer models shown in Table 1. An arrow from node  $m$  to a node  $n$  means that  $m \sqsubseteq n$ .

For example the interpretation of the term **RelationalDatabases** of the mediator shown in Figure 11, in each of the models  $I_{l-}^-$ ,  $I_{l+}^-$ ,  $I_{u-}^-$ ,  $I_{u+}^-$  is as follows:

$$\begin{aligned}
 I_{l-}^- (\text{RelationalDatabases}) &= \{1\} \\
 I_{l+}^- (\text{RelationalDatabases}) &= \{1, 2\} \\
 I_{u-}^- (\text{RelationalDatabases}) &= \{1, 2, 3\} \\
 I_{u+}^- (\text{RelationalDatabases}) &= \{1, 2, 3, 4\}
 \end{aligned}$$



**Fig. 10** The ordering ( $\sqsubseteq$ ) of the eight answer models of the mediator



**Fig. 11** A mediator over one source

Another example of mediator operation is given in Figure 12. Figure 12.(a) shows a mediator having an articulation to a source  $S_1$  and Figure 12.(b) shows two tables. The table at the upper part of the figure shows the interpretation  $I_1$  of source  $S_1$  and the corresponding (sure and possible) models. The first column of the table at the bottom part shows three queries which are actually the three terms of  $T$ . The subsequent columns show what the mediator returns in each of the first four operation modes.

The operation modes of the mediator can either be decided (and fixed) by the mediator designer at design time, or indicated by the mediator users at query time. We can distinguish at least three approaches:

- *Fixed Approach.* The mediator designer selects and fixes one of the eight possible modes of operation for the mediator and the sources, and users simply submit their queries to the mediator without any further indication.
- *Variable Approach.* The mediator users submit their queries along with a specification for the query evaluation mode they wish. This is done by providing values to the mediator for selecting one of the eight operation modes from Table 1. For example, the following user specification selects the operation mode number 3 from Table 1:

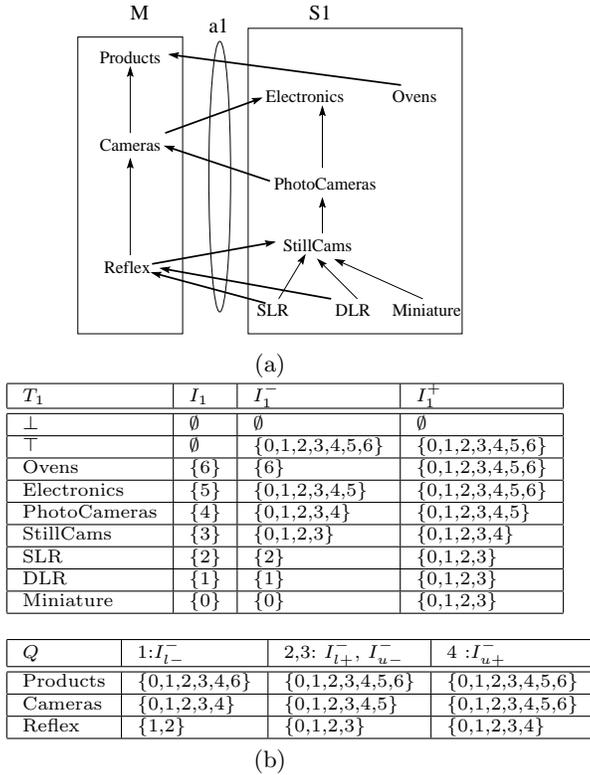


Fig. 12 A mediator with one articulation to a source  $S_1$

term approximation at mediator = **upper**  
 query evaluation at source = **sure**  
 query evaluation at mediator = **sure**

- *Mixed Approach.* The mediator designer selects and fixes some of the attributes of Table 1, and the user provides the remaining ones. For example, the designer may select and fix the query evaluation mode at source (i.e. sure or possible) and the kind of term approximation at the mediator (i.e. lower or upper approximation), during design time, while the users select the query evaluation mode at the mediator, during query time.

Clearly, selecting one of the above approaches depends on several factors, such as the reliability of the sources, the level of expertise of the users, and so on. One can even think of more sophisticated modes of mediator operation than those presented in Table 1. For example, the mediator designer may assign a degree of reliability to each source and then ask sources to evaluate queries in a mode depending on their degree of reliability. In this paper, however, we do not pursue this idea any further.

We have seen so far how the mediator communicates with the sources through the articulations. In fact, the articulations are the *only* means of communication between the sources and the mediator. Now, certain kinds of articulation are better than others. One kind of articulations that are of particular interest are those that

ensure what we call “compatibility” between the sources and the mediator.

**Definition 12** A source  $S_i$  is *compatible* with the mediator  $M$  if for any terms  $s, t$  in  $T_i$ , if  $sa_it$  then  $s \preceq_i t$ .

That is,  $S_i$  is compatible with the mediator whenever the following condition holds: for all terms  $s$  and  $t$  in  $T_i$ , if  $s$  is subsumed by  $t$  in the articulation  $a_i$  then  $s$  is also subsumed by  $t$  in  $\preceq_i$ .

For example, the source  $S_1$  of Figure 7 is compatible with the mediator since we have  $\text{Miniature } a_1 \text{ PhotoCameras}$  and  $\text{Miniature } \preceq_1 \text{ PhotoCameras}$ ,  $\text{Instant } a_1 \text{ PhotoCameras}$  and  $\text{Instant } \preceq_1 \text{ PhotoCameras}$ ,  $\text{Reflex } a_1 \text{ PhotoCameras}$  and  $\text{Reflex } \preceq_1 \text{ PhotoCameras}$ .

An interesting consequence of compatibility is that if a source  $S_i$  is compatible with the mediator, then in every model  $I_i$  of  $S_i$  the following condition holds:  $I_i(t_l^i) \subseteq I_i(t_u^i)$  for each mediator term  $t$ , where  $t_l^i$  is the lower approximation of  $t$  and  $t_u^i$  is the upper approximation of  $t$ . From this property we infer that if all sources are compatible with the mediator then the ordering relation over the eight answer models of the mediator (see Figure 10), is enriched as stated by the following proposition.

**Proposition 4** If all sources are compatible with the mediator then:

- (1)  $I_{l-}^- \sqsubseteq I_{u-}^-$
- (2)  $I_{l+}^- \sqsubseteq I_{u+}^-$
- (3)  $I_{l-}^+ \sqsubseteq I_{u-}^+$
- (4)  $I_{l+}^+ \sqsubseteq I_{u+}^+$

As a result, the two diagrams of Figure 10 are now connected in a single diagram as shown in Figure 13.

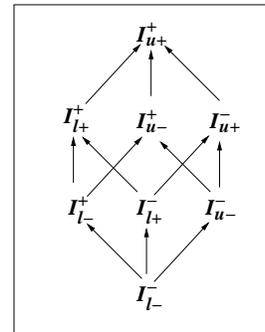


Fig. 13 The ordering ( $\sqsubseteq$ ) of the eight answer models of the mediator in the case where all sources are compatible with the mediator

Note that the above ordering relationships do not hold necessarily if the sources are not compatible with the mediator. For example, consider a source  $S_1$  with terminology  $T_1 = \{b, b'\}$  and no subsumption relationships.

Suppose that the source has a stored interpretation  $I_1$  defined as follows:  $I_1(b) = \{1\}$  and  $I_1(b') = \{2\}$ . Now consider a mediator connected to source  $S_1$  through the articulation  $a_1 = \{b \preceq t, t \preceq b'\}$ , where  $t$  is a term of the mediator. Notice that  $S_1$  is not compatible with the mediator because  $b$  is subsumed by  $b'$  in  $a_1$  while  $b$  is not subsumed by  $b'$  in  $\preceq_1$ , i.e.  $ba_1b'$  and  $b \not\preceq_1 b'$ . Here we have  $t_l^1 = b$  and  $t_u^1 = b'$ , thus  $I_1^-(t_l^1) = \{1\}$  and  $I_1^-(t_u^1) = \{2\}$ . It follows that  $I_1^-(t_l^1) \not\subseteq I_1^-(t_u^1)$ , which implies  $I_{l-}^-(t) \not\subseteq I_{u-}^-(t)$ . From this example we see that if the underlying sources are not compatible with the mediator then  $I_{l-}^- \subseteq I_{u-}^-$  does not hold.

Another interesting implication of compatibility concerns the efficiency of query evaluation. Let  $s, t$  be two terms in  $T_i$  which are known to the mediator (through  $a_i$ ) and assume that the mediator knows that source  $S_i$  is compatible. In this case if  $sa_it$  then  $s \preceq_i t$ . From this knowledge the mediator can conclude that  $I_i(s) \subseteq I_i(t)$ , in every model  $I_i$  of  $T_i$ , and thus  $I_i(s) \cap I_i(t) = I_i(s)$  and  $I_i(s) \cup I_i(t) = I_i(t)$ . This means that the mediator can retain only the minimal elements of the set  $head_i(t)$  and still obtain the same answer for the query  $t_u^i$  from source  $S_i$ . Therefore, if the mediator knows that source  $S_i$  is compatible, then instead of sending to source  $S_i$  the query  $\bigwedge head_i(t)$ , the mediator can send the query  $\bigwedge min(head_i(t))$ . Similarly, in the set  $tail_i(t)$  the mediator can retain only the maximal elements and still obtain the same answer for the query  $t_l^i$  from source  $S_i$ , i.e., instead of sending the query  $\bigvee tail_i(t)$  to source  $S_i$ , the mediator can send the query  $\bigvee max(tail_i(t))$ .

For example, in Figure 7, as source  $S_1$  is compatible with the mediator, the lower approximation of the term **Camera** is the term **PhotoCameras**. If  $S_1$  were not compatible then the lower approximation of **Camera** would be the disjunction **PhotoCameras**  $\vee$  **Miniature**  $\vee$  **Instant**  $\vee$  **Reflex**.

Thus if  $S_i$  is compatible then  $t_u^i = \bigwedge min(head_i(t))$  and  $t_l^i = \bigvee max(tail_i(t))$ . In this case the evaluation of  $t_u^i$  and  $t_l^i$  can be done more efficiently without having to compute the transitive closure of  $a_i$ . Specifically, for evaluating  $max(tail_i(t))$  we traverse in depth-first-search the relation  $a_i$  starting from the term  $t$ . If an element  $t'$  of  $T_i$  is reached then this term is “collected” and the algorithm does not traverse any other element subsumed by  $t'$  (in  $a_i$ ). All elements of  $T_i$  which were collected during the traversal are then returned. We can evaluate  $min(head_i(t))$  analogously. We conclude that if a source is compatible then the approximation of a term for that source can be done more efficiently especially when the articulation to that source is big. Moreover the resulting approximations are shorter, which implies that their transmission requires less time and the underlying source can evaluate these queries more efficiently.

Note that maintaining compatibility is not an easy task. Of course, the designer of the mediator can initially design articulations such that the underlying sources are compatible. However, an update at a source  $S_i$  or at the

mediator (changing either  $T$  or  $a_i$ ) may destroy compatibility. Therefore the mediator should (periodically) check the compatibility of its sources, e.g. by submitting to them queries allowing to check whether  $t \preceq_i t'$ .

#### 4 Query Evaluation at the Mediator

We have seen how the two possible approximations at the mediator (lower or upper) and the two possible query evaluation modes at the sources (sure or possible) give rise to four possible interpretations at the mediator:  $I_{l-}$ ,  $I_{l+}$ ,  $I_{u-}$  and  $I_{u+}$ . If these four interpretations were stored at the mediator then the interaction between a user and the mediator would be straightforward, i.e.

- the user submits a query to the mediator (as if it were a usual source)
- the mediator and/or the user specifies the answer model to be used
- the mediator uses the specified model to provide a sure or possible answer to the query (as it is done in a usual source)

However, there is *no* interpretation actually stored at the mediator, so to answer queries the mediator has to call on the underlying sources, submit to them appropriate queries, then merge the results to produce the final answer for the user. Therefore, the crucial tasks for the evaluation of user queries at the mediator can be summarized as follows:

- translate the user’s query into a set of queries to the underlying sources, i.e. determine *what* queries to send to *which* sources;
- merge the results returned by the sources in order to produce the answer to the user’s query.

Clearly, the complexity of these tasks depends on the nature of the user query, i.e.

- the form of the query (single term, disjunction of terms, etc.),
- the answer model used by the mediator.

In what follows we analyze the complexity of query evaluation at the mediator with respect to the form that a user query can have, and the answer model used by the mediator for evaluating the query.

The complexity measure that we use in our analysis is the *number of queries* that the mediator sends to the sources in order to answer the user’s query, and the *execution time* expressed in terms of several parameters, such as the size of the mediator terminology, the size of the articulations, the number of sources, the length of the query and the size of the domain. However, we believe that the number of queries that the mediator needs to send is the most important measure, as the mediator spends a lot of time waiting for the answers of the sources.

We are aware that, in doing so, we do not take into account the complexity of query evaluation at each source. However, the mediator has little or no control over how queries are evaluated at individual sources. This is especially true for the applications that we have in mind (Web environment), where the mediator is set up by individual users who have no control over the underlying sources (which are Web catalogues).

In the complexity analysis that follows we consider a mediator over  $k$  sources,  $S_1, \dots, S_k$ . Note that we write  $I_l$  instead of  $I_{l-}$  or  $I_{l+}$ , and  $I_u$  instead of  $I_{u-}$  or  $I_{u+}$ , since the translation and the evaluation of queries at the mediator does not depend on the evaluation of queries at the underlying sources. At first we describe the evaluation of queries in the sure models of the mediator, i.e. in the models  $I_l^-$  and  $I_u^-$ .

An interesting remark here is that the mediator *will not* necessarily query all sources. A source is queried only if the evaluation of the answer requires sending a sub-query to that source, otherwise the source is not queried. Thus query translation also determines the selection of the sources.

We study separately the following forms of queries: *single term queries*, *disjunctive queries*, *conjunctive queries*, *CNF queries* and *DNF queries*.

– *Single Term Queries*

**Proposition 5** If the query is a single term, i.e.  $q = t \in T$ , then  $I_l^-(t)$  and  $I_u^-(t)$  can be evaluated as follows:

$$I_l^-(t) = \bigcup_{i=1..k} I_i(q_l^i(t)) \quad \text{where} \quad q_l^i(t) = \bigvee \{s_i^i \mid s \preceq t\}$$

$$I_u^-(t) = \bigcup_{i=1..k} I_i(q_u^i(t)) \quad \text{where} \quad q_u^i(t) = \bigvee \{s_u^i \mid s \preceq t\}$$

◇

This means that the mediator  $M$  can evaluate the query by sending at most one query to each source. Thus  $M$  will send at most  $k$  queries. Note that if  $q_l^i(t) = \epsilon$  (or  $q_u^i(t) = \epsilon$ ) then  $M$  does not have to send any query to source  $S_i$ .

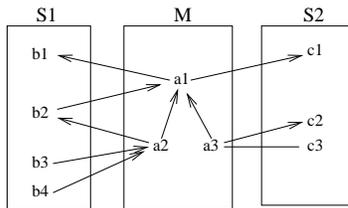


Fig. 14 A mediator over two sources

For example consider a mediator over two sources as shown in Figure 14. The answer in  $I_l^-$  of the query

$q = a_1$  can be evaluated as follows:

$$\begin{aligned} I_l^-(\mathbf{a}_1) &= I_1(q_l^1(\mathbf{a}_1)) \cup I_2(q_l^2(\mathbf{a}_1)) \quad \text{where} \\ q_l^1(\mathbf{a}_1) &= \mathbf{b}_2 \vee (\mathbf{b}_3 \vee \mathbf{b}_4) \\ q_l^2(\mathbf{a}_1) &= \mathbf{c}_3 \end{aligned}$$

while the answer  $I_u^-$  can be evaluated as follows:

$$\begin{aligned} I_u^-(\mathbf{a}_1) &= I_1(q_u^1(\mathbf{a}_1)) \cup I_2(q_u^2(\mathbf{a}_1)) \quad \text{where} \\ q_u^1(\mathbf{a}_1) &= \mathbf{b}_1 \vee \mathbf{b}_2 \\ q_u^2(\mathbf{a}_1) &= \mathbf{c}_1 \vee (\mathbf{c}_2 \wedge \mathbf{c}_3) \end{aligned}$$

If the mediator knows that a source  $S_i$  is *compatible* then the mediator can set

$$q_l^i(t) = \bigvee \max(\bigcup \{tail_i(s) \mid s \preceq t\})$$

Note that if the entire articulation  $a_i$  is stored (including the transitive links), then the computation of  $t_l^i$  can be done in  $O(|a_i|)$  time. The same holds for  $t_u^i$ . Thus the computation of  $q_l^i(t)$  can be done in  $O(|T| * |a_i|)$  time. The same holds for  $q_u^i(t)$ . This means that the computation of all  $q_l^i(t)$ , or  $q_u^i(t)$ , for  $i = 1..k$  can be done in  $O(|T| * |a|)$  where  $a$  denotes the union of all articulations, i.e.  $a = a_1 \cup \dots \cup a_k$ .

Now, the set operations over the answers returned by the sources that are needed for computing  $I_l^-(t)$ , can be performed in  $O(k * U)$  time.

Thus the total computation needed by the mediator can be done in  $O(|T| * |a| + k * U)$  time.

– *Disjunctive Queries*

If the query is a disjunction of terms, i.e.  $q = t_1 \vee \dots \vee t_n$  then

$$I_l^-(t_1 \vee \dots \vee t_n) = \bigcup_{i=1..k} I_i(q_l^i(t_1) \vee \dots \vee q_l^i(t_n))$$

$$I_u^-(t_1 \vee \dots \vee t_n) = \bigcup_{i=1..k} I_i(q_u^i(t_1) \vee \dots \vee q_u^i(t_n))$$

Again, the mediator can evaluate the query by sending at most one query to each source.

If, furthermore, a source  $S_i$  is *compatible* then the mediator can send to  $S_i$  the query:

$$\bigvee \max(\bigcup_{j=1..n} (\cup \{tail_i(s) \mid s \preceq t_j\}))$$

Clearly, the computation of each  $q_l^i(t_1) \vee \dots \vee q_l^i(t_n)$  can be done in  $O(|T| * |a_i| * n)$  time. Thus, the computation of all  $q_l^i(t_1) \vee \dots \vee q_l^i(t_n)$ , for  $i = 1..k$  can be done in  $O(|T| * |a| * n)$  time.

The set operations for computing  $I_l^-(t)$  can be performed in  $O(k * U)$  time.

Thus the total computation needed by the mediator can be done in  $O(|T| * |a| * n + k * U)$  time.

– *Conjunctive Queries*

If the query is a conjunction of terms, i.e.  $q = t_1 \wedge \dots \wedge t_n$ , then

$$I_l^-(t_1 \wedge \dots \wedge t_n) = \bigcap_{j=1..n} \left( \bigcup_{i=1..k} I_i(q_l^i(t_j)) \right)$$

$$I_u^-(t_1 \wedge \dots \wedge t_n) = \bigcap_{j=1..n} \left( \bigcup_{i=1..k} I_i(q_u^i(t_j)) \right)$$

Thus for evaluating the query the mediator has to send at most one query to each source for each term which appears in the conjunction. This means that the mediator will send at most  $k * n$  queries.

The computation of all  $q_l^i(t_j)$ , for  $j = 1..n$ , can be done in  $O(|T| * |a| * n)$  time.

The set operations for computing  $I_l^-(t)$  can be performed in  $O(k * n * U)$  time.

Thus the total computation needed by the mediator can be done in  $O(|T| * |a| * n + k * U * n)$  time.

– *Conjunctive Normal Form Queries (CNF Queries)*

A CNF query is a conjunction of maxterms where each maxterm is either a single term or a disjunction of distinct terms ([30]), i.e.  $q = d_1 \wedge \dots \wedge d_m$  where  $d_j = t_{j1} \vee \dots \vee t_{jn_j}$ ,  $j = 1..m, n_j \leq |T|$ . In this case:

$$I_l^-(q) = \bigcap_{j=1..m} \left( \bigcup_{i=1..k} I_i(q_l^i(t_{j1}) \vee \dots \vee q_l^i(t_{jn_j})) \right)$$

$$I_u^-(q) = \bigcap_{j=1..m} \left( \bigcup_{i=1..k} I_i(q_u^i(t_{j1}) \vee \dots \vee q_u^i(t_{jn_j})) \right)$$

The mediator first evaluates each maxterm (disjunction) by sending at most one query to each source and then it takes the intersection of the returned results. This means that the mediator will send at most  $k * m$  queries where  $m$  is the number of maxterms.

Let  $l$  be the length of the query, that is the number of term appearances in the query, i.e.  $l = \sum_{j=1..m} n_j$ .

The computation of  $q_l^i(t)$ ,  $i = 1..k$ , for all  $t$  that appear in  $q$ , can be done in  $O(|T| * |a| * l)$  time.

The set operations for computing  $I_l^-(t)$  can be performed in  $O(k * m * U)$  time.

Thus the total computation needed by the mediator can be done in  $O(|T| * |a| * l + k * m * U)$  time.

– *Disjunctive Normal Form Queries (DNF Queries)*

A DNF query is a disjunction of minterms where a minterm is either a single term or a conjunction of distinct terms, i.e.  $q = c_1 \vee \dots \vee c_m$  where  $c_j = t_{j1} \wedge \dots \wedge t_{jn_j}$ ,  $j = 1..m, n_j \leq |T|$ . In this case:

$$I_l^-(q) = \bigcup_{j=1..m} \left( \bigcap_{h=1..n_j} \left( \bigcup_{i=1..k} I_i(q_l^i(t_{jh})) \right) \right)$$

$$I_u^-(q) = \bigcup_{j=1..m} \left( \bigcap_{h=1..n_j} \left( \bigcup_{i=1..k} I_i(q_u^i(t_{jh})) \right) \right)$$

Thus  $M$  will send at most  $k * l$  queries, where  $l$  is the length of the query.

The computation of all  $q_l^i(t)$ , for  $i = 1..k$ , for all  $t$  that appear in  $q$  can be done in  $O(|T| * |a| * l)$  time.

The set operations for computing  $I_l^-(t)$  can be performed in  $O(k * l * U)$  time.

Thus the total computation needed by the mediator can be done in  $O(|T| * |a| * l + k * l * U)$  time.

Table 2 summarizes the *number of calls* complexity and Table 3 the time complexity. Note that any query that contains the logical connectives  $\wedge$  and  $\vee$  can be converted to DNF or CNF by using one of the existing algorithms (e.g. see [30]). In our case CNF is preferred to DNF since the evaluation of a query in CNF requires sending a smaller number of queries to the sources. For this reason, the mediator first converts the user query in CNF and then it evaluates the CNF query by sending queries to the sources.

Query Form		Max. num. of calls
single term	$t$	$k$
disjunction	$t_1 \vee \dots \vee t_n$	$k$
conjunction	$t_1 \wedge \dots \wedge t_n$	$k * n$
CNF	$d_1 \wedge \dots \wedge d_m$ where $d_j = t_{j1} \vee \dots \vee t_{jn_j}$	$k * m$
DNF	$c_1 \vee \dots \vee c_m$ where $c_j = t_{j1} \wedge \dots \wedge t_{jn_j}$	$k * \sum_{j=1..m} n_j$

**Table 2** The *number of calls* complexity of query evaluation at the mediator (for the sure model) assuming  $k$  sources

Query Form	Time Complexity (wrt $ T ,  a , k, U$ )
$t$	$O( T  *  a  + k * U)$
$t_1 \vee \dots \vee t_n$	$O( T  *  a  * n + k * U)$
$t_1 \wedge \dots \wedge t_n$	$O( T  *  a  * n + k * U * n)$
$d_1 \wedge \dots \wedge d_m$ where $d_j = t_{j1} \vee \dots \vee t_{jn_j}$	$O( T  *  a  * l + k * m * U)$
$c_1 \vee \dots \vee c_m$ where $c_j = t_{j1} \wedge \dots \wedge t_{jn_j}$	$O( T  *  a  * l + k * l * U)$

**Table 3** The time complexity of query evaluation at the mediator (for the sure model)

We conclude this section by describing the evaluation of queries in the possible models of the mediator, i.e. in the models  $I_l^+$  and  $I_u^+$ . The evaluation of a single term query in  $I_l^+$  or  $I_u^+$  is done by evaluating a conjunction of terms in  $I_l^-$  or  $I_u^-$ , respectively:

$$I^+(t) = \bigcap \{ I^-(u) \mid u \in \text{head}(t) \text{ and } u \not\sim t \}$$

$$= I^-(\bigwedge \{ u \mid u \in \text{head}(t) \text{ and } u \not\sim t \})$$

where  $I^+(t)$  stands for  $I_l^+(t)$  or  $I_u^+(t)$ , and  $I^-$  stands for  $I_l^-$  or  $I_u^-$ , respectively. Therefore the complexity analysis of evaluating  $I^+(t)$  can be done using Tables 2 and 3.

Finally, the evaluation of a disjunction in  $I^+$  is done by evaluating a DNF query in  $I^-$ , and the evaluation of a conjunction in  $I^+$  is done by evaluating a conjunction in  $I^-$ .

## 5 Enhancing the Quality of Answers with Object Descriptions

We have just seen how to compute several kinds of answers at the mediator. We shall now see how to improve the “quality” of the answers by providing additional information on the objects returned.

First, let us see an example in the context of a single source. Consider a source  $S$  that contains an object 1 indexed under two terms, **Cameras** and **Underwater**, and an object 2 also indexed under two terms, **Cameras** and **Miniature**. Next, assume that  $S$  receives the query  $q = \text{Cameras}$  and is asked to return both the sure and the possible answer to that query. Clearly in both cases  $S$  will return the set  $\{1, 2\}$ . However, instead of just returning the set  $\{1, 2\}$ , the source could return the following set

$$\{(1, \{\text{Cameras}, \text{Underwater}\}), (2, \{\text{Cameras}, \text{Miniature}\})\}$$

In this set each object is accompanied by the set of *all* terms under which the object is indexed. This information could provide valuable help to the user. Indeed, the user of our example may have actually been looking for miniature cameras, but he only used the term **Cameras** in his query for one of several reasons. For example,

- the user may have forgotten to use the term **Miniature**;
- or the user did not know that the term **Miniature** was included in the terminology of the source;
- or the user did not know that the objects of the source were indexed in such specificity.

We believe that including in the answer all terms under which each object returned is indexed might aid the user in selecting the objects that are most relevant to his information need. In addition, such terms could aid the user in getting better acquainted with the taxonomy of the source. Indeed, more often than not, users are not familiar with the source taxonomy and know little about its specificity and coverage (see [51]). As a result user queries are often imprecise and do not reflect the real user needs. We believe that familiarity with the source taxonomy is essential for a precise formulation of user queries. Therefore we extend the notion of answer to be a set of objects each accompanied by its *index*, i.e. by the set all terms under which the object is indexed.

**Definition 13** The *index* of an object  $o$  with respect to an interpretation  $I$ , denoted by  $D_I(o)$ , is the set of all terms that contain  $o$  in their interpretation, i.e.  $D_I(o) = \{t \in T \mid o \in I(t)\}$ .

For brevity hereafter we shall sometimes write  $D(o)$  instead of  $D_I(o)$ ,  $D^-(o)$  instead of  $D_{I^-}(o)$ , and  $D^+(o)$  instead of  $D_{I^+}(o)$ , when the interpretation  $I$  is clear from the context. Clearly the index of an object depends on the interpretation  $I$ , so the same object can have different indexes under different interpretations. Here are some examples of indexes in the source shown in Figure 6:

$$\begin{aligned} D(1) &= \{\text{StillCams}\} \\ D^-(1) &= \{\text{StillCams}, \text{Cameras}\} \\ D^+(1) &= \{\text{StillCams}, \text{Cameras}, \text{Reflex}, \text{Miniature}, \\ &\quad \text{MovingPictureCams}, \text{UnderwaterDevices}\} \\ D(2) &= \{\text{Cameras}\} \\ D^-(2) &= \{\text{Cameras}\} \\ D^+(2) &= \{\text{Cameras}, \text{StillCams}, \text{MovingPictureCams}, \\ &\quad \text{UnderwaterDevices}\} \end{aligned}$$

We have seen earlier that the user of a source can submit a query and ask for a sure or a possible answer. Following our discussion on indexes, the user can now also ask for the sure or possible index for each object in the answer. This means that the answer returned by the source to a given query  $q$  can have one of the forms shown in Table 4. It is up to the user to specify the desired form of the answer.

	object set	object index	answer returned
1	sure	sure	$\{(o, D^-(o)) \mid o \in I^-(q)\}$
2	sure	possible	$\{(o, D^+(o)) \mid o \in I^-(q)\}$
3	possible	sure	$\{(o, D^-(o)) \mid o \in I^+(q)\}$
4	possible	possible	$\{(o, D^+(o)) \mid o \in I^+(q)\}$

**Table 4** The answers to a query

Note that if  $I^-$  is stored at the source then the evaluation of  $D^-$  for an object  $o$  is straightforward. If however only the interpretation  $I$  is stored at the source then we can compute  $D^-(o)$  as follows:

**Proposition 6**  $D^-(o) = \bigcup \{\text{head}(t) \mid o \in I(t)\}$ , or equivalently,  $D^-(o) = \bigcup \{\text{head}(t) \mid t \in D(o)\}$ .

If we have computed  $D^-(o)$ , then we can compute  $D^+(o)$  as follows:

**Proposition 7**  $D^+(o) = \{t \mid \text{head}(t) \setminus \{t' \mid t' \sim t\} \subseteq D^-(o)\}$

By analogy to the single source case, a mediator can return answers consisting of objects which are accompanied by their indexes. In other words, a mediator can return a set of pairs  $(o, D_I(o))$ , where  $I$  is the model used by the mediator for answering queries. For example consider two sources,  $S_1$  and  $S_2$ , providing information about animals (e.g. photos) as shown in Figure 15. The terms of source  $S_1$  are in English, while the

terms of source  $S_2$  are in French. Moreover a mediator  $M$  integrates the information of the two sources and provides a unified access through a taxonomy with English terms. Assume now that the mediator receives the query  $q = \text{Animal}$  in which case the mediator sends the query  $q_1 = \text{Animal} \vee \text{Dog}$  to source  $S_1$ , and the query  $q_2 = \text{Mammifère} \vee \text{Chat}$  to source  $S_2$ . Moreover, assume that the sources  $S_1$  and  $S_2$  return objects accompanied by their sure indexes. Then, the source  $S_1$  will return the answer

$$\{ (1, \{\text{Dog}, \text{Animal}\}), (2, \{\text{Canis}, \text{Dog}, \text{Animal}\}) \}$$

and the source  $S_2$  the answer

$$\{ (1, \{\text{Mammifère}\}), (3, \{\text{Chat}, \text{Mammifère}\}) \}$$

Next, assume that the mediator operates under operation mode 1 (see Table 1), that is, the mediator uses the model  $I_{l-}$  for answering queries. Moreover assume that the mediator returns objects accompanied by their sure indexes (in  $I_{l-}$ ). In this case the mediator will return the following answer:

$$\{ (1, \{\text{Dog}, \text{Mammal}, \text{Animal}\}), (2, \{\text{Dog}, \text{Mammal}, \text{Animal}\}), (3, \{\text{Mammal}, \text{Animal}\}) \}$$

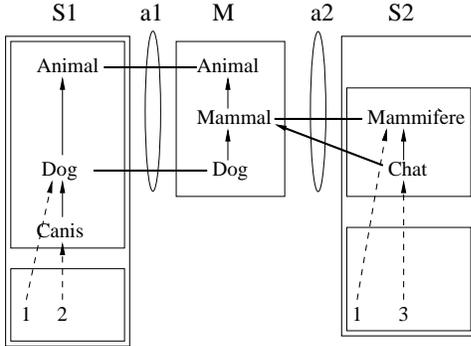


Fig. 15 A mediator over two sources

Let  $I$  denote any of the four interpretations  $I_{l-}$ ,  $I_{l+}$ ,  $I_{u-}$  and  $I_{u+}$  of the mediator, and assume that we want to compute  $D^-(o)$ , i.e. the sure index of some object  $o$ , at the mediator. Since the interpretation  $I$  is not stored at the mediator we cannot compute  $D^-(o)$  like we do for a source (see Prop. 6 above). Instead, we must exploit the articulations  $a_i$  and the indexes  $D_i(o)$  returned by the sources. Specifically, the mediator can compute  $D_l(o)$  (i.e. the index of  $o$  with respect to  $I_l$ ) and  $D_u(o)$  (i.e. the index of  $o$  with respect to  $I_u$ ) as stated by the following proposition. Note that again we write  $I_l$  instead of  $I_{l-}$  or  $I_{l+}$ , and  $I_u$  instead of  $I_{u-}$  or  $I_{u+}$ , since the computation of the object indexes at the mediator does not depend on the evaluation of queries at the underlying sources.

### Proposition 8

$$D_l(o) = \bigcup_{i=1..k} D_l^i(o), \text{ where}$$

$$D_l^i(o) = \{t \in T \mid t_i \in D_i(o) \text{ and } t_i a_i t\}$$

$$D_u(o) = \bigcup_{i=1..k} D_u^i(o), \text{ where}$$

$$D_u^i(o) = \{t \in T \mid (\text{head}_i(t) \neq \emptyset \text{ and } \text{head}_i(t) \subseteq D_i(o)) \text{ or } (\text{head}_i(t) = \emptyset \text{ and } t_i \in D_i(o) \text{ and } t_i a_i t)\}$$

Now,  $D_l^-(o)$  and  $D_u^-(o)$  can be computed by applying Prop. 6 to  $D_l(o)$  and  $D_u(o)$  respectively. Similarly,  $D_l^+(o)$  and  $D_u^+(o)$  can be computed by applying Prop. 7 to  $D_l^-(o)$  and  $D_u^-(o)$  respectively.

## 6 Extending our Model

In this section we discuss various extensions of our model. Specifically, in Section 6.1 we extend the form of our articulations, in Section 6.2 we discuss mediators which also have a stored interpretation of their terminology, in Section 6.3 we describe how our mediators can be combined with information retrieval systems, and in Section 6.4 we discuss how our approach can lead to a network of articulated sources.

### 6.1 Extending the Form of Articulations

According to Section 3 an articulation  $a_i$  consists of subsumption relationships between terms only. However we can extend the definition of an articulation to include subsumption relationships between terms and queries as well. This extension is useful, because now the designer of the mediator can define articulations containing more complex relationships, as in the following examples:

- $\text{Electronics}_M \succeq (\text{TV}_i \vee \text{Mobiles}_i \vee \text{Radios}_i)$
- $\text{DBArticles}_M \sim (\text{Databases}_i \wedge \text{Articles}_i)$

In the first example, the users of the mediator can use the term **Electronics** instead of a long disjunction of terms at source  $S_i$  (benefit: brevity), while in the second they can use the term **DBArticles** instead of the conjunction of two terms at source  $S_i$  (note, however, that this is useful only if  $S_i$  supports multiple classification).

**Definition 14** Let  $(T, \preceq)$  be the taxonomy of a mediator and let  $(T_i, \preceq_i)$  be the taxonomy of source  $S_i$ . An articulation  $a_i$  is a subsumption relation over  $T \cup Q_{T_i}$ , where  $Q_{T_i}$  is the set of all queries over  $T_i$ .

Let us now discuss the consequences of this extension with regard to the functionality of the mediators. For each term  $t \in T$  the *tail* and *head* of  $t$  with respect to  $a_i$  can be defined as follows:

**Definition 15** Given a term  $t \in T$  and articulation  $a_i$  we define

$$\text{tail}_i(t) = \{s \in Q_{T_i} \mid sa_it\} \text{ and } \text{head}_i(t) = \{u \in Q_{T_i} \mid ta_iu\}$$

Note that now the tail and head of a term are not sets of terms of  $T_i$ , but sets of *queries* over  $T_i$ .

The *lower* and *upper* approximation of  $t$  with respect to  $a_i$  are defined as in Section 3. The four interpretations and the eight answer models of the mediator are defined in the same way too.

In this framework the concept of compatibility is now redefined as follows:

**Definition 16** A source  $S_i$  is *compatible* with the mediator  $M$  if for any queries  $s, t$  in  $Q_{T_i}$ , if  $sa_it$  then  $s \preceq_i t$ .

As mentioned earlier, maintaining compatibility is not an easy task. The mediator should (periodically) check the compatibility of its sources, e.g. by submitting to them queries allowing to check whether  $t \preceq_i t'$ . However, now  $t$  and  $t'$  are queries, thus the sources should support subsumption checking over queries.

In general, an articulation may contain relationships between terms and arbitrary queries. For example, consider a source  $S_i$  implemented in the relational model (as described in section 2) and suppose that this source can answer only pure SQL queries. In this case the articulation  $a_i$  may contain relationships of the form  $\text{Cameras} \preceq_{a_i} q_i$  where

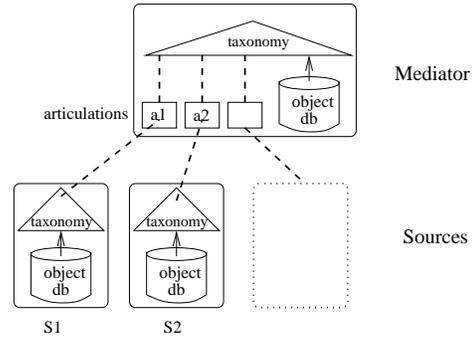
$$q_i = \Pi_{\text{object}}(\sigma_{\text{term-name}=\text{"Cameras"}}(\text{INTERPRETATION} \bowtie \text{TERMINOLOGY}))$$

## 6.2 Mediators with Stored Interpretations

We can easily extend a mediator so as to *also* store an interpretation of its terminology  $T$ . Figure 16 shows graphically the architecture of a mediator of this kind. Such an extension can prove quite useful in the context of the Web: a Web user can define his own mediator consisting of a taxonomy that is familiar to him, a set of articulations to other Web catalogues, *and* a stored interpretation of the mediator's taxonomy. Note that the taxonomy of the mediator and its stored interpretation resembles the bookmarks facility of Web browsers. However, the addition of articulations now allows the user to browse and query remote catalogues.

Let  $I_M$  denote the stored interpretation of  $T$ . When a user sends a query to the mediator he has three choices:

- he can ask for an answer derived from  $I_M$ ,
- he can ask for an answer derived from the interpretations of the remote sources, or
- he can ask for an answer derived from both  $I_M$  and the interpretations of the remote sources.



**Fig. 16** The architecture of a mediator with a stored interpretation

In the first case the mediator operates as a source (see section 2), in the second case it operates like the mediators described earlier, while in the third case it again operates like the mediators described earlier but with one difference: the interpretations  $I_{l-}$ ,  $I_{l+}$ ,  $I_{u-}$ ,  $I_{u+}$  are now defined by taking the union of  $I_M$  and the interpretations of the sources. For instance the interpretation  $I_{l-}$  is now defined as:

$$I_{l-}(t) = I_M(t) \cup \left( \bigcup_{i=1}^k I_i^-(t_i) \right)$$

In other words, in the third case, the mediator operates as usual, except that now, in addition to the  $k$  external sources  $S_1, \dots, S_k$ , we have the mediator's own source  $S_M = \langle (T, \preceq), I_M \rangle$  acting as a  $(k+1)$ -th source.

In the case where the mediator also stores an interpretation  $I_M$  of  $T$  then the mediator's ability to "translate" the descriptions of the objects returned by the underlying sources drives to an interesting scenario for the Web. Consider a user who has submitted a query to the mediator and assume that the mediator has returned a set of objects to the user. If some of these objects are of real interest to the user (e.g. a set of beautiful images, good papers etc) then the user can store these objects in the database of the mediator. These objects will be stored under terms of the mediator's taxonomy, i.e. in the interpretation  $I_M$  of  $T$ .

For example, consider the mediator shown in Figure 15. The mediator can store objects 1 and 2 under the terms **Dog**, **Mammal** and **Animal**, and object 3 under the terms **Mammal** and **Animal**

However one can easily see that it suffice to store the objects 1 and 2 under the term **Dog** and the object 3 under the term **Mammal**. More formally for storing an object  $o$  in  $I_M$  the mediator associates this object with the following terms of  $T$ :

$$\min_{\preceq_M} D_l(o) \quad \text{or} \quad \min_{\preceq_M} D_u(o)$$

### 6.3 Mediators Over Hybrid Sources

Let us use the term *free retrieval source* to refer to a source that indexes the objects of interest using an *uncontrolled* vocabulary. In this case, the objects of the domain have textual content and the vocabulary that is used for indexing them consists of those words that appear in the objects. These sources usually accept natural language queries and return a set of objects ordered according to their relevance to the query. Text retrieval systems (the typical case of “Information Retrieval systems”), as well as the search engines of the Web, fall into this category.

We can now use the term *hybrid source* to refer to a source that is both taxonomy-based and free retrieval source. A hybrid source accepts *two* kinds of queries: queries over a controlled vocabulary and natural language queries. A source whose functionality moves towards this direction is Google. Using Google, one can first select a category, e.g. `Sciences/CS/DataStructures`, from the taxonomy of Open Directory and then submit a natural language query, e.g. “Tree”. The search engine will compute the degree of relevance with respect to the natural language query, “Tree”, only of those pages that fall in the category `Sciences/CS/DataStructures` in the catalog of Open Directory. Clearly, this enhances the precision of the retrieval and is computationally more economical.

Our approach can be used for building mediators over hybrid sources whose functionality extends the functionality offered by the existing meta-searchers of the Web (e.g. MetaCrawler [63], SavvySearch [39], Profusion [28]). The user of a hybrid mediator can use the taxonomy of the mediator in order to browse or query those parts of the sources that are of interest to him. Moreover he is able to query the databases of these sources using natural language queries. This implies that the mediator will send two kinds of queries to the sources: (a) queries which are evaluated based on the indexing of the objects with respect to the taxonomy of the source, and (b) queries which are evaluated based on the contents of the objects (pages). Figure 17 describes this architecture graphically.

The functionality of the mediators described in this paper presumes that each source can provide a sure answer and a possible answer. However, the taxonomy-based sources that can be found in the Web, e.g. Yahoo! or ODP, do not currently provide such answers. This means that the functionality of our mediators cannot be implemented straightforwardly. Nevertheless, we can implement the functionality of our mediators over such sources by employing appropriate *wrappers*.

First note that the taxonomy and the interpretation of a Web catalog is published as a set of Web pages. For each term  $t$  of the taxonomy there is a separate Web page. This page contains the name of the term, and links

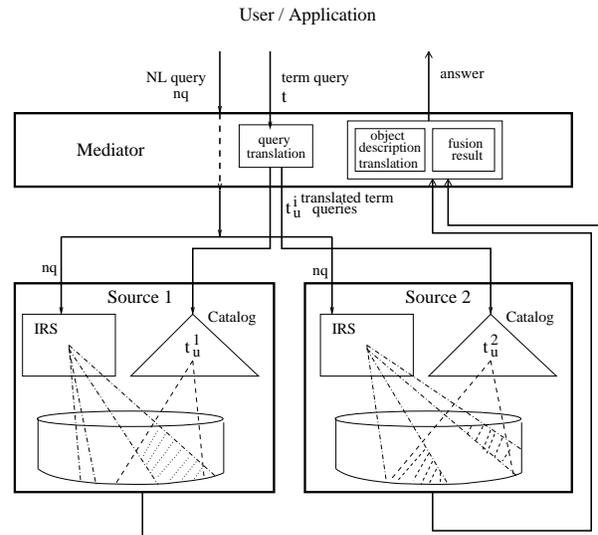


Fig. 17 Building mediators over hybrid sources

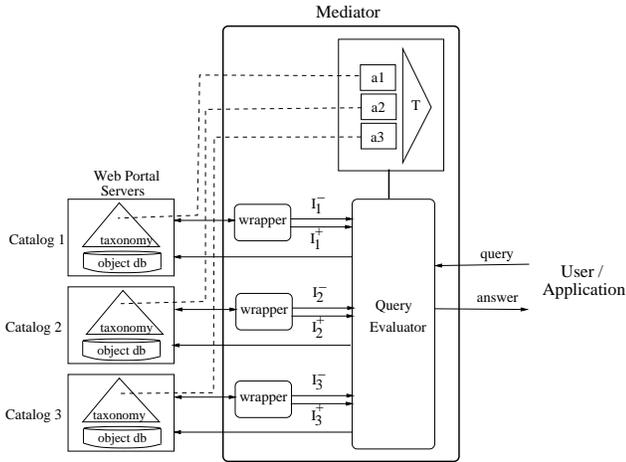
pointing to pages which correspond to the terms which are subsumed by  $t$ . In addition, the page contains links pointing to the objects, here Web pages, which have been indexed under the term  $t$ . However we can employ a wrapper in order to parse each such page and extract the name of the term, the subsumed terms and the indexed objects.

Now, the architecture for implementing our mediators over Web catalogues is shown in Figure 18. The key point is that the interpretation of a term  $t$  of a source  $S_i$  in the sure model  $I_i^-$  and in the possible model  $I_i^+$  can be computed at the mediator side. This can be achieved by building an appropriate wrapper for that source. In particular, for computing  $I_i^+(t)$  the wrapper will fetch the pages of all terms  $t'$  such that  $t \preceq_i t'$  and then it will derive  $I_i^+(t)$  by computing the intersection  $\cap \{I_i^-(t') | t \preceq t'\}$ . According to this architecture our mediators can be implemented by using the standard HTTP protocol. A prototype version of our mediators over the Web has already been implemented at Université de Paris-Sud.

### 6.4 Networks of Articulated Sources

One can easily see how our approach can be used for creating a complex information network, comprising sources and mediators, in a natural and straightforward manner. Indeed,

- in order to add a mediator to such a network one has to (a) design the mediator taxonomy  $(T, \preceq)$  based on the domain of interest, (b) select the sources to be mediated, and (c) design the articulations  $a_i$  based on the known/observed relations between terms of the mediator and terms of the selected sources;
- in order to remove a mediator from the network one just has to disconnect the mediator from the network;



**Fig. 18** An architecture for implementing our mediators over the catalogues of the Web

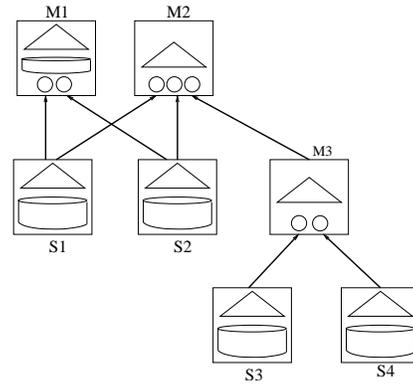
- moreover, in order to add a source to the information network all one has to do is (a) select one or more mediators in the network and (b) design an articulation between the source and each mediator;
- finally, to remove a source from the network one simply has to remove the corresponding articulation from the mediator(s) to which the source is connected and disconnect the source. Note that as each mediator has one articulation for each underlying source, the deletion of an articulation does not affect the rest of the articulations.

A significant consequence of this approach is that network evolution can be *incremental*. Indeed, new relationships between terms of the mediator and terms of the sources can be added with minimum effort as soon as they are observed, and relationships that are seen to be no more valid can be removed just as easily: simply add/remove the relationships at the appropriate articulation in the mediator database storing the articulation.

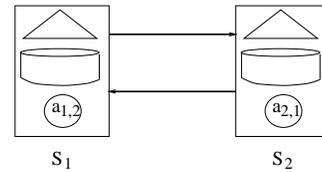
For example, Figure 19 shows a network consisting of four primary sources  $S_1, S_2, S_3, S_4$  and three mediators  $M_1, M_2$  and  $M_3$ . A line segment connecting a mediator  $M$  to a source or to a mediator  $S$ , means that  $M$  is a mediator over  $S$ , and circles denote articulations. For example  $M_2$  is a mediator over the sources  $S_1$  and  $S_2$  and the mediator  $M_3$ . Note that mediator  $M_1$  can be also considered as a primary source, because it has a stored interpretation.

Also note that our approach allows *mutually articulated* sources as shown in Figure 20. In this case we can no longer distinguish sources into primary and secondary.

Query evaluation and updating in a network of articulated sources raises several interesting questions. For example, a query to a source may trigger an infinite number of calls between the sources, if the network is cyclic, e.g. like the one shown in Figure 20. This and other re-



**Fig. 19** A network consisting of primary and secondary taxonomy-based sources



**Fig. 20** Mutually articulated sources

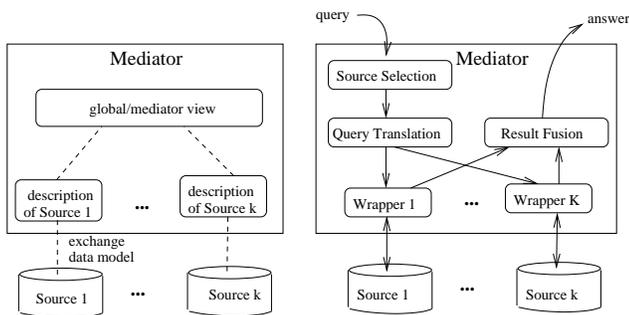
lated problems go beyond the scope of this paper, and are treated in [71, 70].

## 7 Related Work

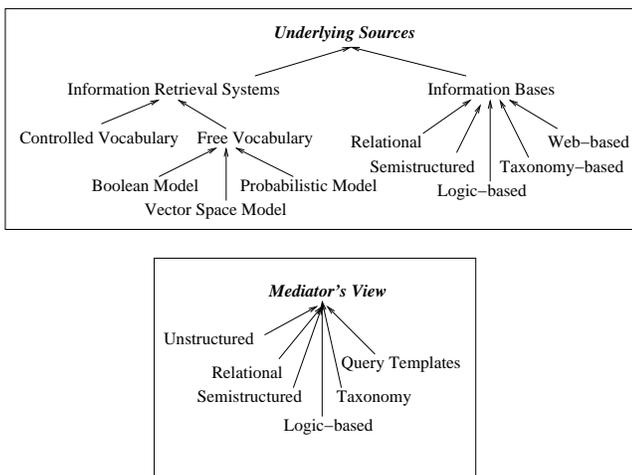
The need for integrated and unified access to multiple information sources has stimulated the research on *mediators*. The concept of mediator was initially proposed by Wiederhold [78]. Since then it has been applied and specialized for several kinds of sources and applications needs. Nevertheless, in every instance we can identify a number of basic architectural components. Specifically, in most of the cases the mediator architecture consists of a mediator’s view (usually in the form of a conceptual model), *source descriptions* and *wrappers* that describe the contents and/or the querying capabilities of each source with respect to the mediator’s view, and an *exchange data model*, which is used to convey information between the mediator and the sources. The mediator accepts queries expressed in the mediator view. Upon receiving a user query, the mediator selects the sources to be queried and formulates the query to be sent to each one of them. These tasks are accomplished based on the source descriptions that encode what the mediator “knows” about the underlying sources. Finally, the mediator appropriately combines the returned results and delivers the final answer to the user. Figure 21 shows the general architecture and the functional overview of the mediator.

In this section we compare our approach with other existing mediator approaches. Our objective is to iden-

tify the basic differences and analogies and identify issues that worth further research, rather than presenting a complete survey of this very wide area. Figure 22 shows a rough taxonomy of the mediator approaches according to two criteria: (a) the kind of the underlying sources, and (b) the kind of the mediator’s view. According to the first, we can partition sources into two broad categories: *information retrieval systems* and *information bases*. The former provide content-based access to a set of (text) documents, while the latter store structured data. We will use this dichotomy in our sequel discussion.



**Fig. 21** Architecture and functional overview of the mediator



**Fig. 22** A rough taxonomy of the mediator approaches

### 7.1 Mediators over Information Retrieval Systems

Information Retrieval Systems (IRS) provide content-based access to a set of (text) documents. The content of the documents (as well as the user queries) is described using an “indexing language” which can be either

- (a) a “free” vocabulary consisting of the words that appear in the documents of the collection, excluding those words that carry no information (such as articles) and reducing words to their grammatical root (a task called “stemming”), or
- (b) a “controlled” vocabulary which may be different from the set of words that appear in the documents. This vocabulary may be structured by a small set of relations like hyponymy and synonymy.

For the relative merits of each of these approaches see [62, 22]. One could say that taxonomy-based sources resemble those IRSs which employ the boolean retrieval model (see [5] for a review) and exploit lexical ontologies or word-based thesauri [40] (like WordNet [19] and Roget’s thesaurus) for query expansion, i.e. for expanding queries with synonyms, hyponyms and related terms in order to improve recall (e.g. see [56], [50], [51], [36]). However note that the IRS techniques are applicable only if the objects of the domain have a textual content (this is not a prerequisite of our approach). Another remarkable difference with our sources is that the taxonomies employed by IRSs usually do not accept the semantic interpretation that we describe in this paper. Lexical ontologies like WordNet [19] are structured using lexical relations (synonymy, hyponymy, antonymy) which are not semantic relations. For instance, according to Wordnet, *window* is subsumed by *opening* and by *panel*. However, every *window* is not a *panel* and an *opening*, thus extensional subsumption does not hold here. The justification of the possible answer (in sources) and the eight answer models (in mediators) does not apply to such ontologies (for more about this problem see [35]). Instead, techniques like spreading activation [56] are more appropriate if lexical ontologies are employed.

By consequence, our approach is quite different from the mediator approaches that have emerged within the IR community. Specifically, a mediator over IRSs that employ free vocabularies, does not have to translate the mediator queries, as each source accepts the same set of queries, i.e. natural language queries. As a consequence, mediators over such sources, mainly focus on issues like *source selection* and *result fusion* (meta-ranking) (e.g. see [76, 77, 11, 34, 29, 21, 6, 67]). On the other hand, mediators over systems which employ controlled and structured vocabularies have not received adequate attention until now. To the best of our knowledge, all of the existing approaches focus on ontology merging and not on ontology articulation. Moreover, as they mainly employ lexical ontologies the mappings between two ontologies consist of lexical relationships too (in many cases one term is associated with a set of terms of the other ontology [3]).

Although the controlled indexing languages that are used for information retrieval usually consist of a set of terms structured by a small number of relations (such as hyponymy and synonymy), there are cases where the indexing of the objects is done (especially in the case of

a manual indexing process) with respect to more expressive conceptual models representing domain knowledge in a more detailed and precise manner. Such conceptual models can be represented using logic-based languages, and the corresponding reasoning mechanisms can be exploited for retrieving objects. There are several works that take this conceptual modeling and reasoning approach to information retrieval (e.g. relevance terminological logics [52], four-valued logics [59]). This conceptual modeling approach is useful and effective if the domain is narrow. If the domain is too wide (e.g. the set of all Web pages) then the problem is that it is hard to conceptualize the domain; actually there are many different ways to conceptualize it, so it is hard to reach a conceptual model of wide acceptance. Thus a mediator over such sources has to tackle complex structural differences (recall the example of Section 2). For this purpose, even today, ontologies that have simple structure, like the one that we consider, are usually employed for retrieving objects from large collections of objects ([61]).

## 7.2 Mediators over Information Bases

We use the term *information bases* to refer to sources that store structured data, not documents. Relational, semistructured, logic-based and Web-based sources belong to this category. Indeed, there are several approaches for building mediators over relational databases (e.g. see [48,31,32,79]), SGML documents (e.g. see [17]), and Web-based sources (e.g. see [4,14,15]). We include this discussion here because our sources can be considered as information bases as we do not presuppose that the objects of the domain have a textual content.

Concerning the kind of the mediator's view, several approaches have come up (as the rightmost taxonomy of Figure 22 illustrates). Indeed, we have seen mediators whose unified view has the form of a relational schema (e.g. Infomaster [33,26]), a semantic network (e.g. SIMS [43]), a F-logic schema (e.g. OntoBroker [7,23]), a Description Logics schema (e.g. Information Manifold [48], OBSERVER [53,41,42], PICSEL [46]), a set of query templates (e.g. TSIMMIS [16,31], HERMES [66]). Furthermore, several data models have been used for conveying information between the mediator and the sources (as shown in Figure 21) including relational tuples (e.g. in Infomaster, SIMS, Information Manifold, OBSERVER), tuples that encode graph data structures (like the OEM in TSIMMIS, or the YAT in [17]), and HTML pages (in Web-oriented approaches).

In order identify similarities, differences and analogies between the above approaches and the one presented in this paper, we first describe a set of layers from which we can view an information base, then we use these layers in order to discuss the kinds of heterogeneity that

may exist between two information bases, and finally, we give a number of remarks.

## A layered view of an information base

We could view a source at five different layers: the *domain*, the *conceptualization*, the *conceptual model*, the *data model*, and the *query language*. We are aware that these distinctions are not crystal clear or widely accepted, however they enable us to discuss systematically a number of issues and draw analogies. There are dependencies among these layers as shown in Figure 23, e.g. the query language layer of a source depends on the data model layer of the source, and so on.

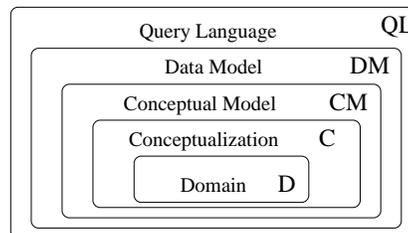
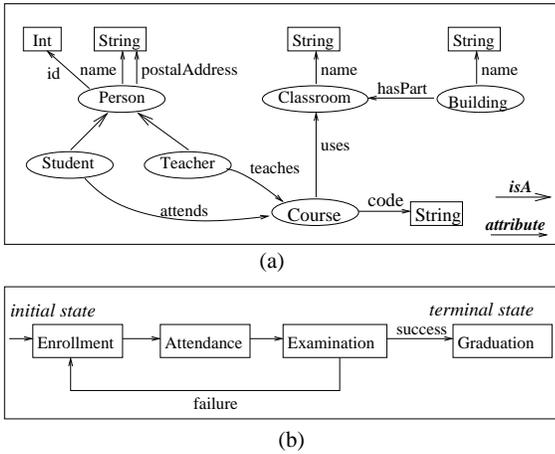


Fig. 23 The layers of a source

Each source stores information about a part of the real world which we call the *domain* layer of the source. For example, the domain of a source can be the set of all URLs, or the set of all universities, or the set of Greek universities, or the Computer Science Department of the University of Crete (that we call *CSD domain* in the sequel). The *conceptualization* of a domain is the intellectual lens through which the domain is viewed. For example, one conceptualization of the CSD domain may describe its *static* aspects, i.e. what entities or things exist (e.g. persons, buildings, classrooms, computers), their attributes and their interrelationships. Another conceptualization may describe its *dynamic* aspects in terms of states, state transitions and processes (e.g. enrollments, graduations, attendances, teaching). A *conceptual model* is used to describe a particular conceptualization of a domain in terms of a set of (widely accepted) structuring mechanisms which are appropriate for the conceptualization. For example, a conceptual model that describes the static aspects of the CSD domain, using generalization and attribution, is shown in Figure 24.(a), while a conceptual model that describes the dynamic aspects of the CSD domain, using states and state transitions, is shown in Figure 24.(b). The representation of a conceptual model in a computer is done according to a specific *data model* (e.g. relational, object-oriented, semantic network-based, semistructured). For example, the class `Person` of the conceptual model of Figure 24 can be represented in the relational model by a relation scheme as follows: `Person(id: Int, name: Str, postalAddress: Str)`. Alternatively, in a different source,

it could be also represented using two relation schemes:  $\text{PERSON}(\text{id}:\text{Int}, \text{name}:\text{Str}, \text{addressId}:\text{Int})$  and  $\text{POSTALADDRESS}(\text{id}:\text{Int}, \text{address}:\text{Str})$ . However, there are also some data models that allow a straightforward representation of the conceptual model, e.g. the semantic network-based data model of SIS-Telos [44,20]. Finally, each source can answer queries expressed in a particular query language. For example, a source may respond to Datalog queries, while another may respond only to SQL queries. In this case we say that the *query language layers* of these sources are different.



**Fig. 24** Two conceptual models of the CSD domain: one for the *static* and one for the *dynamic* aspects.

### Kinds of heterogeneity

Given a source  $S_i$ , we will use  $D_i$  to denote the domain,  $C_i$  the conceptualization,  $CM_i$  the conceptual model,  $DM_i$  the data model, and  $QL_i$  the query language layer of  $S_i$ . Consider now two sources  $S_1$  and  $S_2$ . We may have several forms of heterogeneity between these sources, specifically there are  $2^5 = 32$  different cases (due to the 5 layers). For example, the case  $D_1 = D_2$ ,  $C_1 = C_2$ ,  $CM_1 \neq CM_2$ , means that  $S_1$  and  $S_2$  have the same conceptualization of the (same) domain, but they employ different conceptual models. Even if the conceptual models are expressed using the same structuring mechanisms (e.g. generalization, attribution), they may differ due to:

- *different naming conventions* (also called naming conflicts). A frequent phenomenon is the presence of homonyms and synonyms.
- *different scaling schemes*. They occur when different reference systems are used to measure a value, e.g. 1 foot vs 0.304 meter,  $23^\circ\text{C}$  vs  $73^\circ\text{F}$ .
- *different levels of granularity*. For example  $CM_1$  may contain only a class **Cameras**, while  $CM_2$  may contain the classes **StillCameras** and **MovingPictureCameras**.

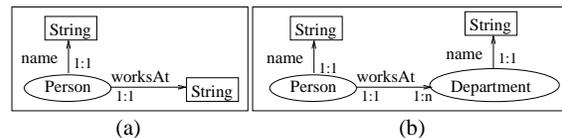
- *structural differences*, e.g.  $CM_1$  may contain a class **Person** having an attribute **owns** with range the class **ArtificialObject**, and a class **Car** defined as a specialization of the class **ArtificialObject**, while  $CM_2$  may contain a class **Car** having an attribute **owner** with range the class **Person**.

As another example, the case  $D_1 = D_2$ ,  $C_1 = C_2$ ,  $CM_1 = CM_2$ ,  $DM_1 \neq DM_2$ , means that  $S_1$  and  $S_2$  have same conceptual model but these models are represented differently in the data models layer. Note that even if  $S_1$  and  $S_2$  employ the same data model, e.g. the relational,  $DM_1$  and  $DM_2$  may differ in that they represent the conceptual model differently.

### Remarks and Analogies

Let us now give some remarks and discuss some analogies between our approach and other mediator approaches that have emerged.

- An important remark is that, given an existing source, we usually have in our disposal only its data model and query language layer, and more often than not, from these two layers we cannot infer the conceptual model or the conceptualization layer of the source. For example, consider the following relation scheme:  $\text{PERSON}(\text{name}:\text{Str}, \text{worksAt}:\text{Str})$ . The underlying conceptual model could be any of the ones shown in Figure 25, as the translations of both (a) and (b) to the relational model (by using an algorithm such as the one described in [9]) are identical. Note that according to (a) the domain consists of entities of one kind, i.e. persons, while according to (b) the domain consists of two kinds of entities: persons and departments. Moreover, although two sources may have the same conceptual model, e.g. the conceptual model (a), their representation in the data model may differ. For example the conceptual model (a) could be represented in the relational model by one relation scheme (as we saw before), or by the following two relation schemes:  $\text{PERSON}(\text{name}:\text{Str}, \text{worksAt}:\text{Int})$  and  $\text{DEP}(\text{depId}:\text{Int}, \text{name}:\text{Str})$ . We believe that this is the basic reason why information integration is a difficult and laborious task.



**Fig. 25** Two conceptual models of the CSD domain

- According to the layered view described above, the sources of our mediators (a) may have different domains (i.e. may index different sets of objects), (b) conceptualize their domains similarly (i.e. all  $C_i$  are

denumerable sets of objects), (c) may have different conceptual models (i.e. different taxonomies), (d) may have different query languages (recall the remark at the end of Section 6.1).

- In relational mediators (see [32] for a review) the mediator view is represented as a relational database schema. Relational mediators have some critical differences with our mediators. Relational mediators and their sources are *schema-based* while our mediators and their sources are *taxonomy-based*. Also recall that the relational model is value-based, not object-based. This implies that the conceptualization and the conceptual model of a relational source is hidden, or unclear. Therefore mediators over such sources “work” on the data model layer. Instead, we propose a totally different conceptual modeling approach for both sources and mediators.

Concerning source descriptions, we can distinguish the *local-as-view* (LAV) and the *global-as-view* (GAV) approach (see [13,47] for a comparison). In the LAV approach the source relations are defined as relational views over the mediator’s relations, while in the GAV approach the mediator relations are defined as views of the source relations. The former approach offers flexibility in representing the contents of the sources, but query answering is “hard” because this requires answering queries using views ([25],[37], [74]). On the other hand, the GAV approach offers easy query answering (expansion of queries until getting to source relations), but the addition/deletion of a source implies updating the mediator view, i.e. the definition of the mediator relations. It worths mentioning here that as our articulations contain relationships between single terms these kinds of mappings enjoy the benefits of both GAV and LAV approaches, i.e. they have (a) the query processing simplicity of the GAV approach, as query processing basically reduces to unfolding the query using the definitions specified in the mapping, so as to translate the query in terms of accesses (i.e. queries) to the sources, and (b) the modeling scalability of the LAV approach, i.e. the addition of a new underlying source does not require changing the previous mappings. On the other hand, term-to-query articulations (presented in Section 6.1) resemble the GAV approach.

Concerning the translation facilities, relational mediators attempt to construct *exact translations* of SQL queries while our mediators allow *approximate translations* of boolean expressions through their articulations. We might say that the answers returned by a relational mediator, correspond to the answers returned by a taxonomy-based mediator in the  $I_L^-$  model.

Moreover, in several approaches (e.g. in Infomaster) a predicate corresponding to a source relation, can appear only in the head or in the tail of a rule.

This means that granularity heterogeneities cannot be tackled easily.

- A different approach to mediators can be found in [12] which presents the fundamental features of a declarative approach to information integration based on Description Logics. The authors describe a methodology for integrating relational sources and they resort to very expressive logics in order to bridge the heterogeneities between the unified view of the mediator and the source views. However the reasoning services for supporting translations have exponential complexity, as opposed to the complexity of our mediators which is polynomial. In addition, the eight possible answers of our approach allow providing a novel query relaxation facility.
- One difference between our approach and the system OBSERVER [53,41,42] is that OBSERVER requires merging the ontologies of all underlying sources. Instead, we just articulate the taxonomies of the sources with the taxonomy of the mediator. Moreover, the compatibility condition introduced here allows the mediator to draw conclusions about the structure of a source taxonomy without having to store that taxonomy.
- In the approximate query mapping approach of [14, 15] the translated queries minimally subsume the original ones. However, the functionality offered by our mediators is different, firstly because we support negation while they do not, and secondly because our mediators support multiple operation modes, one of which is the case where the translated queries subsume the original ones.
- An alternative solution to the problem of query relaxation in mediators is the *query repairing* described in [8]. If the submitted query yields no answer then the mediator provides to the user an answer to a “similar” query. The selection of this query is based on a measure of similarity between the concepts and the predicates which is based on the taxonomic structure of the mediator’s ontology. According to our opinion, the eight answer models of our mediators offer a better founded approach to query relaxation.

## 8 Concluding Remarks

We have presented an approach for providing uniform access to multiple taxonomy-based sources through mediators that render the heterogeneities (naming, contextual, granularity) of the sources transparent to users. This paper integrates and extends the work presented in [72] and [73] and it was inspired by the approach presented in [65].

A user of the mediator, apart from being able to pose queries in terms of a taxonomy that was not used to index the objects of the sources being searched, gets an answer comprised of objects which are accompanied by

descriptions over the mediator’s taxonomy. A mediator is seen as just another source but *without* stored interpretation. An interpretation for the mediator is defined based on the interpretations stored at the sources and on the *articulations* between the mediator and the sources; and in fact, we have seen *eight* different ways for defining a mediator interpretation depending on the nature of the answers that the mediator provides to its users (see Table 1). Since the resulting mediator models are ordered they can be used in order to support a form of *query relaxation*.

Articulations can be defined by humans, but they can also be constructed automatically or semi-automatically in some specific cases, following a model-driven approach (e.g. [64,2,54]) or a data-driven approach (e.g. [3,38,24,45,60,69]).

The distinctive features of our approach are the following:

- We assume that all sources have the same *domain* and the *same* conceptualization of that domain. The intended domain is the Web, and each source views the Web as a set of objects *Obj* (URLs), and stores information about a subset of it (i.e.  $O_i \subseteq Obj$ ). This means that each object has a *unique identity* over all sources. From this point of view, we could call our mediators object-oriented as opposed to mediators over relational sources, which we could call value-oriented.
- We consider that the conceptual layer of each source is a triple  $(T, \preceq, I)$ . This conceptual modeling approach has two main advantages: (a) it is easy to create the conceptual model of a source or a mediator, and (b) the integration of information from multiple sources can be done easily. Indeed, articulations offer a *uniform* method to bridge naming, contextual and granularity heterogeneities between the conceptual models of the sources. Given this conceptual modeling approach, the mediator does not have to tackle complex structural differences between the sources (as happens with relational mediators). Moreover, it allows the integration of *schema* and *data* in a uniform manner. For example consider a source  $S$  having the conceptual model shown in Figure 3.(a), and a source  $S'$  having the conceptual model shown in Figure 3.(b), and suppose that both sources are implemented in the relational model. In source  $S$  the concept *wood* will be represented at the data level (it would be an element of the domain of an attribute), while in  $S'$  it would be a relation. Furthermore, this approach makes the automatic construction of articulations feasible [69].

Summarizing, the taxonomy-based mediation approach presented here offers the following advantages:

- *Easy construction of mediators*  
A mediator can be easily constructed even by ordi-

nary Web users. Indeed, the simple conceptual modeling approach that we adopt makes the definition of the mediator’s taxonomy and articulations very easy.

- *Query Relaxation*  
Often a query to a mediator yields no answer. The sure and the possible answers of sources, as well as the several modes of operation of a mediator, offer a solution to this problem.
- *Efficient Query Evaluation*  
The time complexity of query translation at the mediator is linear with respect to the size of the subsumption relations of the mediator.
- *Scalability*  
Articulation (instead of merging) enables a natural, incremental evolution of a network of sources. The taxonomies employed by Web catalogues contain very large numbers of terms (e.g. the taxonomy of Open Directory contains 450.000 terms). Therefore the *articulation* of taxonomies has several advantages compared to taxonomy *merging*. First, merging would introduce storage and performance overheads. Second, full merging is a laborious task which in many cases does not pay-off because the integrated taxonomy becomes obsolete when the taxonomies involved change. Another problem with full merging is that it usually requires full consistency, which may be hard to achieve in practice, while articulation can work on locally consistent parts of the taxonomies involved. However, note that the taxonomies considered here present no consistency problems. There may only be long cycles of subsumption relationships, which induce big classes of equivalent terms.
- *Applicability*  
The taxonomy-based approach presented provides a flexible and formal framework for integrating data from several sources, and/or for personalizing the contents of one or more sources. The taxonomies considered fit quite well with the content-based organizational structure of Web catalogues (e.g. Yahoo!, Open Directory), keyword hierarchies (e.g. ACM’s thesaurus), XFML [1] taxonomies and personal bookmarks. By defining a mediator, the user can employ his own terminology in order to access and query several Web catalogues, specifically those parts of the catalogues that are of interest to him. Moreover, as a mediator can also have a stored interpretation, our approach can lead to a network of articulated sources. Recall that a mediator can translate the descriptions of the objects returned by the underlying sources. This implies that all (or some) of these objects can be straightforwardly stored in the mediator base (under terms of the mediator taxonomy).

An interesting line of research is to investigate query evaluation and updating in a network of articulated sources. Another interesting issue is to investigate how the me-

diator can exploit the object indexes that are returned by a compatible source in order to check whether that source remains compatible. If we consider sources that answer queries by returning an ordered set of objects, the mediator should also return ordered sets of objects. It would then be interesting to investigate whether the work presented in this paper can be integrated with the work presented in [67] and [27].

*Acknowledgements* Part of this work was conducted while the first two authors were visiting with Meme Media Laboratory, Hokkaido University, Sapporo, Japan. Many thanks to Professor Tanaka, director of the Meme Media Laboratory, and to the University of Hokkaido for their hospitality. We also want to thank the anonymous referees for their comments that improved this paper.

## References

1. “XFML: eXchangeable Faceted Metadata Language”. <http://www.xfml.org>.
2. Bernd Amann and Iriini Fundulaki. “Integrating Ontologies and Thesauri to Build RDF Schemas”. In *Proceedings of the Third European Conference for Digital Libraries ECDL’99*, Paris, France, 1999.
3. S. Amba. “Automatic Linking of Thesauri”. In *Proceeding of SIGIR’96*, Zurich, Switzerland, 1996. ACM Press.
4. José Luis Ambite, Naveen Ashish, Greg Barish, Craig A. Knoblock, Steven Minton, Pragnesh J. Modi, Ion Muslea, Andrew Philpot, and Sheila Tejada. “Ariadne: a system for constructing mediators for Internet sources”. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 561–563, 1998.
5. Ricardo Baeza-Yates and Berthier Ribeiro-Neto. *Modern Information Retrieval*. ACM Press, Addison-Wesley, 1999.
6. Christoph Baumgarten. *Probabilistic Information Retrieval in a Distributed Heterogeneous Environment*. PhD thesis, Technical University of Dresden, February 1999.
7. V. Richard Benjamins and Dieter Fensel. “Community is Knowledge! in (KA)<sup>2</sup>”. *Proceedings of KAW’98*.
8. Alain Bidault, Christine Froidevaux, and Brigitte Safar. “Repairing Queries in a Mediator Approach”. In *Proceedings of the ECAI’02*, Lyon, France, 2002.
9. Magnus Boman, Janis A. Bubenko, Paul Johannesson, and Benkt Wangler. *Conceptual Modelling*. Prentice-Hall, 1997.
10. George Boolos. *Logic, Logic and Logic*. Harvard University Press, 1998.
11. J. P. Callan, Z. Lu, and W. B. Croft. “Searching Distributed Collections with Inference Networks”. In *18th Int. Conf. on Research and Development in Information Retrieval*, 1995.
12. D. Calvanese, G. de Giacomo, M. Lenzerini, D. Nardi, and R. Rosati. “Description Logic Framework for Information Integration”. In *Proceedings of the 6th Int. Conf. on the Principles of Knowledge Representation and Reasoning (KR-98)*, 1998.
13. Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini. A framework for ontology integration. In *Proc. of the 2001 Int. Semantic Web Working Symposium (SWWS 2001)*, pages 303–316, 2001.
14. Chen-Chuan K. Chang and Héctor García-Molina. “Mind Your Vocabulary: Query Mapping Across Heterogeneous Information Sources”. In *Proc. of the ACM SIGMOD*, pages 335–346, 1999.
15. Chen-Chuan K. Chang and Héctor García-Molina. “Approximate query mapping: Accounting for translation closeness”. *VLDB Journal*, 10(2-3), 2000.
16. Sudarshan Chawathe, Hector Garcia-Molina, Joachim Hammer, Kelly Ireland, Yannis Papanikolaou, Jeffrey Ullman, and Jennifer Widom. “The TSIMMIS project: Integration of Heterogeneous Information Sources”. In *Proceedings of IPSJ*, Tokyo, Japan, October 1994.
17. Sophie Cluet, Claude Delobel, Jérôme Siméon, and Katarzyna Smaga. “Your mediators need data conversion!”. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 1998.
18. E. F. Codd. “A Relational Model of Data for Large Shared Data Banks”. *Communications of the ACM*, 13(6):377–387, 1970.
19. Princeton University Cognitive Science Laboratory. “WordNet: A Lexical Database for the English Language”. (<http://www.cogsci.princeton.edu/wn>).
20. Panos Constantopoulos, Martin Doerr, and Yannis Vassiliou. “Repositories for Software Reuse : The Software Information Base”. In *Proceedings IFIP WG 8.1 Conference on Information System Development Process*, pages 285–307, Como, Italy, September 1993.
21. N. Craswell, D. Hawking, and P. Thistlewaite. “Merging Results from Isolated Search Engines”. In *Proceedings of the Tenth Australasian Database Conference*, 1999.
22. Bruce Croft. “Knowledge-based and Statistical Approaches to Text Retrieval”. *IEEE Expert*, 9:8–12, April 1993.
23. S. Decker, M. Erdmann, D. Fensel, and R. Studer. “Ontobroker: Ontology based Access to Distributed and Semi-Structured Information”. In *Semantic Issues in Multimedia Systems*. Kluwer Academic Publisher, 1999.
24. A. Doan, J. Madhavan, P. Domingos, and A. Halevy. “Learning to Map between Ontologies on the Semantic Web”. In *Proceedings of the World-Wide Web Conference (WWW-2002)*, 2002.
25. Oliver M. Duschka and Michael R. Genesereth. “Answering Recursive Queries Using Views”. In *Proceedings of AAAI, 1997.*, 1997.
26. Oliver M. Duschka and Michael R. Genesereth. “Query Planning in Infomaster”. In *Proceedings of the Twelfth Annual ACM Symposium on Applied Computing, SAC’97*, San Jose, February 1997.
27. Ronald Fagin. “Combining Fuzzy Information From Multiple Systems”. *Journal of Computer and System Sciences*, 58, 1999.
28. Yizhong Fan and Susan Gauch. “Adaptive Agents for Information Gathering from Multiple, Distributed Information Sources”. In *1999 AAAI Symposium on Intelligent Agents in Cyberspace*, Stanford University, March 1999.

29. Norbert Fuhr. "A Decision-Theoretic Approach to Database Selection in Networked IR". *ACM Transactions on Information Systems*, 17(3), July 1999.
30. Antony Galton. "Logic for Information Technology". John Wiley & Sons, 1990.
31. Hector Garcia-Molina, Yannis Papakonstantinou, Dallan Quass, Anand Rajaraman, Yehoshua Sagiv, Jeffrey Ullman, Vasilis Vassalos, and Jennifer Widom. "The TSIMMIS Approach to Mediation: Data Models and Languages". In *Proceedings of IPSJ*, Tokyo, Japan, October 1994.
32. Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. "Database System Implementation", chapter 11. Prentice Hall, 2000.
33. M. R. Genesareth, A. M. Keller, and O. Duschka. "Infomaster: An Information Integration System". In *Proceedings of 1997 ACM SIGMOD Conference*, May 1997.
34. L. Gravano and H. Garcia-Molina. "Generalizing GLOSS To Vector-Space Databases and Broker Hierarchies". In *Proc 21st VLDB Conf.*, Zurich, Switzerland, 1996.
35. Nicola Guarino. "Some Ontological Principles for Designing Upper Level Lexical Resources". In *Proceedings of first int. Conf. on Language Resources and Evaluation*, Granada, Spain, May 1998.
36. Nicola Guarino, Claudio Masolo, and Guido Vetere. "OntoSeek: Content-based Access to the Web". *IEEE Intelligent Systems*, pages 70–80, May, June 1999.
37. Alon Y. Halevy. "Answering Queries Using Views: A Survey". *VLDB Journal*, 10(4):270–294, 2001.
38. Heiko Helleg, Jurgen Krause, Thomas Mandl, Jutta Marx, Matthias Muller, Peter Mutschke, and Robert Strogon. "Treatment of Semantic Heterogeneity in Information Retrieval". Technical Report 23, Social Science Information Centre, May 2001. ([http://www.gesis.org/en/publications/reports/iz\\_working\\_papers/](http://www.gesis.org/en/publications/reports/iz_working_papers/)).
39. A. Howe and D. Dreilinger. "SavvySearch: A MetaSearch Engine that Learns Which Search Engines to Query". *AI Magazine*, 18(2), 1997.
40. International Organization For Standardization. "Documentation - Guidelines for the establishment and development of monolingual thesauri", 1986. Ref. No ISO 2788-1986.
41. Vipul Kashyap and Amit Sheth. "Semantic and Schematic Similarities between Database Objects: A Context-based Approach". *VLDB Journal*, 5(4), 1996.
42. Vipul Kashyap and Amit Sheth. "Semantic Heterogeneity in Global Information Systems: the Role of Metadata, Context and Ontologies". In *Cooperative Information Systems: Trends and Directions*. Academic Press, 1998.
43. C. Knoblock, Yigal Arens, and Chun-Nan Hsu. "Cooperating Agents for Information Retrieval". In *Proceedings of the Second International Conference on Cooperative Information Systems*, Toronto, Ontario, Canada, 1994.
44. Information Systems Laboratory. "The Semantic Index System (SIS)". Institute of Computer Science Foundation for Research and Technology Hellas. ([http://zeus.ics.forth.gr/forth/ics/isl/r-d-activities/semantic\\_index\\_system.html](http://zeus.ics.forth.gr/forth/ics/isl/r-d-activities/semantic_index_system.html)).
45. M. Lacher and G. Groh. "Facilitating the Exchange of Explicit Knowledge Through Ontology Mappings". In *Proceedings of the 14th Int. FLAIRS Conference*, 2001.
46. Veronique Lattes and Marie-Christine Rousset. "The use of CARIN language and algorithms for Information Integration: the PISCEL project". In *Proceedings of Second International and Interdisciplinary Workshop on Intelligent Information Integration*, Brighton Centre, Brighton, UK, August 1998.
47. Maurizio Lenzerini. Data integration: A theoretical perspective. In *Proc. ACM PODS 2002*, pages 233–246, Madison, Wisconsin, USA, June 2002.
48. Alon Y. Levy, Divesh Srivastava, and Thomas Kirk. "Data Model and Query Evaluation in Global Information Systems". *Journal of Intelligent Information Systems*, 5(2), 1995.
49. Sean Luke, Lee Spector, David Rager, and Jim Hendler. "Ontology-based Web Agents". In *Proceedings of First International Conference on Autonomous Agents*, 1997. (<http://www.cs.umd.edu/projects/plus/SHOE/>).
50. Zygmunt Mazur. "Models of a Distributed Information Retrieval System Based on Thesauri with Weights". *Information Processing and Management*, 30(1):61–77, 1994.
51. Deborah L. McGuinness. "Ontological Issues for Knowledge-Enhanced Search". In *Proceedings of FOIS'98*, Trento, Italy, June 1998. Amsterdam, IOS Press.
52. Carlo Meghini and Umberto Straccia. "A Relevance Terminological Logic for Information Retrieval". In *Proceedings of SIGIR'96*, Zurich, Switzerland, August 1996.
53. E. Mena, V. Kashyap, A. Sheth, and A. Illarramendi. "OBSERVER: An Approach for Query Processing in Global Information Systems based on Interoperation across Preexisting Ontologies.". In *Proceedings of the First IFCIS International Conference on Cooperative Information Systems (CoopIS'96)*, Brussels, Belgium, June 1996. IEEE Computer Society Press.
54. P. Mitra, G. Wiederhold, and J. Jannink. "Semi-automatic Integration of Knowledge sources". In *Proc. of the 2nd Int. Conf. On Information FUSION*, 1999.
55. Esko Nuutila. "Efficient Transitive Closure Computation in Large Digraphs". PhD thesis, Acta Polytechnica Scandinavica, Helsinki, 1995. (url = <http://www.cs.hut.fi/~enu/thesis.html>).
56. C. Paice. "A Thesaural Model of Information Retrieval". *Information Processing and Management*, 27(5):433–447, 1991.
57. Ruben Prieto-Diaz. "Implementing Faceted Classification for Software Reuse". *Communications of the ACM*, page 88, 1991.
58. S. R. Ranganathan. "The Colon Classification". In Susan Artandi, editor, *Vol IV of the Rutgers Series on Systems for the Intellectual Organization of Information*. New Brunswick, NJ: Graduate School of Library Science, Rutgers University, 1965.
59. Thomas Rolleke and Norbert Fuhr. "Retrieval of Complex Objects Using a Four-Valued Logic". In *Proceedings of SIGIR'96*, Zurich, Switzerland, August 1996.
60. I. Ryutaro, T. Hideaki, and H. Shinichi. "Rule Induction for Concept Hierarchy Allignment". In *Proceedings of the 2nd Workshop on Ontology Learning at the 17th Int. Conf. on AI (IJCAI)*, 2001.
61. Giovanni M. Sacco. "Dynamic Taxonomies: A Model for Large Information Bases". *IEEE Transactions on Knowledge and Data Engineering*, 12(3), May 2000.

62. G. Salton. “Introduction to Modern Information Retrieval”. McGraw-Hill, 1983.
63. E. Selberg and O. Etzioni. “Multi-Service Search and Comparison Using the MetaCrawler”. In *Proceedings of the 1995 World Wide Web Conference*, December 1995.
64. Marios Sintichakis and Panos Constantopoulos. “A Method for Monolingual Thesauri Merging”. In *Proceedings of 20th International Conference on Research and Development in Information Retrieval, ACM SIGIR '97*, Philadelphia, PA, USA, July 1997.
65. Nicolas Spyrtos. “The Partition Model: A Deductive Database Model”. *ACM Transactions on Database Systems*, 12(1):1–37, 1987.
66. V. S. Subrahmanian, S. Adah, A. Brink, R. Emery, A. Rajput, R. Ross, T. Rogers, and C. Ward. “HERMES: A Heterogeneous Reasoning and Mediator System”, 1996. ([www.cs.umd.edu/projects/hermes/overview/paper](http://www.cs.umd.edu/projects/hermes/overview/paper)).
67. Yannis Tzitzikas. “Democratic Data Fusion for Information Retrieval Mediators”. In *ACS/IEEE International Conference on Computer Systems and Applications*, Beirut, Lebanon, June 2001.
68. Yannis Tzitzikas, Anastasia Analyti, Nicolas Spyrtos, and Panos Constantopoulos. “An Algebraic Approach for Specifying Compound Terms in Faceted Taxonomies”. In *13th European-Japanese Conference on Information Modelling and Knowledge Bases*, Kitakyushu, Japan, June 2003.
69. Yannis Tzitzikas and Carlo Meghini. “Ostensive Automatic Schema Mapping for Taxonomy-based Peer-to-Peer Systems”. In *Seventh International Workshop on Cooperative Information Agents, CIA-2003*, number 2782 in Lecture Notes on Artificial Intelligence, pages 78–92, Helsinki, Finland, August 2003. (Best Paper Award).
70. Yannis Tzitzikas and Carlo Meghini. “Query Evaluation in Peer-to-Peer Networks of Taxonomy-based Sources”. In *Proceedings of 19th Int. Conf. on Cooperative Information Systems, CoopIS'2003*, Catania, Sicily, Italy, November 2003.
71. Yannis Tzitzikas, Carlo Meghini, and Nicolas Spyrtos. “Taxonomy-based Conceptual Modeling for Peer-to-Peer Networks”. In *Proceedings of 22th Int. Conf. on Conceptual Modeling, ER'2003*, Chicago, Illinois, October 2003.
72. Yannis Tzitzikas, Nicolas Spyrtos, and Panos Constantopoulos. “Mediators over Ontology-based Information Sources”. In *Proceedings of the 2nd International Conference on Web Information Systems Engineering, WISE 2001*, pages 31–40, Kyoto, Japan, December 2001.
73. Yannis Tzitzikas, Nicolas Spyrtos, and Panos Constantopoulos. “Query Evaluation for Mediators over Web Catalogs”. *International Journal on Information Theories and Applications*, 9(2), 2002.
74. Jeffrey D. Ullman. “Information integration using logical views”. In *In Proc. of the 6th Int. Conf. on Database Theory (ICDT-97)*, 1997.
75. Frank van Harmelen and Dieter Fensel. “Practical Knowledge Representation for the Web”. In *Workshop on Intelligent Information Integration, IJCAI'99*, 1999.
76. E. Vorhees, N. Gupta, and B. Johnson-Laird. “The Collection Fusion Problem”. In *Proceedings of the Third Text Retrieval Conference (TREC-3)*, Gaithersburg, MD, 1995.
77. Ellen Vorhees. “Multiple Search Engines in Database Merging”, 1997. DL 97.
78. G. Wiederhold. “Mediators in the Architecture of Future Information Systems”. *IEEE Computer*, 25:38–49, 1992.
79. R. Yerneni, Chen Li, H.Garcia-Molina, and J.Ullman. “Computing capabilities of mediators”. In *Proceedings of ACM SIGMOD'99*, Philadelphia, 1999.

## Appendix: Proofs

**Prop. 1.** If  $I$  is an interpretation of  $T$  then  $I^-$  is the unique minimal model of  $T$  which is greater than or equal to  $I$ .

**Proof:**

( $I^-$  is a model of  $T$ )

$t \preceq t' \Rightarrow \text{tail}(t) \subseteq \text{tail}(t') \Rightarrow \bigcup\{I(s) \mid s \in \text{tail}(t)\} \subseteq \bigcup\{I(s) \mid s \in \text{tail}(t')\} \Rightarrow I^-(t) \subseteq I^-(t')$ . Thus  $I^-$  is a model of  $T$ .

( $I^-$  is the unique minimal model of  $T$  which is greater than  $I$ )

Let  $I'$  be a model of  $T$  which is larger than  $I$ . Below we prove that  $I^- \subseteq I'$ . By the definition of  $I^-(t)$ , if  $o \in I^-(t)$  then either  $o \in I(t)$  or  $o \in I(s)$  for a term  $s$  such that  $s \preceq t$ . However, if  $o \in I(t)$  then  $o \in I'(t)$  too because  $I'$  is larger than  $I$ , and if  $o \in I(s)$  for a term  $s$  such that  $s \preceq t$  then  $o \in I'(t)$  too because  $I'$  is a model of  $T$ . We conclude that for every  $o \in I^-(t)$  it holds  $o \in I'(t)$ . Thus  $I^-$  is the unique minimal model  $T$  which is larger than  $I$ .

◇

**Prop. 2.** If  $I$  is an interpretation of  $T$  then  $I^+$  is a model of  $T$  and  $I \subseteq I^- \subseteq I^+$ .

**Proof:**

( $I^+$  is a model of  $T$ )

$t \preceq t' \Rightarrow \{u \mid u \in \text{head}(t)\} \supseteq \{u \mid u \in \text{head}(t')\} \Rightarrow \{u \mid u \in \text{head}(t) \text{ and } u \not\prec t\} \supseteq \{u \mid u \in \text{head}(t') \text{ and } u \not\prec t'\} \Rightarrow \bigcap\{I(u) \mid u \in \text{head}(t) \text{ and } u \not\prec t\} \subseteq \bigcap\{I(u) \mid u \in \text{head}(t') \text{ and } u \not\prec t'\} \Rightarrow I^+(t) \subseteq I^+(t')$ .

( $I^- \subseteq I^+$ )

Clearly, if  $t \in T$ ,  $u \in \text{head}(t)$  and  $u \not\prec t$  then in every model  $I$  of  $T$  we have  $I(t) \subseteq I(u)$ . Thus this also holds in the model  $I^-$ , i.e.  $I^-(t) \subseteq I^-(u)$ . From this we conclude that for every  $t \in T$ :

$I^-(t) \subseteq \bigcap\{I^-(u) \mid u \in \text{head}(t) \text{ and } u \not\prec t\} = I^+(t)$ . Thus  $I^- \subseteq I^+$ .

◇

**Prop. 3.** The answer models of the mediator are ordered as follows:

$$(a) I_{i-}^- \subseteq I_{i+}^-$$

$$(b) I_{u-}^- \subseteq I_{u+}^-$$

$$(c) I_{i-}^+ \subseteq I_{i+}^+$$

$$(d) I_{u-}^+ \subseteq I_{u+}^+$$

$$(e) I_{i-}^- \subseteq I_{i-}^+$$

$$(f) I_{i+}^- \subseteq I_{i+}^+$$

$$(g) I_{u-}^- \subseteq I_{u-}^+$$

$$(h) I_{u+}^- \sqsubseteq I_{u+}^+$$

**Proof:**

The proofs of the propositions (a)-(d) come easily from the fact that in every model  $I_i$  of a source  $S_i$  it holds:  $I_i^- \sqsubseteq I_i^+$ .

The proofs of the propositions (e)-(h) come easily from the fact that in every model  $I$  of the mediator it holds:  $I^- \sqsubseteq I^+$ .

◇

**Prop. 4.** If all sources are compatible with the mediator then:

- (1)  $I_l^- \sqsubseteq I_{u-}^-$
- (2)  $I_{l+}^- \sqsubseteq I_{u+}^-$
- (3)  $I_{l-}^+ \sqsubseteq I_{u-}^+$
- (4)  $I_{l+}^+ \sqsubseteq I_{u+}^+$

**Proof:**

Let  $t$  be a term of  $T$ . Clearly, for every  $s \in \text{tail}_i(t)$  and  $u \in \text{head}_i(t)$  it holds  $sa_iu$  (because  $sa_it$  and  $ta_iu$ ). Since the source  $S_i$  is compatible we know that  $sa_iu \Rightarrow s \preceq_i u$ . This implies that in every model  $I_i$  of  $T_i$  it holds:

$$\bigcup \{I_i(s) \mid sa_it\} \subseteq \bigcap \{I_i(u) \mid ta_iu\} \Leftrightarrow I_i(t_i^s) \subseteq I_i(t_u^t)$$

From  $I_i(t_i^s) \subseteq I_i(t_u^t)$  we infer that  $I_i^-(t_i^s) \subseteq I_i^-(t_u^t)$  and  $I_i^+(t_i^s) \subseteq I_i^+(t_u^t)$ . From this we obtain propositions (1)-(4).

For example, the proof of proposition (1) i.e.  $I_l^- \sqsubseteq I_{u-}^-$ , and proposition (3) i.e.  $I_{l-}^+ \sqsubseteq I_{u-}^+$ , comes as follows: Since  $\forall t \in T$  and  $\forall i = 1..k$  it holds  $I_i^-(t_i^s) \subseteq I_i^-(t_u^t)$ , we conclude that:

$$\bigcup_{i=1..k} I_i^-(t_i^s) \subseteq \bigcup_{i=1..k} I_i^-(t_u^t) \Leftrightarrow I_l^-(t) \subseteq I_{u-}^-(t) \Rightarrow \begin{array}{l} I_l^- \sqsubseteq I_{u-}^- \quad (I_1 \sqsubseteq I_3) \\ I_{l-}^+ \sqsubseteq I_{u-}^+ \quad (I_5 \sqsubseteq I_7) \end{array}$$

◇

**Prop. 5.** If  $q = t \in T$ , then  $I_l^-(t)$  and  $I_u^-(t)$  can be evaluated as follows:

$$\begin{aligned} I_l^-(t) &= \bigcup_{i=1..k} I_i(q_i^i(t)) \quad \text{where} \quad q_i^i(t) = \bigvee \{s_i^i \mid s \preceq t\} \\ I_u^-(t) &= \bigcup_{i=1..k} I_i(q_u^i(t)) \quad \text{where} \quad q_u^i(t) = \bigvee \{s_u^i \mid s \preceq t\} \end{aligned}$$

**Proof:**

$$\begin{aligned} I_l^-(t) &= \bigcup \{I_l(s) \mid s \in \text{tail}(t)\} = \bigcup \{I_l(s) \mid s \preceq t\} \\ &= \bigcup \{\bigcup_{i=1..k} I_i(s_i^i) \mid s \preceq t\} = \bigcup_{i=1..k} \{I_i(s_i^i) \mid s \preceq t\} \\ &= \bigcup_{i=1..k} I_i(\bigvee \{s_i^i \mid s \preceq t\}) = \bigcup_{i=1..k} I_i(q_i^i(t)) \end{aligned}$$

Analogously, we prove that  $I_u^-(t) = \bigcup_{i=1..k} I_i(q_u^i(t))$ . ◇

**Prop. 6.**  $D^-(o) = \bigcup \{\text{head}(t) \mid o \in I(t)\}$

**Proof:**

$$\begin{aligned} D^-(o) &= \{t \in T \mid o \in I^-(t)\} = \{t \in T \mid o \in \bigcup \{I(s) \mid s \preceq t\}\} \\ &= \{t \in T \mid o \in I(s) \text{ and } s \preceq t\} = \bigcup \{\text{head}(s) \mid o \in I(s)\} \end{aligned}$$

◇

**Prop. 7.**  $D^+(o) = \{t \mid \text{head}(t) \setminus \{t' \mid t' \sim t\} \subseteq D^-(o)\}$

**Proof:**

$$\begin{aligned} t \in D^+(o) &\Leftrightarrow o \in I^+(t) \Leftrightarrow o \in \bigcap \{I^-(u) \mid u \in \text{head}(t) \text{ and } u \not\sim t\} \\ &\Leftrightarrow o \in I^-(u) \quad \forall u \in \text{head}(t) \text{ s.t. } u \not\sim t \Leftrightarrow u \in D^-(o) \quad \forall u \in \text{head}(t) \text{ s.t. } u \not\sim t \\ &\Leftrightarrow \text{head}(t) \setminus \{t' \mid t' \sim t\} \subseteq D^-(o). \end{aligned}$$

◇

**Prop. 8.**

$$D_l(o) = \bigcup_{i=1..k} D_l^i(o) \quad \text{where}$$

$$D_l^i(o) = \{t \in T \mid t_i \in D_i(o) \text{ and } t_i a_i t\}$$

$$D_u(o) = \bigcup_{i=1..k} D_u^i(o) \quad \text{where}$$

$$\begin{aligned} D_u^i(o) &= \{t \in T \mid (\text{head}_i(t) \neq \emptyset \text{ and } \text{head}_i(t) \subseteq D_i(o)) \text{ or} \\ &\quad (\text{head}_i(t) = \emptyset \text{ and } t_i \in D_i(o) \text{ and } t_i a_i t)\} \end{aligned}$$

**Proof.**

Consider a mediator over a single source  $S_i$  and let  $o$  be an object stored at that source. Let  $t_i \in D_i(o)$  where  $I_i$  is the answer model of the source  $S_i$  that it is used by the mediator. If  $\exists t \in T$  such that  $t_i a_i t$  then certainly  $o \in I_l(t)$  (since  $I_l(t) = \bigcup \{I(t_i) \mid t_i a_i t\}$ ), thus  $t \in D_l(o)$ . Hence  $D_l(o) = \{t \in T \mid t_i \in D_i(o) \text{ and } t_i a_i t\}$ . However, since there are many sources, we denote the right part of the above formula by  $D_l^i(o)$ , and since an object may belong to more than one source, we reach to the following:  $D_l(o) = \bigcup_{i=1..k} D_l^i(o)$ .

Consider again a mediator over a single source  $S_i$  and let  $o$  be an object stored at that source. If  $\exists t \in T$  such that  $\text{head}_i(t) \neq \emptyset$  and  $\text{head}_i(t) \subseteq D_i(o)$  then certainly  $o \in I_u(t)$  (since  $I_u(t) = \bigcap \{I(t_i) \mid t a_i t_i\}$ ), thus certainly  $t \in D_u(o)$ . If  $\exists t \in T$  such that  $\text{head}_i(t) = \emptyset$  and there is a  $t_i \in D_i(o)$  and  $t_i a_i t$  then certainly  $o \in I_u(t)$  (since in this case  $I_u(t) = \bigcup \{I(t_i) \mid t_i a_i t\}$ ). Hence  $D_u(o) = \{t \in T \mid (\text{head}_i(t) \neq \emptyset \text{ and } \text{head}_i(t) \subseteq D_i(o)) \text{ or } (\text{head}_i(t) = \emptyset \text{ and } t_i \in D_i(o) \text{ and } t_i a_i t)\}$ . However, since there are many sources, we denote the right part of the above formula by  $D_u^i(o)$ , and since an object may belong to more than one source, we reach to the following:  $D_u(o) = \bigcup_{i=1..k} D_u^i(o)$ .

◇