

# Semantic-based Connectivity in a Peer Data Management System

Carlos Eduardo Santos Pires<sup>1</sup>, Ana Carolina Salgado<sup>1</sup>, Bernadette Farias Lóscio<sup>2</sup>

<sup>1</sup>Centro de Informática – Universidade Federal de Pernambuco (UFPE)  
Caixa Postal 7851 - 50732-970 - Recife, PE - Brazil

<sup>2</sup>Departamento de Computação - Universidade Federal do Ceará (UFC)  
Fortaleza, CE - Brazil

{cesp,acs}@cin.ufpe.br, bernafarias@lia.ufc.br

**Student:** Carlos Eduardo Santos Pires

**Advisor:** Ana Carolina Salgado

**Co-Advisor:** Bernadette Farias Lóscio

**Program:** Ph.D. Program in Computer Science - UFPE, Brazil

**Admission Year in the Ph.D. Degree Program:** 2005

**Conclusion Expected by:** February 2009

**Concluded Stages:** Credits in disciplines; raising of the state of the art on data management in P2P systems; delimitation of thesis research scope; definition of a Peer Data Management System (PDMS) architecture; writing and submission of scientific papers about the proposed PDMS; writing and presentation of qualifying exam and thesis proposal.

***Abstract.** In the last years, Peer Data Management Systems (PDMS) came into the focus of research as a natural extension to distributed databases in the peer-to-peer (P2P) context. PDMS are P2P networks where every contributing peer has its own data and intends to share parts of its data with other peers. The main data management issues that a PDMS must deal with when sharing structured and semi-structured data are the identification of schema mappings and the query processing. Such issues can be facilitated if an efficient strategy for peer connectivity is employed, for example, grouping semantically similar peers within semantic communities. In this case, better quality semantic mappings can be generated, enhancing query results. Also, queries can be addressed only to relevant peers, minimizing network traffic and improving system scalability. In this work, we propose an approach for semantic-based peer connectivity in a PDMS. Such PDMS provides a P2P infrastructure that facilitates identification of semantic mappings and query processing. The data shared by peers are represented through ontologies which are employed for grouping peers within semantic communities and clusters. Due to the dynamic behavior of peers, approaches for load balancing of semantic clusters and fault tolerance are also proposed.*

**Keywords:** Peer-to-Peer, PDMS, Connectivity, Ontology, Semantic Similarity

## 1. Introduction

In the last years, the development of solutions for data integration has been important in several environments, including distributed systems in the Web and, most recently, peer-to-peer (P2P) systems. The main goal of initial P2P systems was the sharing of unstructured data, such as music files. Recently, a new category of P2P systems, named Peer Data Management Systems (PDMS) [Tatarinov et al. 2003], has emerged to increase the functionality of initial P2P systems. PDMS enable the sharing of structured and semi-structured data. Also, they offer a richer representation for shared data as well as functionalities for query processing rather than simple searches through keywords. PDMS are considered the result of blending the benefits of P2P networks, e.g. lack of a centralized authority, with the richer semantics of databases.

PDMS do not consider a single global schema [Valduriez and Pacitti 2004]. Instead, each peer represents an autonomous data source and exports its data schema. Such schema, named exported schema, represents the data to be shared with the other peers of the system. Among those exported schemas, semantic mappings are generated. Thus, query processing is accomplished by traversing such semantic mappings, rewriting the queries, executing them on the peers and gathering the results at the peer that requested data.

The main data management issues that a PDMS must deal with are the identification of schema mappings and the query processing [Heese et al. 2005]. These issues can be facilitated if an efficient strategy for peer connectivity is employed. One of such strategies consists in grouping semantically similar peers, i.e., peers sharing similar data, within semantic communities. A semantic community is “*a set of peers with common interests about a specific topic which are organized according to a particular topology*” [Castano and Montanelli 2005]. Each community treats of a specific interest, for example education or health. An interest can be formalized through keywords or ontologies, and should be generic enough to include relevant peers. The usage of semantic communities in a PDMS can improve the generation of semantic mappings, enhancing query results. Also, queries can be addressed only to relevant peers, minimizing network traffic and increasing system scalability.

Several issues should be considered when employing semantic communities in a PDMS: shared data representation format, clustering policies, semantic grouping level, and mostly, the discovery, formation and maintenance of semantic communities. In this sense, the main goal of this research is to propose an approach for semantic-based peer connectivity in a PDMS. Such PDMS provides a P2P infrastructure that facilitates identification of semantic mappings and query processing. The data shared by peers are represented through ontologies which are employed for grouping semantically similar peers within semantic communities and clusters. Due to the dynamic behavior of peers, approaches for load balancing of semantic clusters and fault tolerance are also proposed.

The remainder of this paper is organized as follows. Section 2 presents a detailed description of the proposed PDMS architecture. Section 3 details our strategy for semantic-based peer connectivity, load balancing, and fault tolerance in the PDMS. Section 4 outlines the main contributions of the thesis. Section 5 presents the current stage of the work. Related work is discussed in Section 6. Finally, Section 7 presents the concluding remarks.

## 2. System Architecture

In this section, we present a PDMS, named SPEED (Semantic PEEr-to-Peer Data Management System), in which peers are clustered according to their shared data [Pires et al. 2006, Pires 2007]. The system employs a mixed network topology, particularly DHT [Stoica et al. 2001] and super-peer [Yang and Garcia-Molina 2003], in order to exploit the strengths

of both topologies. A DHT network is used to assist peers with common interests to find each other and form semantic communities. Within a community, peers are arranged according to a super-peer topology. The combination of distinguishing network topologies facilitates the identification of semantic mappings and improves query processing.

## 2.1. Architecture Overview

As shown in Figure 1, three distinct types of peers are considered in SPEED: data peers, integration peers and semantic peers. A **data peer** represents a data source sharing structured or semi-structured data with other data peers in the system. In Figure 1,  $I_1D_1$  and  $I_1D_2$  are examples of data peers. Data peers are grouped within **semantic clusters** according to their respective exported schema. The exported schema corresponds to an ontology description of the data shared by the data peer, and is named **local peer ontology**.

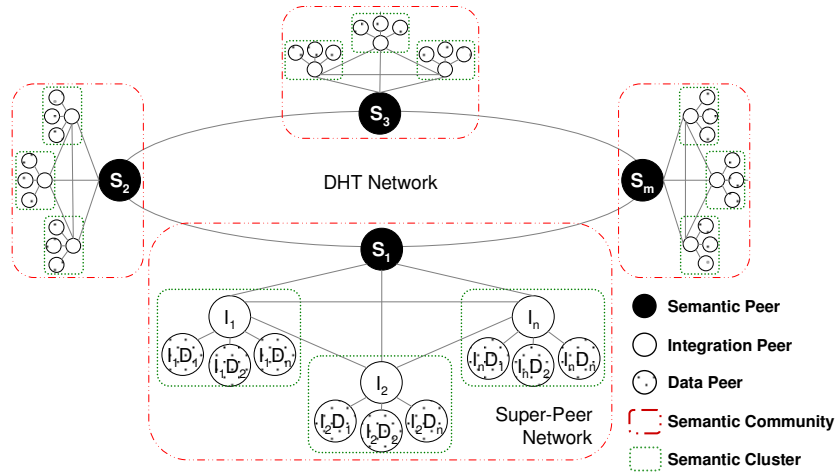


Figure 1. Overview of the SPEED's architecture

Each semantic cluster has a special type of peer with high computational capacity, named **integration peer**. In fact, integration peers are data peers with high availability, network bandwidth, processing power and storage capacity. Such peers are responsible for tasks like management of data peers' metadata, query processing and data integration. In Figure 1,  $I_1$  is the integration peer of the semantic cluster formed by the data peers  $I_1D_1$ ,  $I_1D_2$ , and  $I_1D_n$ . An integration peer maintains a **cluster ontology**, which is obtained through the merging of the local ontologies representing data peers' and integration peer's exported schemas. It acts as shared vocabulary inside a semantic cluster, inter-relating semantically similar ontology concepts. Integration peers communicate with a **semantic peer**, which is responsible for storing and offering a **community ontology** containing concepts and properties of a particular knowledge domain. Also, semantic peers are responsible for managing integration peers' metadata. A set of clusters sharing semantically similar interests forms a **semantic community**. In this sense, a data peer may participate in one or more semantic clusters within the same semantic community. In Figure 1,  $S_1$  is an example of a semantic peer.

In a PDMS, peer connectivity is considered dynamic and ad-hoc. In SPEED, the connection of requesting peers<sup>1</sup> starts through the DHT network to facilitate resource discovery by assisting them to efficiently find other related peers and form semantic communities. A requesting peer is initially connected as a data peer. As DHT networks are characterized by efficient searches and sensibility to changes in their structure, SPEED's DHT

<sup>1</sup> A requesting peer is a peer wishing to connect to the system.

network is composed only by semantic peers, i.e., peers with high reliability, network bandwidth, and availability. Excluding dynamic peers from the DHT network avoids unnecessary maintenance costs. In addition, it helps to forward requesting peers to adequate communities which are more likely to be achieved with a smaller number of hops. In this sense, the semantic associated to the content shared by peers is a crucial aspect for the formation of semantic communities.

Within semantic communities, peers are organized according to the super-peer topology. Clustering peers according to their semantic interest provides an environment that is better suited to ontology matching techniques [Shvaiko and Euzenat 2005]. If we consider the use of semantic clusters, then semantic mappings are established between semantically similar peers. Besides, as each integration peer maintains an index of its attached data peers, query routing can be efficiently carried out. Furthermore, the physical heterogeneity of participating peers is also exploited.

### 3. Semantic-based Peer Connectivity in SPEED

In SPEED, the semantic domain dictates the nature of peer connectivity. A semantic-based peer connectivity strategy is utilized to associate requesting peers to adequate communities and clusters. In addition, load-balancing and fault tolerance approaches are available in such a way that, in the event of an integration peer disconnection, for example, a new integration peer is chosen and data peers are efficiently redistributed among other clusters. In the following subsections, we explain these approaches in more details.

#### 3.1. Semantic Community Discovery

Basically, the discovery of a semantic community is performed in a two-fold way: search through keywords and ontology comparison. Since the latter seems to be an expensive process in P2P networks, the former intends to minimize the number of ontology comparisons by discarding irrelevant semantic peers. Initially, the requesting peer sends a description of its knowledge domain to an arbitrary semantic peer within the DHT network. Such description is represented by a set of keywords extracted from its local ontology. Semantic peers are searched according to a particular DHT protocol (for example, Chord). Afterwards, the semantic matchmaker module located in each semantic peer found performs an ontology comparison between the local ontology (requesting peer) and the community ontology (Figure 2). The comparison consists in identifying a semantic similarity degree among concepts stored in both ontologies. A numeric measurement of the similarity between both ontologies is generated [Wang and Ali 2005].

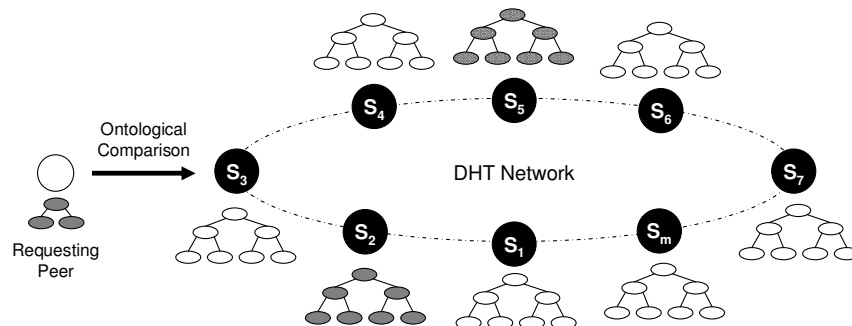


Figure 2. Semantic community discovery through ontology comparison between requesting peer and semantic peers (S<sub>2</sub> and S<sub>5</sub>)

The requesting peer is relevant for the community if the semantic matchmaker produces a value higher than a *community threshold*. A community threshold specifies the minimum semantic similarity value required to consider the local ontology and the

community ontology as similar ontologies. In other words, if the comparison result is greater than the community threshold, then the requesting peer will make part of that community. Otherwise, the requesting will be addressed to another semantic peer according to a semantic-based DHT protocol [Sangpachatanaruk and Znati 2004]. In the current version of SPEED, a peer can participate in only one community. The community threshold is defined by a system administrator and can vary from community to community.

### 3.2. Semantic Cluster Discovery and Formation

Once the semantic community has been discovered, the requesting peer should find out appropriate semantic clusters within the super-peer network. Similarly to the community discovery, cluster discovery is also performed through an ontology comparison (Figure 3). However, at this moment the semantic matchmaker module which is located at the integration peer performs an ontology comparison between the requesting peer's local ontology and the cluster ontology, producing a similarity degree between both ontologies.

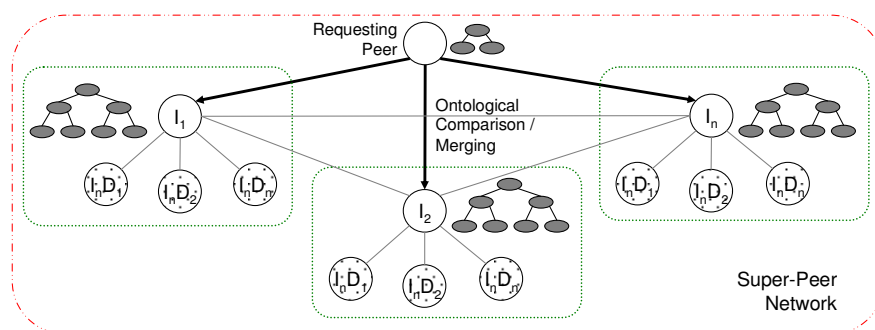


Figure 3. Cluster discovery (ontology comparison) and cluster formation (ontology merging)

The requesting peer is relevant for the cluster if the semantic matchmaker produces a value higher than a *cluster threshold*. A cluster threshold specifies the minimum semantic similarity value required to consider the local ontology and the cluster ontology as matching ontologies. That is, if the comparison result is higher than the cluster threshold, then the requesting peer will make part of that cluster. Otherwise, the requesting will be addressed to next integration peer within the same semantic community, and so on.

If the comparison result is lower than the cluster threshold for all clusters, then the requesting peer will form a new cluster. Exceptionally, the first peer in a cluster is connected as an integration peer. In this sense, the cluster ontology is created when the first data peer is connected to the cluster. As long as other data peers join that cluster, the integration peer can extend the cluster ontology by adding new concepts through ontology merging [Bruijn et al. 2004]. Similarly, the cluster ontology can be expanded by enriching existing concept descriptions in terms of new attributes and of new relationships acquired by other data peer's local ontology. For ontology merging techniques we rely on existing solutions, by adapting them to the problem of ontology merging in P2P environments.

### 3.3. Load Balancing of Semantic Clusters

The dynamic behavior of data peers and integration peers can lead to situations where a cluster may have more data peers than another. Therefore, it is necessary to redistribute data peers among other clusters within the same community. An adequate **cluster average size** is of great importance in the sense that, if most of the connected peers are integration peers, the system is more like a pure P2P system and several peers will participate in query processing. On the other hand, if too few integration peers are available, the system is more like a centralized system. The optimal value for cluster average size depends on the system or

application to be developed. A threshold can be used to determine lower and upper limits for a cluster size. The lower limit is defined as *cluster average size – threshold*, whereas upper limit is *cluster average size + threshold*. In this work, we assume that both parameters are defined by a system administrator. Therefore, considering a cluster average size and a threshold of 10 and 5, respectively, then a cluster is considered balanced if its actual size is between 5 (lower limit) and 15 (upper limit).

In our approach event-condition-action (ECA) rules are used to maintain cluster average size and, if possible, preserve the semantic distribution of data peers. Depending on events and conditions, one or more actions can be triggered. Three different types of events are considered: (i) a data peer requests a connection to a cluster; (ii) a data peer disconnects from a cluster; and (iii) an integration peer disconnects from a cluster. The condition concerns in verifying whether cluster actual size  $\pm 1$  is between lower and upper limits ( $\text{lower limit} \leq \text{cluster actual size} \pm 1 \leq \text{upper limit}$ ). Several types of actions are taken into account, for example, a data peer is chosen for integration peer, a cluster is split into two clusters, and two or more clusters are merged forming a new single cluster.

### 3.4. Integration Peer Replacement

As discussed in Section 2.1, each semantic cluster must have an integration peer. When an integration peer fails or disconnects, a fault tolerance approach must be available in order that the other data peers keep connected. Therefore, we employ a pro-active solution in which one of the data peers of a particular cluster is previously elected as a candidate integration peer. It acts as a redundant integration peer and keeps a copy of the actual integration peer’s knowledge base. The knowledge base is periodically replicated from the actual integration peer to the candidate integration peer. If the actual integration peer fails, then the candidate integration peer assumes its role and another data peer is chosen as candidate integration peer.

Since integration peers are responsible for executing important issues within a cluster, several characteristics need to be considered so a data peer can become an integration peer candidate. Such characteristics include physical resources available such as physical memory, disk space, CPU powerfulness, and network bandwidth. Additionally, the behavior of a data peer, while it is connected to the system, should be an essential factor when determining an integration peer candidate. Thus, subjective characteristics are also taken into account, for example, availability, accuracy, response time, completeness, and amount of data.

In order to measure the eligibility of candidate peers we use the capacity metric proposed in [Zhuang et al. 2004]. The function  $capacity(p)$  refers to the capacity value of a data peer  $p$ . The value of capacity is computed through the following formula, where  $v_i(p)$  is the value of the  $i_{th}$  metric for data peer  $p$  and  $w_i(d)$  is the weight of the respective metric. The weights are adjusted by a system administrator.

$$capacity(p) = \sum_{i=1}^n w_i * v_i(p)$$

## 4. Main Contributions

In this work, we are mainly concerned with peer connectivity, load balancing and fault tolerance issues in a PDMS. Although related, other data management issues such as the identification of schema mappings and the query processing do not make part of the scope of this work and are being develop in parallel [Souza 2007].

The primary contribution of this thesis is an approach for semantic-based peer connectivity in a PDMS. In addition, we have envisioned more specific contributions: (i) specification and formalization of the SPEED architecture; (ii) analysis of semantic-based

peer clustering policies; (iii) definition of a semantic-based protocol for peer connectivity in a PDMS, including techniques for manipulating ontologies as a semantic representation of peers' shared content; (iv) definition of a load balancing approach to efficiently redistribute peers among semantic clusters; and (v) definition of a fault tolerance approach to replace integration peers in case of failure or disconnection.

## 5. Current Stage of Work

At the present moment, the SPEED architecture has been specified and is currently being formalized. Since SPEED is based on an existing data integration system, named Integra [Lóscio 2003], we are currently analyzing some of the Integra's components, especially those concerned with query processing, in order to identify the modules that can be reused in SPEED. In addition, to enable network and lower-level services as well as to realize tests, we are investigating network platforms such as JXTA. Semantic-based protocols for P2P networks and techniques for ontology manipulation in P2P systems are also being studied.

A prototype is being developed to simulate peer connectivity and analyze the formation of semantic communities and clusters. In the literature, we have found only P2P simulators for file sharing among peers. Therefore, we are presently adapting a P2P simulator [PlanetSim 2007] to simulate data sharing. The main issues to be analyzed in the prototype are:

- *DHT network*: test the efficiency of the semantic protocol defined for the DHT network. In fact, we are interested in finding out answers to questions such as: (i) *Are requesting peers being connected to a correct semantic community?* (ii) *How many semantic peers are contacted on average before a requesting peer discovers an appropriate community?*
- *Performance and scalability issues*: ontology manipulation in P2P systems seems to be an expensive solution. For example, due to the dynamic and ad-hoc connectivity of peers, it may be impracticable to constantly update the cluster ontology. Test results could reveal the most appropriate events that could trigger an update in the cluster ontology.

## 6. Related Work

Recently, P2P systems have exploited the semantic properties of peers' shared content to cluster semantically similar peers and, consequently, improve querying and searching performance. [Li and Vuong 2005] propose an ontology-based community routing architecture to optimize search in P2P file sharing systems. Such architecture integrates different types of network topologies: an upper-level DHT-based category network and multiple lower-level decentralized unstructured community networks. Within each community peers are organized in an unstructured pure topology. In each community only one peer participates in the DHT network. Similarly, [Comito et al. 2006] propose PARIS, a PDMS whose goal is to develop a decentralized network of semantically related schemas that enables the formulation of queries over autonomous, heterogeneous, and distributed data sources. PARIS also uses a hybrid topology that mixes pure unstructured and DHT network topologies. However, in each community several peers participate in the DHT network.

In the literature, several load balancing approaches have been proposed to adapt the size of clusters in super-peer systems. [Zhuang et al. 2004] propose a dynamic layer management algorithm, which can adaptively elect peers and maintain a given size ratio of super-layer to leaf-layer. Such algorithm ignores the semantic associated to peers' shared content. In the solution proposed by [Brito and Moura 2005], when a cluster is unbalanced, peers are randomly redistributed among any other unbalanced clusters.

Differently from the previous systems, SPEED utilizes DHT and super-peer topologies, since we believe that pure unstructured networks suffer from serious scalability

problems as the number of peers in network increases. Additionally, in SPEED we avoid unnecessary maintenance costs by excluding dynamic peers from the DHT network. In our connectivity approach peers are grouped in a finer granularity level (cluster level) to improve the identification of semantic mappings and query processing. Regarding load balancing, we utilize a semantic-based approach to redistribute data peers within semantic clusters and, if possible, preserve the semantic distribution of data peers.

## 7. Concluding Remarks

In this work, we have presented an approach for semantic-based peer connectivity in a PDMS. Such approach intends to facilitate important peer data management issues in PDMS, such as identification of semantic mappings and query processing on a large number of data sources. In this sense, a PDMS, named SPEED, has been proposed which utilizes a mixed network topology. In SPEED, peers are grouped within semantic communities and clusters according to their shared data, represented through ontologies. Due to the dynamic behavior of peers, approaches for load balancing of semantic clusters and fault tolerance are also proposed.

## References

- Brito, G., Moura, A. M. (2005) "ROSA-P2P: a Peer-to-Peer System for Learning Objects Integration on the Web". In Proc. of the 11<sup>th</sup> WebMedia'05, Poços de Caldas, Minas Gerais, Brazil, 1-9.
- Bruijn, J., Martin-Recuerda, F., Manov, D., Ehrig, M. (2004) "State-of-the-art-survey on Ontology Merging and Aligning". v1. SEKT Project deliverable D4.2.1.
- Castano, S., Montanelli, S. (2005) "Semantic Self-Formation of Communities of Peers". In Proc. of the ESWC Workshop on Ontologies in Peer-to-Peer Communities, Heraklion, Greece.
- Comito, C., Patarin, S., Talia, D. (2006) "A Semantic Overlay Network for P2P Schema-Based Data Integration". In Proc. of the 11<sup>th</sup> IEEE Symposium on Computers and Communications (ISCC'06), 88-94.
- Heese R, Herschel S, Naumann F, Roth A. (2005) "Self-Extending Peer Data Management". In GI-Fachtagung für Datenbanksysteme in Business, Technologie und Web (BTW'05), Karlsruhe, Germany.
- Li, J., Vuong, S. (2005) "Ontology-Based Clustering and Routing in Peer-to-Peer Networks". In Proc. of the 6<sup>th</sup> Int. Conf. on Parallel and Distributed Computing, Applications and Technologies, Dalian, China.
- Lóscio, B. (2003) "Managing the Evolution of XML-based Mediation Queries". Ph.D. Thesis. UFPE, Brazil.
- Pires, C. E. S. (2007) "Um Sistema P2P de Gerenciamento de Dados com Conectividade Baseada em Semântica". Ph.D. Thesis Proposal. UFPE, Brazil.
- Pires, C. E. S., Lóscio, B. F., Salgado, A. C. (2006) "Data Management in P2P Systems". In Proc. of the 21<sup>st</sup> Brazilian Symposium on Databases (SBBD'06), Florianópolis, Brazil. pp. 310.
- PlanetSim P2P Simulator (2007), <http://planet.urv.es/trac/planetsim/wiki/PlanetSim>, last date of access May 25, 2007.
- Sangpachatanaruk, C. and Znati, T. (2004) "Semantic Driven Hashing (SDH): An Ontology-based Search Scheme for the Semantic Aware Network (SA Net)". In Proc. of the 4<sup>th</sup> Int. Conf. on Peer-to-Peer Computing, Zürich, Switzerland.
- Shvaiko, P., Euzenat, J. (2005) "A Survey of Schema-Based Matching Approaches". Journal on Data Semantics IV: 146-171.
- Souza, D. (2007) "Semantic-based Query Reformulation in PDMS". Ph.D. Thesis Proposal. UFPE, Brazil.
- Stoica, I., Morris, R., Karger, D., Kaashoek, M. F., and Balakrishnan, H. (2001) "Chord: a Scalable Peer-to-Peer Lookup Service for Internet Applications". ACM SIGCOMM'01, San Diego, USA. pp. 149-160.
- Tatarinov, I., Ives, Z., Madhavan, J., Halevy, A., Suci, D., Dalvi, N., Dong, X., Kadiyska, Y., Miklau, G., Mork, P. (2003) "The Piazza Peer Data Management Project". In Proc. of the ACM SIGMOD Record, 32(3):47-52.
- Valduriez, P., Pacitti, E. (2004) "Data Management in Large-Scale P2P Systems". In Proc. of Int. Conf. on High Performance Computing for Computational Science (VecPar'04), Valencia, Spain.
- Wang, J. Z., Ali, F. (2005) "An Efficient Ontology Comparison Tool for Semantic Web Applications". In Proc. of the IEEE/WIC/ACM International Conference on Web Intelligence, France. pp. 372-378.
- Yang, B. and Garcia-Molina, H. (2003) "Designing a Super-Peer Network". In Proc. of International Conference on Data Engineering (ICDE'03), Bangalore, India.
- Zhuang, Z., Liu, Y., and Xiao, L. (2004) "Dynamic Layer Management in Super-Peer Architectures". In Proc. of the International Conference on Parallel Processing (ICPP'04) - Volume 00, pp. 29-36.