

Start making sense: The Chatty Web approach for global semantic agreements^{*}

Karl Aberer, Philippe Cudré-Mauroux, Manfred Hauswirth
École Polytechnique Fédérale de Lausanne (EPFL)
CH-1015 Lausanne, Switzerland
{Karl.Aberer, Philippe.Cudre-Mauroux, Manfred.Hauswirth}@epfl.ch

Abstract

This paper describes a novel approach for obtaining semantic interoperability among data sources in a bottom-up, semi-automatic manner without relying on pre-existing, global semantic models. We assume that large amounts of data exist that have been organized and annotated according to local schemas. Seeing semantics as a form of agreement, our approach enables the participating data sources to incrementally develop global agreement in an evolutionary and completely decentralized process that solely relies on pair-wise, local interactions: Participants provide translations between schemas they are interested in and can learn about other translations by routing queries (gossiping). In previous work we relied on the realistic assumption that such translations would be provided manually only. In contrast to that, we assume in this paper that only some translations exist and generate random translations for reaching overall semantic agreement automatically. To support the participants in assessing the semantic quality of the achieved agreements we develop a formal framework that takes into account both syntactic and semantic criteria. The assessment process is incremental and the quality ratings are adjusted along with the operation of the system. Ultimately, this process results in global agreement, i.e., the semantics that all participants understand. We discuss strategies to efficiently find translations and provide results from our experiments to justify our claims. We specifically focus on semantic analyses and provide pointers to the possible quality that is achievable through semantic analysis only. Our approach applies to any system which provides a communication infrastructure (existing websites or databases, decentralized systems, P2P systems) and offers the opportunity to study semantic interoperability as a global phenomenon in a network of information sharing parties.

Keywords: Semantic integration, semantic agreements, self-organization

1 Introduction

The recent success of peer-to-peer (P2P) systems and the initiatives to create the Semantic Web have emphasized again a key problem in information systems: the lack of semantic interoperability. Semantic interoperability is a crucial element for making distributed information systems usable. It is prerequisite for structured, distributed search and data exchange and provides the foundations for higher level (web) services and processing.

For example, the technologies that are currently in place for P2P file sharing systems either impose a simple semantic structure a-priori (e.g., Napster, Kazaa) and leave the burden of semantic annotation to the user, or do not address the issue of semantics at all (e.g.,

^{*}The work presented in this paper was supported (in part) by the National Competence Center in Research on Mobile Information and Communication Systems (NCCR-MICS), a center supported by the Swiss National Science Foundation under grant number 5005-67322.

the current web, Gnutella, Freenet) but simply support a semantically unstructured data representation and leave the burden of “making sense” to the skills of the user, e.g., by providing pseudo-structured file names such as *Enterprise-2x03-Mine-Field* that encapsulate very simple semantics.

Also, classical attempts to make information resources semantically interoperable, in particular in the domain of database integration, do not scale well to global information systems, like P2P systems. Despite a large number of approaches and concepts, such as federated databases, the mediator concept [26], or ontology-based information integration approaches [12, 21], practically engineered solutions are still frequently hard-coded and require substantial support from human experts. A typical example of such systems are domain-specific portals such as CiteSeer (www.researchindex.com, publication data), SRS (srs.ebi.ac.uk, biology) or streetprices.com (e-commerce). They integrate data sources on the Internet and store them in a central warehouse. The data is converted to a common schema which usually is of simple to medium complexity. This approach adopts a simple form of wrapper-mediator architecture and typically requires substantial development efforts for the automatic or semi-automatic generation of mappings from the data sources into the global schema.

In the context of the Semantic Web, a major effort is devoted to the provision of machine processable semantics expressed in meta-models such as RDF, OIL [7], OWL [5] or DAML+OIL [11], and based on shared ontologies. Still, these approaches rely on commonly agreed upon ontologies, which existing information sources can be related to by proper annotation. This is an extremely important development, but its success will heavily rely on the wide adoption of common ontologies or schemas.

The advent of P2P systems, however, introduces a different view on the problem of semantic interoperability by taking a social perspective which relies on self-organization heavily. We argue that we can see the emerging P2P paradigm as an opportunity to improve semantic interoperability rather than a threat, in particular in revealing new possibilities on how semantic agreements can be achieved. This motivated us to look at the problem from a different perspective and has inspired the approach presented in this paper.

In the following, we abstract from the underlying infra-structure such as federated databases, web sites or P2P systems and regard these systems as graphs of interconnected data sources. For simplicity, but without constraining the general applicability of the presented concepts, we denote these data sources as *peers*. Each peer offers data which are organized according to some schema expressed in a data model, e.g., relational, XML, or RDF. Among the peers, communication is supported via suitable protocols and architectures, for example, HTTP or JXTA.

The first thing to observe is that semantic interoperability is always based on some form of agreement. Ontology-oriented approaches in the Semantic Web represent this agreement essentially *explicitly* through a shared ontology. In our approach, no explicit representation of a globally shared agreement will be required, but agreements are *implicit* and result from the way our (social) mechanism works.

We impose a modest requirement on establishing agreements by assuming the existence of local agreements provided as mappings between different schemas, i.e., agreements established in a P2P manner. These agreements will have to be established in a manual or semiautomatic way since in the near future we do not expect to be able to fully automate the process of establishing semantic mappings even locally. However, a rich set of tools is getting available to support this [24]. Establishing local agreements is a less challenging task than establishing global agreements by means of globally agreed schemas or shared ontologies. Once such agreements exist, we establish on-demand relationships among schemas of different information systems that are sufficient to satisfy information processing needs such as distributed search.

We briefly highlight two of the application scenarios that convinced us (besides the obvious applicability for information exchange on the web) that enabling semantic interoperability in a bottom-up way driven by the participants is valid and applicable: introduction of

meta-data support in P2P applications and support for federating existing, loosely-coupled databases.

Imposing a global schema for describing data in P2P systems is almost impossible, due to the decentralization properties of such systems. It would not work unless all users conscientiously follow the global schema. Here our approach would fit well: We let users introduce their own schemas which best meet their requirements. By exchanging translations between these schemas, the peers can incrementally come up with an implicit “consensus schema” which gradually improves the global search capabilities of the P2P system. This approach is orthogonal to the existing P2P systems and could be introduced basically into all of them.

The situation is somewhat similar for federating existing loosely-coupled databases. Such large collections of data exist for example for biological or genomic databases. Each database has a predefined schema and possibly some translations may already be defined between the schemas, for example data import/export facilities. However, global search, i.e., propagation of queries among the set of databases, is usually not provided and if this feature exists it is usually done in an ad-hoc, non-systematic way, i.e., not reusable and not automated. The more complex these database schemas get, the less likely it is that the schemas partially overlap and the harder it gets to increasingly generate translations automatically.

In our approach, we build on the principle of gossiping that has been successfully applied for creating useful global behaviors in P2P systems. In any P2P system, search requests are routed in a network of interconnected information systems. We extend the operation of these systems as follows: When different schemas are involved, local mappings are used to further distribute a search request into other semantic domains.

For simplicity but without constraining general applicability, we will limit the following discussions to the processing of search requests. The quality of search results in such a gossiping-based approach depends clearly on the quality of the local mappings in the mapping graph. *Our fundamental assumption is that these mappings may be incorrect.* Thus our agreement construction mechanisms try to determine which mappings can be trusted and which not and take this into account to guide the search process.

A main contribution of the paper is to identify the different methods that can be applied to establish global forms of agreement starting from a graph of local mappings among schemas. We elaborate the details of each of these methods for a simple data model, that is yet expressive enough to cover many practical cases. This model is similar to other data models currently considered for semantic annotation in P2P architectures [15]. Three methods will be introduced in particular:

1. A syntactic analysis of search queries after mappings have been applied in order to determine the potential information-loss incurred through the transformation.
2. A semantic analysis of composite mappings along cycles in the mapping graph, in order to determine the level of agreement that peers achieve throughout the cycle.
3. A semantic analysis of search results obtained through composite mappings based on the preservation of data dependencies.

The information obtained by applying these different analyses is then used to direct searches in a network of semantically heterogeneous information sources (e.g, on top of a P2P network). We will provide results from first experiments that have been performed for this setting.

We believe that this radically new approach to semantic interoperability shifts the attention from problems that are inherently difficult to solve in an automated manner at the global level (“How do humans interpret information models in terms of real world concepts?”), to a problem that leaves vast opportunities for automated processing and for increasing the value of existing information sources, namely the processing of existing local semantic relationships in order to raise the level of their use from local to global semantic interoperability. The remaining problem of establishing semantic interoperability at a local

level seems to be much easier to tackle once an approach such as ours is in place.

2 Overview

Before delving into the technical details, this section provides an informal overview of our approach and of the paper.

We assume that there exists a communication facility among the participants that enables sending and receiving of information, i.e., queries, data, and schema information. This assumption does not constrain the approach, but emphasizes that it is independent of the system it is applied to. The underlying system could be a P2P system, a federated database system, the web, or any other system of information sources communicating via some communication protocol. We denote the participants as peers abstracting from the concrete underlying system.

In the system, groups of peers may have agreed on common semantics, i.e., a common schema. We denote these groups as *semantic neighborhoods*. The size of a neighborhood may range from a single individual peer up to any number. If two peers located in two disjoint neighborhoods meet, they can exchange their schemas and provide mappings between them (how peers meet and how they exchange this information depends on the underlying system but does not concern our approach). We assume that skilled experts supported by appropriate mapping tools provide the mappings. The direction of the mapping and the node providing a mapping are not necessarily correlated. For instance, nodes *A* and *B* might both provide a mapping from *schema(A)* to *schema(B)*, and they may exchange this mapping upon discretion. During the life-time of the system, each peer has the possibility to learn about existing mappings and add new ones. This means that a directed graph of mappings as shown in Figure 1 will be built between the neighborhoods along with the normal operation of the system (e.g., query processing and forwarding in a P2P system).

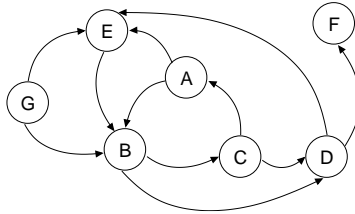


Figure 1: Mapping graph among semantic neighborhoods

This mapping graph has two interesting properties: (1) based on the already existing mappings and the ability to learn about existing mappings, new mappings can be added automatically by means of transitivity, for example, $D \rightarrow E \rightarrow B \Rightarrow D \rightarrow B$ and (2) the graph has cycles. (1) means that we can propagate queries towards nodes for which no direct translation link exists. This is what we call *semantic gossiping*. (2) gives us the possibility to assess the degree of *semantic agreement* along a cycle, i.e., to measure the quality of the translations and the degree of semantic agreement in a community.

In such a system, we expect peers to perform several task: (1) upon receiving a query, a peer has to decide where to forward the query to, based on a set of criteria that are introduced below; (2) upon receiving results or feedback (cycle), it has to analyze the quality of the results at the schema and at the data level and adjust its criteria accordingly; and (3) update its view of the overall semantic agreement.

The criteria to assess the quality of translations—which in turn is a measure of the semantic agreement—can be categorized as *context-independent* and *context-dependent*. Context-independent criteria, discussed in Section 4, are syntactic in nature and relate only

to the processed query and to the required translation. We introduce the notion of *syntactic similarity* to analyze the extent to what a query is preserved after translation.

Context-dependent criteria, which are discussed in Section 5, relate to the degree of agreement that can be achieved among different peers upon specific translations. Such degrees of agreement may be computed using feedback mechanisms (cycles appearing in the translation graph and results returned by different peers). This means that a peer will locally obtain both returned queries and data through multiple cycles. In case a disagreement is detected (e.g., a wrong attribute mapping at the schema level a concept mismatch at the content level), the peer has to suspect that at least some of the mappings involved in the cycle were incorrect, including the mapping it has used itself to propagate the query. Even if an agreement is detected, it is not clear whether this is not accidentally the result of compensating mapping errors along a cycle. Thus, analyses are required that assess which are the most probable sources of errors along cycles, to what extent the own mapping can be trusted and therefore of how to use these mappings in future routing decisions. At a global level, we can view the problem as follows: The translations in between domains of semantic homogeneity (same schemas) form a directed graph. Within that directed graph we find cycles. Each cycle allows to return a query to its originator which in turn can make the analysis described above.

Each of these criteria is applied on an attribute-basis to the transformed queries and results in a *feature vector*. This vector encompasses the outcome of the criterion for each of the attributes concerned. The decision whether or not to forward a query using a translation link then is based on these feature vectors. The details of the query forwarding process are provided in Section 6.

Assuming all the peers implement this approach, we expect the network to converge to a state where a query is only forwarded to the peers most-likely understanding it and where the correct mappings are increasingly reinforced by adapting the per-hop forwarding behaviors of the peers. Implicitly, this is a state where a global agreement on the semantics of the different schemas has been reached. To demonstrate this, we present experimental results where semantic agreement is reached in a network of partially erroneous mappings in Section 8.

3 The Model

3.1 The Data Model

We assume that each peer p is maintaining its database DB_p according to a schema S_p . The peers are able to identify their schema, either by explicitly storing it or by keeping a pseudo unique schema identifier, obtained for example by hashing. The schema consists of a single relational table R , i.e., the data that a peer stores consists of a set of tuples t_1, \dots, t_n of the same type. The attributes have complex data types and NULL-values are possible.

We do not consider more sophisticated data models to avoid diluting the discussion of the main ideas through technicalities related to mastering complex data models. Moreover, many practical applications, in particular in scientific databases, use exactly the type of simplistic data model we have introduced, at least at the meta-data level.

We use a query language for querying and transforming databases. The query language consists of basic relational algebra operators since we do not care about the practical encoding, e.g., in SQL or XQuery. The relational operators that we require are:

- Selection $\sigma_{p(a)}(R)$, where $a = \langle A_1, \dots, A_k \rangle$ is a list of attribute names, and p is any predicate on the attributes a using standard atomic predicates on the respective datatypes, i.e., $p = p(A_1, \dots, A_k)$.
- Projection $\pi_a(R)$, where a is a list of attribute names A_1, \dots, A_k .

- Mapping $\mu_f(R)$, where f is a list of functions of the form $A_0 := F(A_1, \dots, A_k)$ and A_1, \dots, A_k are attribute names occurring in R . The function F is specific to the datatypes of the attributes A_1, \dots, A_k . A special case is renaming of an attribute: $A_0 := A_1$.

We assume that queries can be evaluated against any database irrespective of its schema. Predicates containing attributes not present in the evaluated schema are ignored.¹ Projection attributes which are not present in the current schema return a NULL-value. Mappings applied to non-existing attributes also return NULL-values.

3.2 The Network Model

Let us now consider a set of peers P . Each peer $p \in P$ has a basic communication mechanism that allows it to establish connection with other peers. Without loss of generality, we assume in the following that it is based on the Gnutella protocol [4]. Thus peers can send *ping* messages and receive *pong* messages in order to learn about the network structure. In extension to the Gnutella protocol, peers also send their schema identifier as part of the *pong* message.

Every peer maintains a neighborhood $N(p)$ selected from the peers that it identified through *pong* messages. The peers in this neighborhood are distinguished into those that share the same schema, $N_e(p)$, and those that have a different schema, $N_d(p)$ as shown in Figure 2.

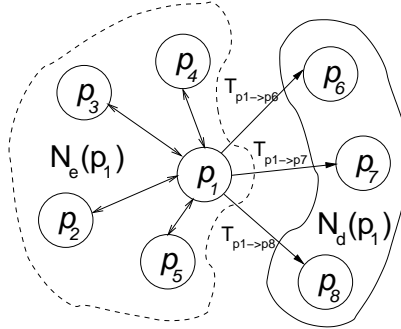


Figure 2: The Network Model

A peer p includes another peer p' with a different schema into its neighborhood if it knows a translation for queries against its own schema to queries against the foreign schema. The query translation operator $T_{p \rightarrow p'}$ is given by a query q_T that takes data structured according to schema $S_{p'}$ and transforms it into data structured according to schema S_p .

Thus $T_{p \rightarrow p'}$ has the property

$$T_{p \rightarrow p'}(q_p)(DB_{p'}) = q_p(q_T(DB_{p'}))$$

We assume that transformations only use a mapping operator followed by a projection on the attributes that are preserved. Thus q_T will always be of the form

$$q_T(DB_{p'}) = \pi_a(\mu_f(DB_{p'}))$$

Furthermore, we assume that the transformation query is normalized as follows: If an attribute A is preserved, it also occurs in the mapping operator as an identity mapping, i.e.,

¹We do not use the same conventions as XPath/XQuery here, but we will make use of additional mechanisms for dropping queries.

$A := A \in f$. This simplifies our subsequent analysis. Note that multiple transformations may be applied to a single query iteratively:

$$T_{n-1 \rightarrow n}(\dots T_{1 \rightarrow 2}(q) \dots) = T_{1 \rightarrow 2, \dots, n-1 \rightarrow n}(q)$$

Such query translations may be implemented easily using various mechanisms, for example XQuery, as done in our case study in Section 8.

Queries can be issued to any peer through a query message. A query message contains a query identifier id , the (potentially transformed) query q , the query message originator p , and the translation trace TT to keep track of the translations already performed. In the subsequent sections we will extend the contents of the query message in order to implement a more intelligent control of query forwarding. The basic query message format is

$$query(id, q, p, TT)$$

The translation trace TT is a list of pairs $\{(p_{from}, S_{p_{from}}), (p_{to}, S_{p_{to}})\}$ keeping track of the peers having sent the request through a translation link (p_{from}) and of the peers having received it after the translation link (p_{to}), along with their respective schema identifiers ($S_{p_{from}}$ and $S_{p_{to}}$). We will call p_{from} the sender, and p_{to} the receiver. For any translation link, we have to record both the sender and the recipient, as after a translation a query might be forwarded without transformation to peers sharing the same schema.

4 Syntactic Similarity

As context-independent criterion to measure the degree of similarity between two queries (in our context, between an original query and a transformed query), we introduce the notion of *syntactic similarity*, which is related to the number of attributes preserved during translation. Note that a high syntactic similarity in terms of number of attributes lost during translation will not ensure that forwarding the query is useful, but conversely a low syntactic similarity implies that it might not be useful to forward the query.

Let us suppose we have a query q , which always has the generic form of a selection-projection-mapping query

$$q = \pi_{ap}(\sigma_{p(as)}(\mu_{fa}(DB)))$$

where as is a list of attributes used in the selection predicates, ap is a list of attributes used in the projection, and fa is a list of functions applied. Again, without loss of generality, we assume that the query is normalized such that all attributes required in as and ap are computed by one of the functions in fa to simplify the subsequent analysis.

Therefore the transformed query will be of the form (this is also true for multiple transformations after normalization)

$$T(q)(DB') = \pi_{ap}(\sigma_{p(as)}(\mu_{fa}(\pi_a(\mu_f(DB')))))$$

It might occur that attributes used in q are no longer available after the transformation T has been applied to q . This can only happen when an attribute needed for the derivation of a new attribute by means of one of the functions in fa and required in ap or as is missing, i.e., not occurring in a .

We now analyze which attributes are exactly needed in order to properly evaluate the query q . We define

$$att_{\sigma}(q) = \{[A_0 : \{A_1, \dots, A_k\}] \mid A_0 \in as, A_0 := F(A_1, \dots, A_k) \in fa\}$$

and similarly

$$att_{\pi}(q) = \{[A_0 : \{A_1, \dots, A_k\}] \mid A_0 \in ap, A_0 := F(A_1, \dots, A_k) \in fa\}$$

Given a transformation T we can define the source of an attribute $source_T(A)$:

If

$$\exists F \in fa \text{ such that } A := F(A_1, \dots, A_k)$$

then

$$source_T(A) = \{A_1, \dots, A_k\}$$

else

$$source_T(A) = \perp.$$

Informally, $source_T(A)$ tells whether and how an attribute is preserved in a transformation T . Then we can define the operation $\omega_T(att_\sigma(q))$ as follows:

If

$$\forall [A_0 : \{A_1, \dots, A_k\}] \in att_\sigma(q) \quad \forall A \in \{A_1, \dots, A_k\} source_T(A) \neq \perp$$

then

$$[A_0 : \bigcup_{A \in \{A_1, \dots, A_k\}} source_T(A)] \in \omega_T(att_\sigma(q))$$

else

$$[A_0 : \perp] \in \omega_T(att_\sigma(q)).$$

This definition extends naturally to multiple transformation. In order to define

$$\omega_{T_n}(\dots(\omega_{T_1}(att_\sigma(q)))\dots)$$

we simply apply the above definition for ω_{T_n} to $\omega_{T_{n-1}}(\dots(\omega_{T_1}(att_\sigma(q)))\dots)$ instead of $att_\sigma(q)$. All definitions are analogous for $\omega_T(att_\pi(q))$.

$\omega_{T_n}(\dots(\omega_{T_1}(att_\sigma(q)))\dots)$ allows to determine which of the required attributes for evaluating queries are at disposal after applying the transformations T_1, \dots, T_n . The definitions are given such that they can be evaluated locally, i.e., for each transformation step in an iterative manner. Using this information we can now define the syntactic similarity between the transformed query and the original query.

The decision on the importance of attributes is query dependent. We have two issues to consider:

1. Not all attributes in as are preserved. Therefore some of the atomic predicates in $p(as)$ will not be correctly evaluated (the atomic predicates will simply be dropped in this case). Depending on the selectivity of the predicate this might be harmful to different degrees. We capture this by calculating a value $fv_{A_i}^\sigma$ for every attribute $A_i \in as \cup ap$ as follows:

If

$$A_i \in as, [A_i : \perp] \in \omega_{T_n}(\dots(\omega_{T_1}(att_\sigma(q)))\dots)$$

then

$$fv_{A_i}^\sigma(T_{1 \rightarrow \dots \rightarrow n}(q)) = 0$$

else

$$fv_{A_i}^\sigma(T_{1 \rightarrow \dots \rightarrow n}(q)) = sel_{A_i}$$

where sel_{A_i} is the selectivity (ranging over the interval $[0, 1]$, with high values indicating highly-selective attributes, i.e., attributes whose predicates select a small proportion of the database) of an attribute A_i . Given the values $fv_{A_i}^\sigma$ for $A_i \in as \cup ap$ we can introduce a feature vector $\overrightarrow{FV}_\sigma$ for the transformed query $T_n(\dots T_1(q) \dots)$ characterizing the syntactic similarity with respect to the selection operator:

$$\overrightarrow{FV}_\sigma(T_{1 \rightarrow \dots \rightarrow n}(q)) = (fv_{A_1}^\sigma, \dots, fv_{A_k}^\sigma)$$

We derive the syntactic similarity between the original query and the transformed query for the selection from this feature vector and from a user-defined weight vector $\overrightarrow{W} = (w_{A_1}, \dots, w_{A_k})$ with $A_i \in as \cup ap$ pondering the importance of the attributes:

$$S_\sigma(q, T_{1 \rightarrow \dots \rightarrow n}(q)) = \frac{\overrightarrow{W} \cdot \overrightarrow{FV}_\sigma}{|\overrightarrow{W}| |\overrightarrow{FV}_\sigma|}$$

where

$$\overrightarrow{W} \cdot \overrightarrow{FV}_\sigma = w_1 fv_{A_1}^\sigma + w_2 fv_{A_2}^\sigma + \dots + w_k fv_{A_k}^\sigma$$

and where

$$|\overrightarrow{X}| = \|\overrightarrow{X}\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_k^2}.$$

This value is normalized on the interval $[0, 1]$. Originally, the similarity will be one, and it will decrease proportionally to the relative weight and selectivity of every attribute lost in the selection operator, until it reaches 0 when all attributes are lost.

2. Not all attributes in ap are found in a or af . Therefore, some of the results may be incomplete or even erroneous (due to the loss of key attributes, for example). Following the method used above for the selection, we propose to measure this for every attribute:

If

$$A_i \in as, [A_i : \perp] \in \omega_{T_n}(\dots (\omega_{T_1}(att_\pi(q)) \dots))$$

then

$$fv_{A_i}^\pi(T_{1 \rightarrow \dots \rightarrow n}(q)) = 0$$

else

$$fv_{A_i}^\pi(T_{1 \rightarrow \dots \rightarrow n}(q)) = 1.$$

The feature vector and the syntactic similarity for the projection operator then are

$$\overrightarrow{FV}_\pi(T_{1 \rightarrow \dots \rightarrow n}(q)) = (fv_{A_1}^\pi, \dots, fv_{A_k}^\pi)$$

and

$$S_\pi(q, T_{1 \rightarrow \dots \rightarrow n}(q)) = \frac{\overrightarrow{W} \cdot \overrightarrow{FV}_\pi}{|\overrightarrow{W}| |\overrightarrow{FV}_\pi|}.$$

Again, this similarity decreases with the number of translations applied to the query, until it reaches 0 when all the projection attributes are lost.

5 Semantic Similarity

The context-independent measure of syntactic similarity is based on the assumption that the query translations are semantically correct, which in general might not be the case. A better way to view semantics is to consider it as an agreement among peers. If two peers agree on the meaning of their schemas, then they will generate compatible translations. From that basic observation, we will now derive context-dependent measures of semantic similarity. These measures will allow us to assess the quality of attributes that are preserved in the translation.

To that end, we introduce two mechanisms for deriving the quality of a translation. One mechanism will be based on analyzing the fidelity of translations at the schema level, the other one will be based on analyzing the quality of the correspondences in the query results obtained at the data level.

5.1 Cycle Analysis

For the first mechanism, we exploit the protocol property that detects cycles as soon as a query reenters a semantic domain it has already traversed (see Section 6 for the exact algorithm). Such a cycle starts with a peer p_1 transmitting a query q_1 to a second peer p_2 through a translation link $T_{1 \rightarrow 2}$ (see Figure 3).

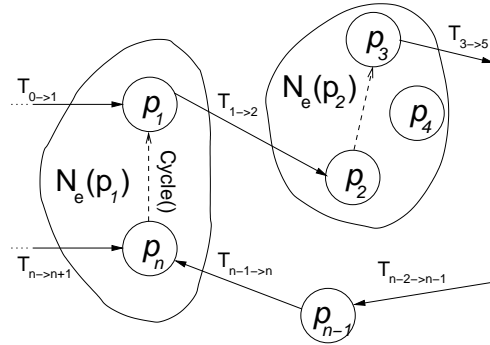


Figure 3: The Feedback Mechanism

In the example, after a few hops, the query is finally sent to a peer p_n which, sharing the same schema as p_1 , detects a cycle and informs p_1 . The returning query q_n is of the form

$$q_n = T_{1 \rightarrow 2, 3 \rightarrow 5, \dots, n-1 \rightarrow n}(q_1)$$

p_1 may now analyze what happened to the attributes $A_1 \dots A_k$ originally present in q_1 . We could attempt to check whether the composed mapping is identity, but the approach we propose here appears more practical. We differentiate three cases:

- Case 1: $source_{T_{1 \rightarrow \dots \rightarrow n}}(A_i) = \{A_i\}$, this means that A_i has been maintained throughout the cycle. It usually indicates that all the peers along the cycle agree on the meaning of the attribute. Such an observation increases the confidence in the correctness of the mapping.
- Case 2: $source_{T_{1 \rightarrow \dots \rightarrow n}}(A_i) = \perp$, this means that someone along the cycle had no representation for A_i . A_i is not part of the common semantics. This leaves the confidence in the mapping unchanged.
- Case 3: Otherwise, if none of the two previous cases occurs, e.g., $source_{T_{1 \rightarrow \dots \rightarrow n}}(A_i) = \{A_j\}, j \neq i$, this indicates some semantic confusion along the cycle. Subcases can occur depending on what happens to A_j . This lowers the confidence in the mapping.

We then derive heuristics for p_1 to assess the correctness of the translation $T_{1 \rightarrow 2}$ it has used based on the different cycle messages it received. Let us consider a translation cycle c_i composed of $\|c_i\|$ translation links. On an attribute basis, c_i may result in *positive* feedback (case 1 above), *neutral* feedback (case 2, not used for the rest of this analysis but taken into account by the syntactic similarity), or *negative* feedback (case 3). We denote by ϵ_s and by ϵ_f the probability of p_1 's translation (i.e., $T_{1 \rightarrow 2}$) and another foreign translation (i.e., $T_{3 \rightarrow 4} \dots T_{n-1 \rightarrow n}$) being wrong for the attribute in question. Considering the foreign error probabilities as being independent and identically distributed random variables, the probability of not having a foreign translation error along the cycle is

$$(1 - \epsilon_f)^{\|c_i\| - 1}$$

Moreover, *compensating errors*, i.e., series of independent translation errors resulting in a correct mapping, may occur along the cycle of foreign links without being noticed by p_1 (which only has the final result q_n as its disposal). Thus, assuming $T_{1 \rightarrow 2}$ correct and denoting by δ the probability of errors being compensated somehow, the probability of a cycle being positive is

$$(1 - \epsilon_f)^{\|c_i\| - 1} + (1 - (1 - \epsilon_f)^{\|c_i\| - 1})\delta = p(\|c_i\|, \epsilon_f, \delta)$$

while, under the same assumptions, the probability of a cycle being negative is

$$(1 - (1 - \epsilon_f)^{\|c_i\| - 1})(1 - \delta) = 1 - p(\|c_i\|, \epsilon_f, \delta).$$

Similarly, if we assume $T_{1 \rightarrow 2}$ to be incorrect, the probability of a cycle being respectively negative and positive are

$$(1 - \epsilon_f)^{\|c_i\| - 1} + (1 - (1 - \epsilon_f)^{\|c_i\| - 1})(1 - \delta) = q(\|c_i\|, \epsilon_f, \delta)$$

and

$$(1 - (1 - \epsilon_f)^{\|c_i\| - 1})\delta = (1 - q(\|c_i\|, \epsilon_f, \delta)).$$

Combining those equations, the likelihood of receiving a set of cycles $C = c_1, \dots, c_k$, for some positive ($c_i \in C^+$), and for some negative ($c_i \in C^-$), is

$$\begin{aligned} l_1(c_1, \dots, c_k) = & (1 - \epsilon_s) \prod_{c_i \in C^+} p(\|c_i\|, \epsilon_f, \delta) \prod_{c_i \in C^-} (1 - p(\|c_i\|, \epsilon_f, \delta)) \\ & + \epsilon_s \prod_{c_i \in C^-} q(\|c_i\|, \epsilon_f, \delta) \prod_{c_i \in C^+} (1 - q(\|c_i\|, \epsilon_f, \delta)) \end{aligned}$$

Now, we integrate over ϵ_f and δ ,² and let ϵ_s tend toward zero and one in order to obtain the likelihood of the translation $T_{1 \rightarrow 2}$ being semantically correct or incorrect respectively:

$$\begin{aligned} p_1 &= \lim_{\epsilon_s \rightarrow 0} \int_{\delta=0}^1 \int_{\epsilon_f=0}^1 l_1(c_1 \dots c_k) d\epsilon_f d\delta \\ p_2 &= \lim_{\epsilon_s \rightarrow 1} \int_{\delta=0}^1 \int_{\epsilon_f=0}^1 l_1(c_1 \dots c_k) d\epsilon_f d\delta \end{aligned}$$

Finally, we define a variable γ for the relative degree of correctness of the translation:

$$\gamma = \frac{p_1}{p_1 + p_2}$$

²We could take into account density functions here if we have any *a priori* knowledge about those two variables.

Such an analysis may be performed for every outgoing link and every attribute independently, resulting in a series of values γ_i^j indicating the likelihood of the translation T_j being correct for the attribute A_i . Examples of such calculations are given in Section 8.

As for the preceding section, we define now a feature vector and a similarity value to capture the semantic losses along the translation links. Let us suppose that a peer p_k issues a query $q = \pi_{ap}(\sigma_{p(as)}(\mu_{fa}(DB)))$ through a translation link $T_{k \rightarrow j}$. p_k computes a feature vector for q based on the cycle messages it has received as follows:

$$\overrightarrow{FV_{\odot}}(T_{k \rightarrow j}(q)) = (fv_{A_1}^{\odot}, \dots, fv_{A_k}^{\odot})$$

where the feature values $fv_{A_i}^{\odot}$ are defined for every attribute $A_i \in as \cup ap$ as

$$fv_{A_i}^{\odot}(T_{k \rightarrow j}(q)) = \gamma_i^j$$

These values are updated by iteratively multiplying the probabilities for each semantic domain traversed. We consider here that if two translations Ta and Tb have probabilities of a and b respectively and are independent, the overall probability for $(Tb \circ Ta)$ to be correct is ab . Thus, when forwarding a transformed query using a link $T_{k \rightarrow j}$, peer p_k updates each value $fv_{A_i}^{\odot}$ of the feature vector $\overrightarrow{FV_{\odot}}$ it has received along with the transformed query $T_{1 \rightarrow \dots \rightarrow k}(q)$ in this way:

$$fv_{A_i}^{\odot}(T_{1 \rightarrow \dots \rightarrow k, k \rightarrow j}(q)) = fv_{A_i}^{\odot}(T_{1 \rightarrow \dots \rightarrow k}(q)) fv_{A_i}^{\odot}(T_{k \rightarrow j}(T_{1 \rightarrow \dots \rightarrow k}(q)))$$

where γ_i^j values for which we did not receive significant feedback (either because p does not have a representation for A_i or because no cycle message has been received so far) are evaluated to 1. The semantic similarity associated with this vector is

$$S_{\odot}(q, T_{1 \rightarrow \dots \rightarrow n}(q)) = \frac{\overrightarrow{W} \cdot \overrightarrow{FV_{\odot}}}{|\overrightarrow{W}| |\overrightarrow{FV_{\odot}}|}$$

This value starts from 1 (in the semantic domain which the query originates from) and decreases as the query traverses more and more semantically heterogeneous domains.

5.2 Results Analysis

The second mechanism for analyzing the semantic quality of the translations is based on the analysis of the results returned. In [1] we have introduced a method using functional dependencies at the data level in order to assess the quality of translations. This method was based on analyzing to which extent integrity constraints are preserved after translation.

Here we present an alternative approach. We assume that peers annotate documents using meta-data expressed according to our data model. Having sent a query, peers start to receive answer documents with semantically rich content. Based on this content they attempt to assess to which extent the queries expressed at the meta-data level were properly translated and thus led other peers to return the right result documents.

Queries in our meta-data model are thus an intensional way of expressing semantic concepts, whereas extensionally the concepts are related to sets of documents. The problem that we address is of how to arrive at agreed annotation schemes at the intensional level that result in concept definitions that are compatible with the extensional notion of concepts that peers have. In the simplest case (on which we will base our subsequent discussion studies) where relationships among different concepts are not further considered, the meta-data model is used to give names to concepts $c \in \mathcal{C}$. These names can be different for different peers, but the peers should be able to properly translate them.

The extensional notion of concept each peer has is based on methods of content analysis. Here, we do not make any assumption about the methods (e.g., layout analysis, lexicographical analysis, contour-detection, etc., or even simple manual classification) used

to extract meaningful features out of the documents and use them for association with a concept; we simply treat them as high-level abstractions used to unambiguously classify any possible retrieved documents $d \in \mathcal{D}$ into concepts $c \in \mathcal{C}$ using a decision rule \mathcal{R} :

$$\mathcal{R}(d) : \mathcal{D} \rightarrow \mathcal{C}.^3$$

Using their local classification schemes, peers can thus determine for every received document the concept it belongs to.

Imagine now a peer p_i classifying documents according to a rule \mathcal{R}_{p_i} . p_i issues a query q_{p_i} for retrieving documents related to the concept c_k . Upon reception of a document d_{p_j} from a foreign peer $p_j \in N_e(p_i)$, p_i performs the classification operation. Different situations may then occur:

- $\mathcal{R}_{p_i}(d_{p_j}) = c_k$: this is the result p_i was expecting; it is an indication that the outgoing translation link used to forward q_{p_i} to p_j was semantically correct for query q_{p_i} . We treat this as positive feedback (F^+).
- $\mathcal{R}_{p_i}(d_{p_j}) = c_l, c_l \neq c_k$: p_i receives a document related to another class than c_k ; considering that the classification is mostly peer independent, it means that some semantic confusion occurred along the path from p_i to p_j . In this case, we consider this as negative feedback (F^-).

If p_i and p_j are directly connected, the situation gives us a clear indication about the semantic (in)correctness of the translation link $T_{p_i \rightarrow p_j}$ for the attributes in question. Evaluating the mean classification error probability to $eClass$, the probability of the link being correct and incorrect in case of positive feedback are respectively $1 - eClass$ and $eClass$. In case of negative feedback, they become $eClass$ and $1 - eClass$. Also, note that in this case (and for sufficiently small $eClass$) we get a good indication for correcting the mapping, since p_j 's documents classified into concept c_l directly relate to the query q_{p_i} with probability $(1 - eClass)$ (see the experimental evaluation where this is used).

If the two peers are separated by one or more semantic domains, the situation is somewhat more complicated since we have to take into account all the successive links used to forward the query from p_i to p_j . Let us suppose that a peer receives some feedback f_i after the query has gone through $\|f_i\|$ different translation links; reusing some of the equations from the cycle analysis, the probability of receiving a positive feedback assuming the link we are analyzing is correct is

$$p(\|f_i - 1\|, \epsilon_f, \delta)(1 - eClass) + (1 - p(\|f_i - 1\|, \epsilon_f, \delta))\Delta eClass$$

where Δ represents the probability of a document being misclassified and taken as belonging to the class related to the query.

Performing an analysis analogous to that given in Section 5.1 and introducing l_2 as the likelihood of receiving a certain combination of responses for a given error model, we obtain again two values p_3 and p_4 for the likelihood of the translation being semantically correct or not:

$$p_3 = \lim_{\epsilon_s \rightarrow 0} \int_{\delta=0}^1 \int_{\epsilon_f=0}^1 l_2(c_1 \dots c_k, e) d\epsilon_f d\delta$$

$$p_4 = \lim_{\epsilon_s \rightarrow 1} \int_{\delta=0}^1 \int_{\epsilon_f=0}^1 l_2(c_1 \dots c_k, e) d\epsilon_f d\delta$$

Defining κ_i^j as the likelihood of the translation T_j being correct for attribute A_i with value

$$\kappa = \frac{p_3}{p_3 + p_4}$$

³In a more general setting this could be a probabilistic rule.

we obtain again another feature vector:

$$\overrightarrow{FV_{\rightleftharpoons}}(T_j(q)) = (fv_{A_1}^{\rightleftharpoons}, \dots, fv_{A_k}^{\rightleftharpoons})$$

whose feature values $fv_{A_i}^{\rightleftharpoons}$ are defined for every attribute $A_i \in as \cup ap$ as

$$fv_{A_i}^{\rightleftharpoons}(T_j(q)) = \kappa_i^j$$

and where, again, we evaluate missing values to 1 and we update the vectors iteratively:

$$fv_{A_i}^{\rightleftharpoons}(T_{1 \rightarrow \dots \rightarrow k, k \rightarrow j}(q)) = fv_{A_i}^{\rightleftharpoons}(T_{1 \rightarrow \dots \rightarrow k}(q))fv_{A_i}^{\rightleftharpoons}(T_{k \rightarrow j}(T_{1 \rightarrow \dots \rightarrow k}(q))).$$

The associated semantic similarity is, as expected:

$$S_{\rightleftharpoons}(q, T_{1 \rightarrow \dots \rightarrow n}(q)) = \frac{\overrightarrow{W} \cdot \overrightarrow{FV_{\rightleftharpoons}}}{|\overrightarrow{W}| |\overrightarrow{FV_{\rightleftharpoons}}|}.$$

6 Gossiping Algorithm

At this point, we have four measures (S_σ , S_π , S_\cup and S_{\rightleftharpoons}) for evaluating the losses due to the translations. We will now make use of these values to decide whether or not it is worth forwarding a query to a foreign semantic domain.

First, we require the creator of a query to attach a few user-defined or generated values to the query it issues:

- The weights \overrightarrow{W} pondering the importance of the attributes in the query
- The respective selectivity of the selection attributes \overrightarrow{sel}
- The minimal values $\overrightarrow{S_{min}}$ for the similarity measures under which a transformed query is so deteriorated that it can no longer be considered as equivalent to the original query.

We extend the format of a query message to include these values as well as the iteratively updated feature vectors:

$$query(id, q, p, TT, \overrightarrow{W}, \overrightarrow{sel}, \overrightarrow{S_{min}}, \overrightarrow{FV_\sigma}, \overrightarrow{FV_\pi}, \overrightarrow{FV_\cup}, \overrightarrow{FV_{\rightleftharpoons}}).$$

Now, upon reception of a query message, we require a peer to perform a series of tasks:

1. detect any semantic cycles
2. check whether or not this query has already been received
3. in case the local neighborhood has not received the query, forward it to the local neighborhood
4. return potential results

and, for each of its outgoing translation links:

5. apply the translation to the query
6. update the similarity measures for the transformed query
7. perform a test for each of the feature vectors: $similar(\overrightarrow{FV_i})$ evaluates to 1 if

$$\frac{\overrightarrow{W} \cdot \overrightarrow{FV_i}}{|\overrightarrow{W}| |\overrightarrow{FV_i}|} \geq S_{min,i}$$

that is if the semantic similarity is greater or equal to the specified minimal value, and to 0 otherwise

8. forward the query using the link if all $similar()$ tests succeed (i.e., evaluate to 1).

This algorithm ensures that queries are forwarded to a sufficiently large set of peers capable of rendering meaningful feedback without flooding the entire network.

7 Case Study

Several experiments were conducted following the approach presented above. This section presents one of them as a case study detailing how the aforementioned heuristics may be deployed in a concrete setting.

Seven people from our laboratory were first asked to design a simple XML document containing some project meta-data. The outcome of this voluntary imprecise task definition was a collection of structured documents lacking common semantics though overlapping partially for a subset of the embraced meta-data (e.g., *name of the project* or *start date*). Viewing these documents as seven distinct semantic domains in a decentralized setting, we then produced a random graph connecting the different domains together with series of translation links (the resulting topology is depicted in Figure 4).

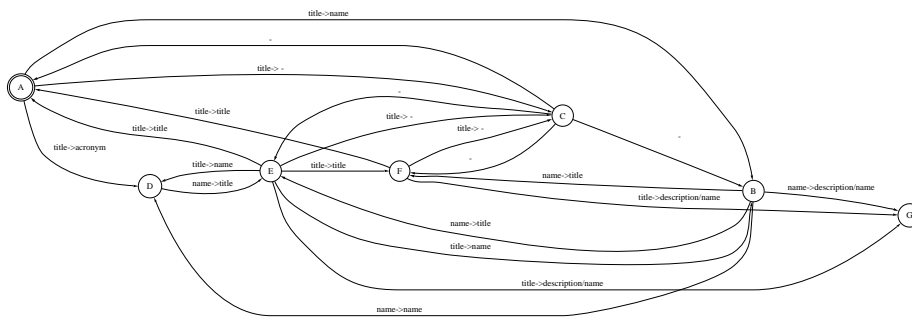


Figure 4: The Semantic Graph

Translations were formulated as XQuery expressions in such a way that they strictly adhere to the principles stipulated above (see Section 3). As an example, Figure 5 presents two different documents as well as a simple query translation using the translation expression *T12*. Providing the authors with the required documents, we asked them to write the translations for every link departing from their domain (thus, p_A was asked to provide us with the translation to p_B , p_C and p_D). Finally, using the IPSI-XQ XQuery libraries [8] and the Xerces [23] XML parser, we built a query translator capable of handling and forwarding the queries following the gossiping algorithm.

We focus now on a single node, p_A , and on a single-attribute query issued by p_A to obtain all the titles of the different projects, namely:

```
Query =
FOR $project IN "project_A.xml" /*
RETURN
<title>$project/title </title>
```

Note that the weight and selectivity values attached to the query do not matter here, as a single attribute is concerned. Moreover we will not consider S_σ and S_{\Rightarrow} for the rest of this study (here, S_σ always evaluates to 1 because there is no selection attribute, and so does S_{\Rightarrow} since we do not return any document). The other minimal values are set to 0.5.

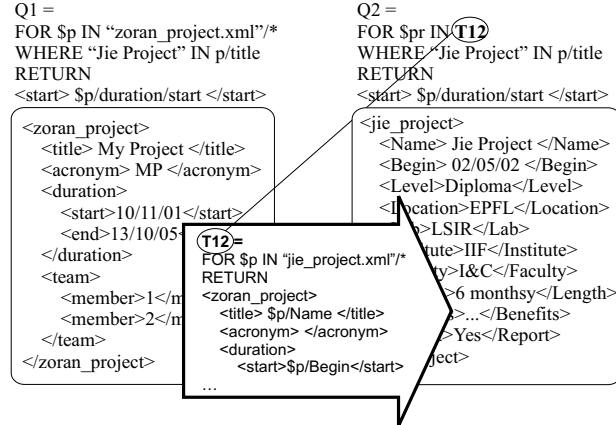


Figure 5: The Translation Mechanism

All the domains have some representation for the title of the project (usually referred to as *name* or *title*, see Figure 4 where the translations for the attribute *title* are represented on top of the link), except p_C which only considers a mere *ID* for identifying the projects. Following the gossiping algorithm, p_A first attempts to transmit the query to its direct neighbors, i.e., p_B , p_C and p_D . p_B and p_D in turn forward the query to the other nodes, but p_C will in fact never receive the query: As p_C has no representation for the *title*, the only projection attribute would be lost in the translation process from p_A to p_C , lowering S_π to 0.

Let us now examine the semantic similarity S_\odot . For the topology considered, thirty-one semantic cycles could be detected by p_A in the best case. As the query never traverses C , only eight cycles remain (Table 1 lists those cycles). We use now the formulas from Section 5; For its first outgoing link (i.e., the link going from p_A to p_B), p_A receives five positive cycles, raising the semantic similarity measure for this link and the attribute considered to 0.79.⁴ p_A does not receive any semantically significant feedback for its second outgoing link $Tp_A \rightarrow p_B$, which is anyway handled by the syntactic analysis. Yet, it receives three negative cycles for its last outgoing link $Tp_A \rightarrow p_D$. This link is clearly semantically erroneous, mapping *title* onto *acronym*. This results in p_A excluding the link for this attribute, the semantic similarity dropping to 0.26.

Cycle	$T_{p_A \rightarrow p_C}$	Erroneous	$T_{p_B \rightarrow p_D}$	Erroneous
A, B, D, E, A	+		-	
A, B, D, E, F, A	+		-	
A, B, E, A	+		+	
A, B, E, F, A	+		+	
A, B, F, A	+		+	
A, D, E, A	-		+	
A, D, E, B, F, A	-		+	
A, D, E, F, A	-		+	

Table 1: Cycles Resulting In Positive(+) or Negative(-) Feedback

The situation may consequently be summarized in this way: p_A restrains from sending the query through p_C because of the syntactic analysis (too much information lost in the translation process) and excludes p_B because of the high semantic dissimilarity.

The situation somewhat changes if we correct the erroneous link and add a mistake for the link $Tp_B \rightarrow p_D$. For the attribute considered, the semantic similarity drops to 0.69 for

⁴Remember that we did not make any assumption regarding the distribution of erroneous links. In this case, the positive feedback received may well come from a series of compensating errors.

the outgoing link to p_B (two long cycles are negative, see third column in Table 1). Even though it is not directly connected to an erroneous link, p_A senses the semantic incompatibilities affecting some of the messages traversing p_B . It will continue to send queries through this link, as long as it receives positive feedback at least.

8 Experimental evaluation

In the preceding section, we have evaluated the Chatty Web by examining query forwarding in a small network of static translations generated by a group of users. In contrast to this, we detail below simulation experiments where semantic gossiping is used to automatically reach semantic agreement in large networks of computer-generated and dynamic translation links. Such an approach in place could for example be used to derive basic, common ontologies from a dynamic system with heterogeneous schemas, or to gradually refine existing networks of translations. The initial results interpreted below provide promising evidence that it is worth pursuing further research along these lines and highlight some of the issues to be addressed in that course.

8.1 Experimental setup

The setup we used in the experiments is as follows: We assume a network of peers representing individual semantic domains. Peers share similar concepts, i.e., operate in a certain semantic domain (for example, biological databases) inside the network. They share annotated documents (or data) relative to those concepts, but refer to concepts using different names (they denominate the concepts differently). From this basic setup, we attempt to create global interoperability by applying semantic gossiping techniques using purely pairwise, local translations.

Here is the exact description of the process: first, we create a topology of $nPeers$ peers $p_1 \dots p_{nPeers}$, each of them connected through translation links to $nTLinks$ other peers. The peers share $nConcepts$ concepts $C_1 \dots C_{nConcepts}$, but use distinct names to refer to them. Thus we study the problem of peers sharing the same concepts but lacking knowledge of how to refer to them by names. This is somewhat similar to the approach taken in [25], without aiming at universally agreed upon names. Each peer p_i uses its own set of names $n_{p_i}^1 \dots n_{p_i}^{nConcepts}$ to identify the concepts. We write $(n_{p_i}^k \mapsto C_l)$ when peer p_i uses name $n_{p_i}^k$ to refer to concept C_l . These names may be seen in our data model, for example, as attribute names indicating the presence of a concept in a document. Also, peers can verify whether a document belongs to a concept or not.

We generate mappings $\mu(n_{p_i}^1 \dots n_{p_i}^{nConcepts})$ for every translation link relating names from the first peer to names from the second peer, with every name used by the first peer mapped onto a distinct name used by the second peer. For every mapping in every translation link, we say that the mapping is correct if and only if the two names bound by the mapping actually refer to the same concept, that is if

$$\mu(n_{p_1}^i) = n_{p_2}^j \wedge n_{p_1}^i, n_{p_2}^j \mapsto C_k.$$

Thus, random mappings would only have a probability of $\frac{1}{nConcepts!}$ of being correct in such a setting. In the experiments, we generate a fraction $eRate$ of erroneous mapping initially.

Unless specified otherwise, we use small-world graphs to interconnect peers with translation links; small-world topologies have been extensively applied to model computer networks or social behaviors. They are typically characterized by high clustering coefficients (average fraction of pairs of neighbors of a node that are also neighbors of each other) and relatively small path length (average minimal distance between two nodes). In the

following, we generate graphs with an average clustering coefficient of 0.1 and with 10% shortcuts (i.e., links rewired to any random peer in the network).

Starting from the original topology, we apply semantic gossiping techniques iteratively in order to detect and rectify erroneous mappings. At every simulation step, each peer selects one of its names randomly and issues a query about this name (i.e., the query consists of a projection on one attribute: the name selected). The query is propagated to the other peers (semantic domains) in a Gnutella fashion with a low time-to-live (TTL) value.

The syntactic analysis for this simplistic type of query is straightforward: peers forward the query through an outgoing translation link if there exists a mapping translating the local name used in the query (projection attribute) into another name for the foreign peer. Now, for detecting and repairing erroneous translation links, we slightly modify the semantic analysis; we forward queries irrespectively of the results of previous semantic analyses in order to get as many evidences as possible, and use these results to reach semantic agreements by gradually modifying mappings.

Before taking a closer look at the final results, we will evaluate in the following sections each of the semantic analyses (cycle and result analysis) separately to underline their specificities.

8.2 Cycle Analysis

Let us start with the cycle analysis. For every iteration step, peers randomly choose a name, send a query for this name and analyze the cycle messages they get in return.

Here, we do not only estimate the correctness of the actual mapping as explained in Section 5.1, but determine also which of the possible mappings is most likely correct and adopt it as a new mapping. Therefore, peers view mappings resulting of returning queries as new mapping candidates. Consider for example Figure 6, where peer p_A systematically receives n_A^1 mapped onto n_B^2 in returning queries (negative feedback). In addition to evaluating the correctness of the current mapping, p_1 considers other mappings as well. It takes the potential mapping receiving the highest probability of being correct and in case this probability is above 50% and the most probable mapping is different from the current mapping it changes it. In this example, p_1 evaluates the correctness of mapping n_A^1 onto n_B^1 , and might consider to modify it to a mapping n_A^1 onto n_B^1 .

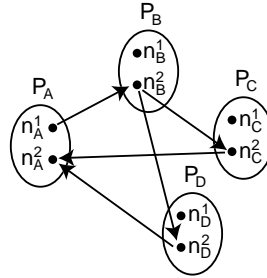


Figure 6: New Mapping Candidates

As indicated in Section 5.1, preexisting knowledge on the distribution of error probabilities δ and ϵ_f may be used in the computation of semantic similarity. δ , the probability of series of different errors to compensate along a cycle, is approximated to $(nConcepts - 1)^{-1}$, which is the probability of the last erroneous link in the cycle to map to the original name and thus to correct previous errors.

We estimate ϵ_f with standard maximum-likelihood techniques applied to the feedback information we receive. From the probability of receiving a positive cycle of length $\|c_i\|$ knowing that the error probability of a translation link is ϵ ,

$$(1 - \epsilon)^{\|c_i\|} + (1 - (1 - \epsilon)^{\|c_i\|})\delta,$$

and from its negative counterpart, we derive the density function for the likelihood of ϵ :

$$L(\epsilon_f \| C) = K \prod_{c_i \in C^+} (1 - \epsilon_f)^{\|c_i\|} + (1 - (1 - \epsilon_f)^{\|c_i\|}) \delta \prod_{c_j \in C^-} (1 - (1 - \epsilon)^{\|c_j\|}) (1 - \delta)$$

where K is a normalizing constant. The local maximum of this function over $[0, 1]$ gives a good approximation of ϵ_f , supposing we have sufficient feedback information.

What is the result of such a process in the long run? It of course depends on the initial setting. In the end, this method attempts to obtain a mapping consensus based on the different feedback cycles detected in the network. Considering a high density of links and relatively few erroneous links, the method converges (i.e., repairs all erroneous mappings) rapidly, since peers can base their decisions on numerous and meaningful feedback cycles. For settings where links are scarce, peers do not have sufficient information for making sensible choices, and results may diverge.

The figures below show experimental results for topologies where $nPeers = 25$, $eRate = 0.1$, $nConcepts = 4$, $TTL = 5$ and $nTlinks = 5$ and where one of those parameters varies. All the curves are actually average curves over ten consecutive runs. At every step, each peer sends a query picking a random concept for every outgoing edge and modifies the mappings depending on the results of the analysis explained above. Steps are represented on the x -axis. The graph shows the evolution of the percentage of erroneous mappings, starting at a rate $eRate$ initially. Clearly, the outcome depends on the density of links, which directly impacts on the number of cycles we have at our disposal for taking mapping decisions (see Figure 7). For $nTlinks = 4$ and the topology considered, we get on average only one positive feedback per mapping candidate, which is obviously insufficient to take sensible decisions. For $nTlinks = 5$ and $nTlinks = 6$, the value raises to respectively 1.8 and 2.9 and most of the erroneous mappings get corrected after ten iterations. Finally, for $nTlinks = 7$, we get enough evidences (4.5 on average per mapping candidate) for correcting all the erroneous links, thus reaching a perfect semantic agreement, in eight steps.

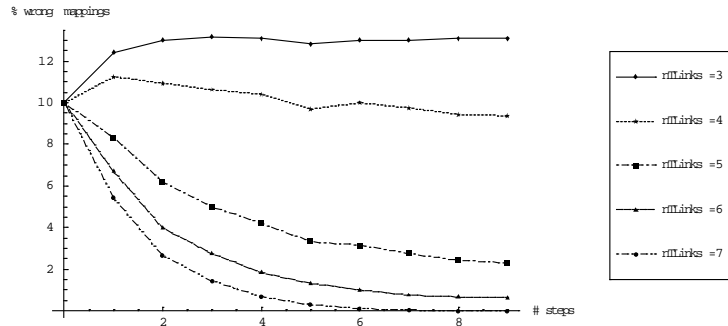


Figure 7: Sensitivity to the number of outgoing edges

Similar considerations may be drawn for variable TTLs. Figure 9 shows results using the same parameters as before, but this time for a fixed number of outgoing edges ($nTlinks = 4$) and TTLs ranging from 3 to 6. Again, for low values, peers do not gain sufficient feedback information to correct mappings. Starting from $TTL = 4$ (1.8 positive feedback per decision), peers receive sufficient information to correct more than 75% of erroneous mappings after nine iterations. Low-connectivity networks may thus benefit from increasing the TTL value of their queries in order to the peers to get sufficient feedback information.

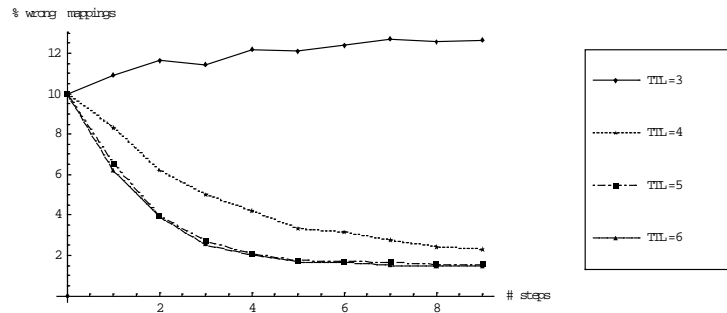


Figure 8: Sensitivity to the TTL

Our approach is rather insensitive to variations of the initial error rate (see Figure 9) until a certain threshold, where too many bad links are present initially to reach a correct consensus based on the feedback cycles. Finally, it is worth mentioning that the approach scales very well with the number of nodes; This is not surprising, considering that the method relies solely on local interactions (no central component or computation) and that the clustering coefficient of the network is relatively high. See Figure 10 where the experiments were conducted for networks ranging from 50 to 800 peers without fundamental results variations. The small deviations are due to the *shortcuts* in the small world topology which connect two random peers in the network. The bigger the graph, the less likely it is that these links can be used to form cycles within a certain neighborhood.

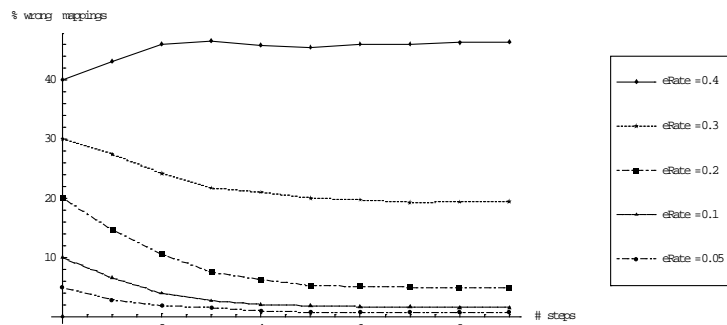


Figure 9: Sensitivity to the initial error rate

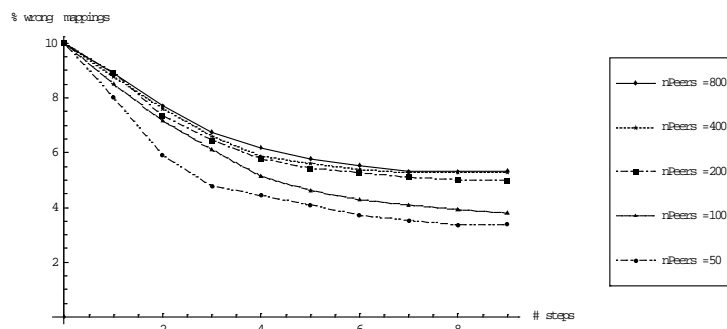


Figure 10: Scalability

8.3 Results Analysis

We consider now the second part of the analysis, where peers analyze and categorize documents they receive. The process is as follows: At every step, peers first issue a couple of queries with a high TTL for estimating the error rate as explained above. Then, for each of their outgoing links, peers pick a concept randomly and issue a query asking for documents relating to that concept. They receive in return series of documents they analyze according to what is described in Section 5.2. They modify the mapping they have used to forward the query with the most probable mapping if it has a likelihood of at least 0.5 of being correct.

For the simulations, we consider a fixed set of documents scattered randomly among the peers. Documents are all assigned to concepts. Each document owner has a probability ($eClass$) of misclassifying a document by relating it to a wrong concept. Peers do not try to evaluate the probability of misclassification, but arbitrarily use a fixed, low value instead (5% in the following experiments). For our setting, Δ , the probability that a misclassified document is seen as relating to another specific concept, is equal to $(nConcepts - 1)^{-1}$.

Unless specified otherwise, we used below a network of 50 peers sharing in total 100 document, 2 outgoing translation links per peer, 4 concepts, a TTL of 3, an initial error rate of 10%, and a probability of 10% of misclassifying documents.

First, it is interesting to remark that this approach is very robust vis--vis the initial error rate, mainly because a few links suffice here to get meaningful results (thus, the very low TTL), while whole link cycles were needed previously. See Figure 11 where the initial percentage of wrong mappings vary from 10% to 50%.

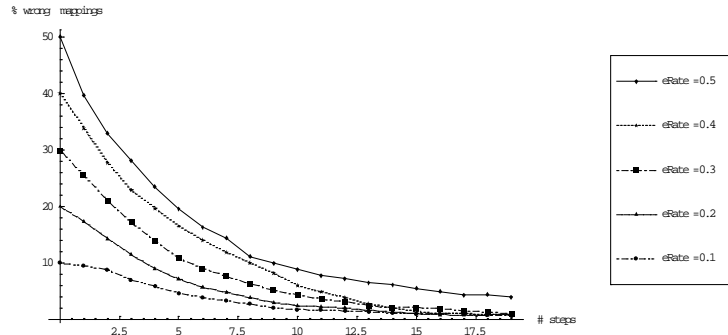


Figure 11: Sensitivity to initial error rate

Nevertheless, the approach is rather sensitive to the rate of misclassification of documents, as shown in Figure 12. This is especially true since we do not try to evaluate this parameter but consider a mere fixed value.

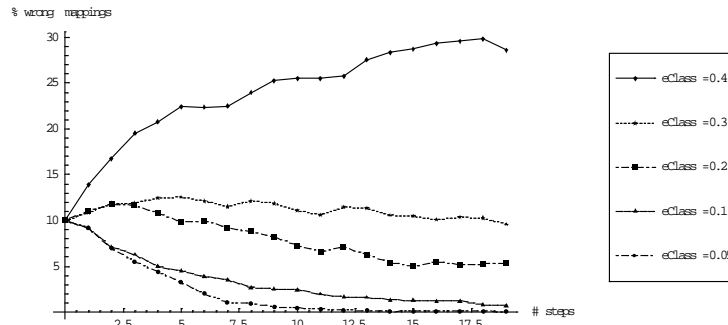


Figure 12: Sensitivity to misclassification rate

The approach taken here is completely local, and does not take into consideration any global behavior, and scales with the number of peers quite naturally (see Figure 13). Here, we increase the number of documents linearly with the number of peers, to keep the average number of documents per peer constant. This number is essential to this analysis, since it is directly proportional to the number of evidences a peer gets for every query. Take a look at Figure 14, where this effect is depicted: Peers start having trouble correcting the mappings as they get less and less documents returned for their queries (documents scarcity).

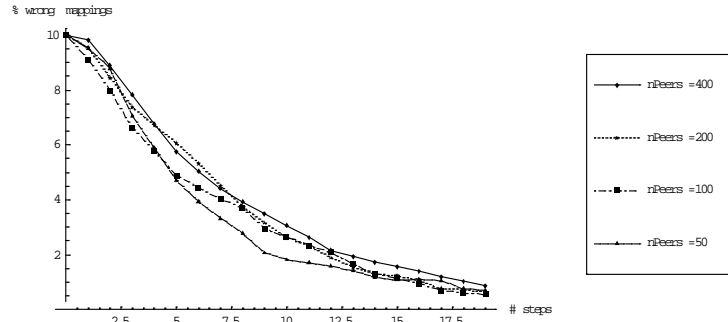


Figure 13: Scalability

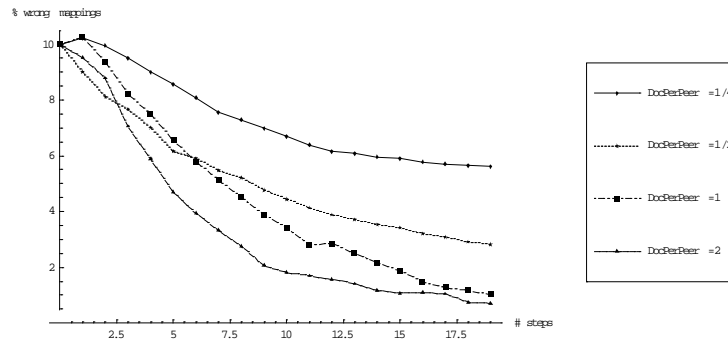


Figure 14: Sensitivity to number of documents

8.4 Combined Results

Below, we show some results where the two mechanisms were used in parallel. Many possibilities exist for combining the two analyses; we chose a very simple one: at each step, every peer performs first a results step (modifying a few mappings depending on the results returned) and sequentially performs then a cycle step (trying to reach some local agreement on mappings based on cycle feedback). The results for topologies with 25 peers, 4 concepts, 2 outgoing edges, TTLs of 3 (results) or 6 (cycles) and a varying error rates on initial mappings are depicted on Figure 15. This method takes more time to converge than the two analyses taken separately; this is because the analyses keep interfering with each other until some state is reached that is consistent from both a cycle and a feedback point of view. Also, note that the overall results outperform in the end the two analysis taken separately (e.g., more than 95% of erroneous mappings corrected after 50 steps with 50% erroneous mappings initially).

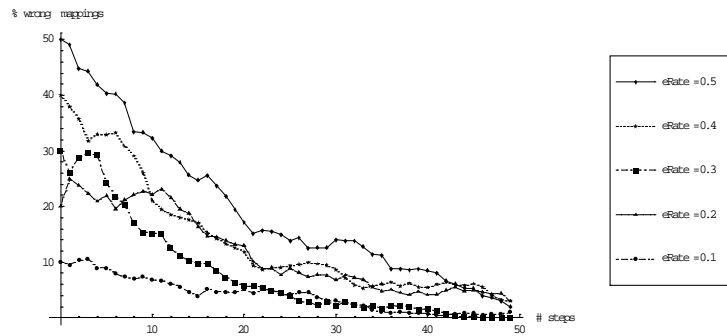


Figure 15: Combined results, varying initial error rate

9 Implementation framework

All the tasks of the Chatty Web approach have been mapped onto an implementation architecture which uses a meta-data model expressed in XML and XQuery as the language to translate among schemas. The framework assumes the availability of a communication infrastructure, for example, simple web access via HTTP or a P2P infrastructure such as JXTA [9]. However, we are not bound to any specific communication infrastructure. All we require is access to the relevant schema data and to query information and results. This can easily be achieved by a standard abstraction layer that maps a specific communication infrastructure's interface to the one we require. Since this is a fairly standard software engineering task we omit it in the following discussion. Based on these assumptions, Figure 16 shows the standard architecture used for semantic gossiping in the Chatty Web.

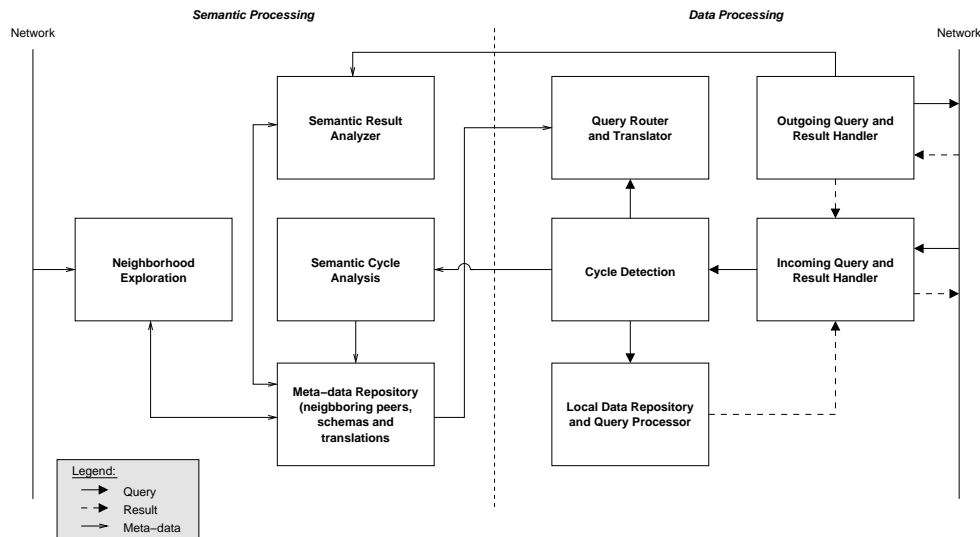


Figure 16: Architecture for semantic gossiping

Incoming queries are registered at and handled by the *Incoming Query and Result Handler* whose task is to communicate with other peers, to forward the query for further processing and to gather partial results which it uses to assemble the final result of a specific query. The next step then is to detect whether a cycle has occurred. If so, semantic analysis of the cycle is triggered. Otherwise, the query is processed, first by querying the local

database and then by handing it over to the *Query Router and Translator* to collect results from other peers.

For this purpose the *Query Router and Translator* inquires for possible translations, evaluates the quality of the resulting queries, and if it is above a defined threshold, forwards the query to the respective peer in a different semantic domain. Queries are forwarded by the *Outgoing Query and Result Handler* which is also in charge of collecting the results and forwarding the results to the *Incoming Query and Result Handler* which returns them to the original requester. Additionally, it provides input data for semantic result analysis.

This is the main data processing flow of the architecture. In parallel, partly triggered by the ongoing data processing, there is also semantic processing as depicted on the right side of Figure 16. Its main tasks are semantic analyses of results based on the existing knowledge of schemas and their relationships and the semantic analyses of detected cycles. The results of these analyses are integrated again into the system's knowledge base and provide the basic decision criteria for query routing.

Additionally, the knowledge base is updated and improved by exploring the peer's neighborhood and detecting new schemas and translations. The meta-data repository will try to infer further translations and present new ones for human analysis or apply for actively detecting semantic agreements in an automatic way.

10 Related Work

A number of approaches for making heterogeneous information sources interoperable are based on mappings between distributed schemas or ontologies without making the canonical assumption on the existence of a global schema.

For example, in OBSERVER [17] each information source maintains an ontology, expressed in description logics, to associate semantics with the information stored and to process distributed queries. In query processing, they use local measures for the loss of information when propagating queries and receiving results. Similarly to OBSERVER, KRAFT [22] proposes an agent-based architecture to manage ontological relationships in a distributed information system. Relationships among ontologies are expressed in a constraint language. [2] propose a model and architecture for managing distributed relational databases in a P2P environment. They use local relational database schemas and represent the relations between those with domain relations and coordination formulas. These are used to propagate queries and updates. The relationships given between the local database schemas are always considered as being correct. In [21] a probabilistic framework for reasoning with assertions on schema relationships is introduced. Thus their approach deals with the problem of having possibly contradictory knowledge on schema relationships. [18] propose an architecture for the use of XML-based annotations in P2P systems to establish semantic interoperability.

An approach to self-organizing vocabularies is described in [25]. A set of agents communicate by randomly associating a fixed set of words to a fixed set of meanings (which is called a vocabulary but in fact is an ontology) and repeatedly evaluate how successful their communicative acts have been. Depending on the success the binding between a word and a meaning is maintained or replaced by a new random coupling. The decision is based on sigmoid functions so that the probability of change quickly decreases if the majority of agents uses the same coupling. This approach is related to the method of cycle analysis we use and simulate in Section 8. However, it does not employ result analysis. Nevertheless [25] shows that semantic agreements are reached rather quickly. The additional result analysis we perform may help to speed up convergence speed and increase the scalability and robustness of the self-organization process. It is interesting to note that [25] shows that an increased numbers of agents, words, and meanings does not lead to combinatorial explosion but implosion. This is due to the fact that the increasing number of words with consistent meaning narrow the selection space drastically. This phenomenon is similar to

the combinatorial implosions described by Kauffman [13] for the clustering and interconnection of autocatalytic networks.

Edutella [19] is a recent approach to apply the P2P architectural principle to build a semantically interoperable information system for the educational domain. The P2P principle is applied at the technical implementation level whereas logically a commonly shared ontology is used. The original design of Edutella which is based on Gnutella is changed to a super peer network approach in [20] which offers better scalability and provides sophisticated routing and clustering strategies based on the meta-data schemas attributes and ontologies used. This approach includes a methodology for mediation between local schemas at super peers which enables super peers to route queries and combine results from different semantic domains into one result. It employs transformation rules, so-called correspondences, which have already been used in mediator-based information systems [26]. *Query Response Assertions* [16] and *Model Correspondences* [3] are used to express correspondences between heterogeneous schemas.

The Piazza system [10] defines a mapping language to specify mappings between sets of XML or RDF data sources that tries to take into account both domain and document structure in the mediation process. The transitive closure of these mappings is used to provide a query answering algorithm over the graph of data source defined by the mappings. Piazza's approach is complementary to our approach since it assumes the existence of pairwise mappings between data sources and uses these mappings for answering queries while we try to detect the quality of mappings in terms of an overall agreement among nodes (which can also be seen as a form of transitive closure). However, the mapping language of Piazza together with its query rewriting and query answering methods could also be used in the Chatty Web approach for more expressive mappings and improved query routing.

Approaches for automatic schema matching—see [24] for an overview—would ideally support the approach we pursue in order to generate mappings in a semi-automated manner. In fact, we may understand our proposal as extending approaches for matching two schemas to an approach matching multiple schemas in a networked environment. One example illustrating how the schema matching process could be further automated at the local level is introduced in GLUE [6] which employs machine learning techniques to assist in the ontology mapping process. GLUE is based on a probabilistic model, employs similarity measures and uses a set of learning strategies to exploit ontologies in multiple ways to improve the resulting mappings.

Finally, we see our proposal also as an application of principles used in Web link analysis, such as [14], in which local relationships of information sources are exploited to derive global assessments on their quality (and eventually their meaning).

11 Conclusions

Semantic interoperability is a key issue on the way to the Semantic Web which can push the usability of the web considerably beyond its current state. The success of the Semantic Web, however, depends heavily on the degree of global agreement that can be achieved, i.e., global semantics. In this paper we have presented an approach facilitating the fulfilment of this requirement by deriving global semantics (agreements) from purely local interactions/agreements. This means that explicit local mappings are used to derive an implicit global agreement. We see our approach as a complementary effort to the on-going standardization in the area of semantics which may help to improve their acceptance and application by augmenting their top-down approach with a dual bottom-up strategy. We have developed our approach in a formal model that is built around a set of instruments which enable us to assess the quality of the inferred semantics. To demonstrate its validity and practical usability, the model is applied in a simple yet practically relevant case study. Also, series of experimental results legitimate our claims and illustrate our interests in pursuing research aiming at a better understanding of network-related properties fostering

semantic interoperability.

References

- [1] Karl Aberer, Philippe Cudré-Mauroux, and Manfred Hauswirth. The Chatty Web: Emergent Semantics Through Gossiping. In *Proceedings of the 12th International World Wide Web Conference*, 2003.
- [2] P.A. Bernstein, F. Giunchiglia, A. Kementsietsidis, J. Mylopoulos, L. Serafini, and I. Zaihrayeu. Data management for peer-to-peer computing: A vision. In *Workshop on the Web and Databases (WebDB)*, 2002.
- [3] S. Busse. *Model Correspondences in Continuous Engineering of MBIS*. PhD thesis, Logos Verlag, 2002.
- [4] Clip2. The Gnutella Protocol Specification v0.4 (Document Revision 1.2), Jun. 2001. http://www9.limewire.com/developer/gnutella_protocol_0.4.pdf.
- [5] Mike Dean, Dan Connolly, Frank van Harmelen, James Hendler, Ian Horrocks, Deborah L. McGuinness, Peter F. Patel-Schneider, and Lynn Andrea Stein. OWL Web Ontology Language 1.0 Reference, 2002. W3C Working Draft 29 July 2002. <http://www.w3c.org/TR/owl-ref/>.
- [6] AnHai Doan, Jayant Madhavan, Pedro Domingos, and Alon Halevy. Learning to map between Ontologies on the Semantic Web. In *WWW2002*, 2002.
- [7] Dieter Fensel, Ian Horrocks, Frank van Harmelen, Stefan Decker, Michael Erdmann, and Michel C. A. Klein. OIL in a Nutshell. In *EKAW*, pages 1–16, 2000.
- [8] FhG-IPSI. IPSI-XQ - The XQuery Demonstrator, 2002.
- [9] Li Gong. JXTA: A Network Programming Environment. *IEEE Internet Computing*, 5(3):88–95, May/June 2001.
- [10] Alon Y. Halevy, Zachary G. Ives, Peter Mork, and Igor Tatarinov. Piazza: Data Management Infrastructure for Semantic Web Applications. In *Proceedings of the Twelfth International World Wide Web Conference*, 2003.
- [11] Ian Horrocks. DAML+OIL: a Description Logic for the Semantic Web. *IEEE Data Engineering Bulletin*, 25(1):4–9, 2002.
- [12] Richard Hull. Managing Semantic Heterogeneity in Databases: A Theoretical Perspective. In *PODS*, pages 51–61, 1997.
- [13] Stuart A. Kauffman. *The Origins of Order - Self-Organization and Selection in Evolution*. Oxford Univ. Press, 1993.
- [14] Jon M. Kleinberg. Hubs, authorities, and communities. *ACM Computing Surveys*, 31(4es), 1999.
- [15] M. Koubarakis, C. Tryfonopoulos, P. Raftopoulou, and T. Koutris. Data Models and Languages for Agent-Based Textual Information Dissemination. In *CIA 2002*, pages 179–193, 2002.
- [16] U. Leser. *Query Planning in Mediator Based Information Systems*. PhD thesis, TU Berlin, 2002.

- [17] Eduardo Mena, Vipul Kashyap, Amit P. Sheth, and Arantza Illarramendi. OB-SERVER: An Approach for Query Processing in Global Information Systems based on Interoperation across Pre-existing Ontologies. *Distributed and Parallel Databases*, 8(2):223–271, 2000.
- [18] A. Mukherjee, B. Esfandiari, and N. Arthorne. A Peer-to-peer System for Description and Discovery of Resource-sharing Communities. In *RESH*, 2002.
- [19] Wolfgang Nejdl, Boris Wolf, Changtao Qu, Stefan Decker, Michael Sintek, Ambjörn Naeve, Mikael Nilsson, Matthias Palmér, and Tore Risch. EDUTELLA: a P2P networking infrastructure based on RDF. In *WWW2000*, pages 604–615, 2000.
- [20] Wolfgang Nejdl, Martin Wolpers, Wolf Siberski, Christoph Schmitz, Mario Schlosser, Ingo Brunkhorst, and Alexander Löser. Super-Peer-Based Routing and Clustering Strategies for RDF-Based Peer-To-Peer Networks. In *Proceedings of the Twelfth International World Wide Web Conference*, 2003.
- [21] Aris M. Ouksel and Iqbal Ahmed. Ontologies are not the Panacea in Data Integration: A Flexible Coordinator to Mediate Context Construction. *Distributed and Parallel Databases*, 7(1):7–35, 1999.
- [22] Alun D. Preece, Kit-Ying Hui, W.A. Gray, Trevor J. M. Bench-Capon Philippe Marti, Zhan Cui, and Dean Jones. Kraft: An Agent Architecture for Knowledge Fusion. *IJCIS*, 10(1–2):171–195, 2001.
- [23] Apache XML Project. Xerces Parser, 2002.
- [24] Erhard Rahm and Philip A. Bernstein. A survey of approaches to automatic schema matching. *VLDB Journal*, 10(4):334–350, 2001.
- [25] L. Steels. Self-organising vocabularies. In *Proceedings of Artificial Life V*, 1996.
- [26] Gio Wiederhold. Mediators in the Architecture of Future Information Systems. *IEEE Computer*, 25(3):38–39, 1992.