

Query Planning in P2P Database Systems

Kjetil Nørnvåg,* Eirik Eide and Odin Hole Standal
Dept. of Computer Science,
Norwegian University of Science and Technology,
Trondheim, Norway

Abstract

The peer-to-peer (P2P) paradigm is emerging as a possible solution to some of the problems in distributed data processing, including scalability, availability, and administrative cost. P2P has already proved to be suitable in contexts like file sharing, distributed computations, and distributed search. In our research we are aiming at using P2P to solve some problems in the domain of distributed databases. In this paper we 1) present PORDaS, a distributed DBMS based on P2P techniques, 2) describe query processing and query planning in PORDaS, and 3) present results from an experimental evaluation of different query planning variants.

1. Introduction

The peer-to-peer (P2P) paradigm is emerging as a possible solution to some of the problems in distributed data processing, including scalability, availability, and administrative cost. P2P has already proved to be suitable and efficient for file sharing, distributed computations, and distributed search.

In our research we are aiming at using P2P to solve some problems in the domain of *distributed databases*, and in particular in the application area of database support for Grid applications. So far, Grid computing has gained some maturity with respect to the actual computation. However, the management of data in Grid networks is still a very immature area. In general, simple files are used. The need for using databases to a larger extent has been identified and there has also been some work on standardized data access services like OGSA-DAI [14].

The goal behind our research is to provide database facilities where distribution (and the availability of the Grid backbone) is transparent to the user. It should also provide

services for metadata discovery and seamless queries between heterogeneous sources. In this paper we will give an overview of PORDaS, which is the distributed database system layer of the DASCOSA Grid database framework [11]. We will also present some details from query processing and planning in the PORDaS prototype.

Many of the architectural decisions of PORDaS are affected by characteristics of the intended application areas: relatively complex and structured data, and the fact that much of the data will be local data that should be made available to the outside world for querying, but for various reasons (including the size of the data volumes) the raw/source data itself should not be distributed. For this reason our system-wide data model is based on the traditional object-relational data model in order to present users/application the same data model as what is used in their applications.

The main contributions of this paper are: 1) a presentation of the PORDaS P2P DBMS, 2) query processing in a P2P DBMS, and 3) query planning in a P2P DBMS.

The organization of the rest of this paper is as follows. In Section 2 we give an overview of related work. In Section 3 we give an overview of PORDaS, and in Section 4 we give a more detailed description of query processing. In Section 5 we present results from an experimental evaluation of different query planning variants. Finally, in Section 6, we conclude the paper.

2. Related work

Much of the previous work on distributed database systems is obviously relevant. For a survey of state of the art in this area we refer to [8]. Recent work in this area includes query processor for Internet data sources, for example ObjectGlobe [4].

Our Grid DBMS PORDaS is based on distributed hash tables (DHT). A number of papers deal with focused issues like query processing in DHT networks [2, 6, 16], and replica management [9]. [13] describes how to use Dis-

* Email of contact author: Kjetil.Norvag@idi.ntnu.no.

tributed Hash Sketches to estimate cardinality of multisets in the context of P2P system based on DHTs.

Three systems, PIER, AmbientDB and PeerDB also aim at providing DBMS support using P2P technology:

- *PIER* [7] has many similarities with PORDaS. It is a middleware query engine built on top of a storage manager and DHT. However, it is not designed to support replication and does not maintain system metadata, and essentially only indexes whatever the applications register in the system.
- *AmbientDB* [3] is a system designed to provide full relational database functionality for stand-alone operation in autonomous devices that may be mobile and disconnected for long periods of time, while enabling them to cooperate in an ad-hoc way with (many) other AmbientDB devices. A DHT is used both as a means for connection peers in a resilient ways as well as supporting indexing of data.
- *PeerDB* [10] is a P2P system supporting queries against data stored on remote nodes. The system is based on an unstructured P2P system, focusing on data retrieval instead of distributed querying. Instead of relying on global schemas or mediators, information retrieval techniques are used to find matching relations. Both relation matching and queries are performed by agents.

A general problem in P2P systems is selfish behavior: most peers want to receive more than they contribute. In order to reduce the impact of this behavior, techniques using accounting [12] and other approaches to enable trust have been developed.

3. PORDaS

PORDaS is a distributed DBMS using P2P techniques to achieve high performance, scalability, and availability. It is built as a database layer on which larger applications can be built, and provides location-transparent storage of data with Grid applications as main application area. Each node in PORDaS is autonomous. Data is created and stored locally but globally available for querying.

In both Grids and other large distributed systems heterogeneous hardware and operating systems will be the case. In order to ease deployment, we have based PORDaS on components written in Java. In the current version, the P2P communication is performed by the FreePastry DHT [5], and for local storage the Derby DBMS [1] is used. It should be mentioned that both the DHT and DBMS are defined as interfaces so that they can easily be replaced with other implementations if desired. They are also both embedded, so that no separate installation of DHT and DBMS is necessary.

In a distributed DBMS, metadata management is an important issue. In PORDaS this is solved by separate table descriptions and data discovery tools, so that queries over related data sources using heterogeneous schema descriptions can be performed. This schema management is orthogonal to the work presented on query processing in this paper, so due to space constraints we will not go into further details on this issue.

In the rest of this section the overall architecture of PORDaS, while storage and query processing will be described in more detail in the subsequent sections.

3.1. Use and data access

PORDaS uses the object-relational data model, i.e., relations of tuples and possibility of tuple identifiers and attributes referencing other tuples.

When a table is created, a hash value based on the table schema is created. This is indexed in the DHT, so that it is possible to find other tables having the same schema signature. It is also possible to annotate a schema by keyword descriptions that can be taxonomy/ontology based. This information is also stored in the DHT. This means that it is possible to find related tables automatically.

Data that is inserted is stored in the local database. A query against local table names will be performed on the local database only. If it is desired that the global database should be queried, i.e., potentially all other local databases in the system, global tables have to be specified in the query. The identifier of global tables are found by using either schema signature or schema description as described above, or provided directly.

Given the identifier of a global table, the DHT can return the identifier of all peers storing elements of this table. It can also indicate value ranges that they store for the key value, which might speed up, e.g., selection queries. Although indexing individual tuples is an alternative using a DHT, this will in general be too fine-grained and having too high maintenance cost.

3.2. Overview of architecture

The architecture of PORDaS is illustrated in Figure 1. An application accesses the databases through the PORDaS API. The application layer can interact both with the local database and query data at other sites transparently, without users having to know where the data is located.

Storage: The storage component keeps track of all the data stored at a node. It holds the local database and a metadata repository that manages information about local tables. The storage component also stores parts of the distributed index, which all the nodes in the system participate in sharing. The

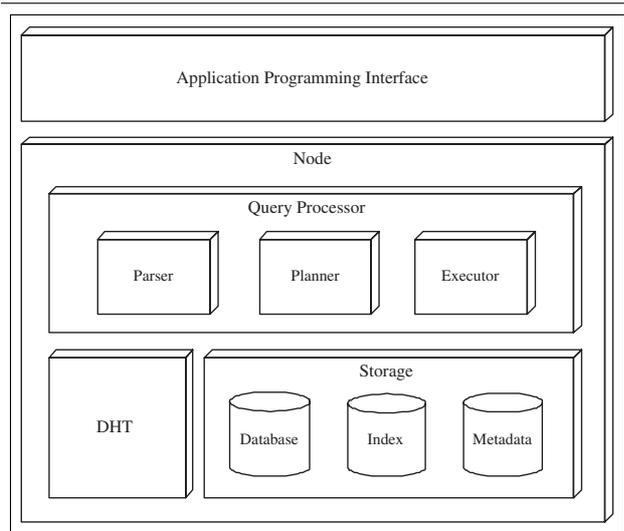


Figure 1. The PORDaS architecture.

index holds information about the contents and location of other tables in the system.

The local database is used for the persistent data stored in the system, i.e., tables and local metadata. PORDaS operates internally on an object relational data model. This is also the model exposed to users/applications of PORDaS. A query could be performed either against local data only, or against global data. In the case of global data, location is transparent.

Query processor: The query processor consists of a parser, a planner and an executor. It takes queries as input from the application layer, creates a plan and executes it. The resulting tuples are pipelined back to the application layer.

Communication: The nodes in a PORDaS network are loosely connected through a DHT. The purpose of the communication module is to enable resource location and to route requests for data. It also enables message sending and reception, and is responsible for unwrapping messages and forwarding them to the appropriate component.

Except when returning query results, every message in PORDaS is routed through the DHT layer. Messages are sent for keyword and query requests, to resolve subqueries and to maintain the distributed index. When returning results, direct connections between the nodes are used in order to improve performance.

The distributed index is realized using a DHT, and among its applications in PORDaS are 1) schema/table discovery, and 2) finding locations of arbitrary tables. The index makes PORDaS location transparent.

Soft state is used to maintain the distributed index. It means that every index record has an associated expiration

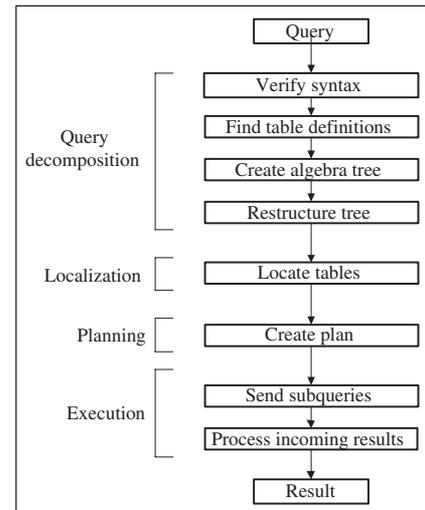


Figure 2. Query processing.

time. When it expires, the record is deleted from the index. Thus, in order for a node to keep its index records in the system, they must be continually refreshed. After a node leaves the system, either voluntarily or because of a failure, its index records will be removed when they expire.

4. Query Processing in PORDaS

Figure 2 shows how queries are resolved in PORDaS, with a structure mostly the same as for traditional distributed systems. The query processor can handle multiple queries at a time. We now give a brief overview of each query processing step, followed by a more detailed description of the query planning step.

Query language: The query and data manipulation language in the PORDaS prototype is a subset of SQL.

Query decomposition: In this step queries are decomposed from a textual representation into an algebra tree.

When a new query is submitted, first the query syntax is verified. Then the next task is to verify type correctness, which is making sure the relations and attributes referenced in the query actually exist. The operations in the query are checked against the type of each attribute as well, making sure they match. Before this can be done, the table definition for each table referenced in the query must be fetched. First, the local metastore is searched. If one or more table definitions are lacking, they must be fetched from the distributed index. This will delay the query until all table definitions have been found.

The next step is creating the initial algebra tree. The final tree is either bushy or linear. The first part of the tree is

always the same for both types. First, the projections are defined as the root, followed by each selection in the query.

PORDaS builds left-deep trees, i.e., the tree is always extended along the leftmost path of each operator, with only base relations as the right child of the operators. Cartesian products are added as the left child of its parent. If there are joins in the query, these are added as the last part of the tree, in the same fashion as cartesian products.

When building a bushy tree, the goal is to balance the tree as much as possible, catering for parallel execution. This is done by maximizing the number of cartesian products and joins with two operators as children.

Restructuring the tree is performed to achieve a better tree. The projections and selections in the algebra tree are pushed as far down as they can be, and in doing so, the size of intermediate results will be reduced during execution.

Localization: Localization means finding every site that has a table referenced in the query. These are found by querying the distributed index. Because of the potentially volatile nature of the P2P network, caching is not used to improve the localization process. Nodes that store one of the tables in the query might have left or joined the network since the last time they were queried.

Planning: When the locations of every table in the query are found, a plan for the execution can be devised. The planner uses heuristics to create plans; these plans are either *centralized* or *distributed*.

A *centralized plan* is a plan where the required base relations are fetched to the initiating site so that the operators in the query can be resolved locally. In the spirit of reducing the amount of network traffic, any selection and projection on base relations are executed at the remote sites before the streaming of results is started. Figure 3 (top) gives an example of a centralized plan. The dotted lines indicate network communication. It is important to note that the data fetched from a base relation in the figure may include contacting one or more sites.

In a *distributed plan*, the responsibility for executing the query tree is distributed among the set of nodes that store tables referenced in the query. If statistics about each table was available, like cardinality, maximum and minimum values, this information could be used to restructure the tree in a beneficial way. As this is not the case, predefined rules are used instead. The rule is best explained by an illustration, see Figure 3 (bottom). It gives an example of a conversion from an algebra tree to a distributed plan. The first rule is to always calculate cartesian products at the initiating node, which avoids sending too much data over the network. It is conceivable that in certain cases it might be better to distribute cartesian products as well, but it is assumed to not be the average case.

The second rule is to delegate the resolution of joins to one of the owners of the leftmost table in the join tree. The

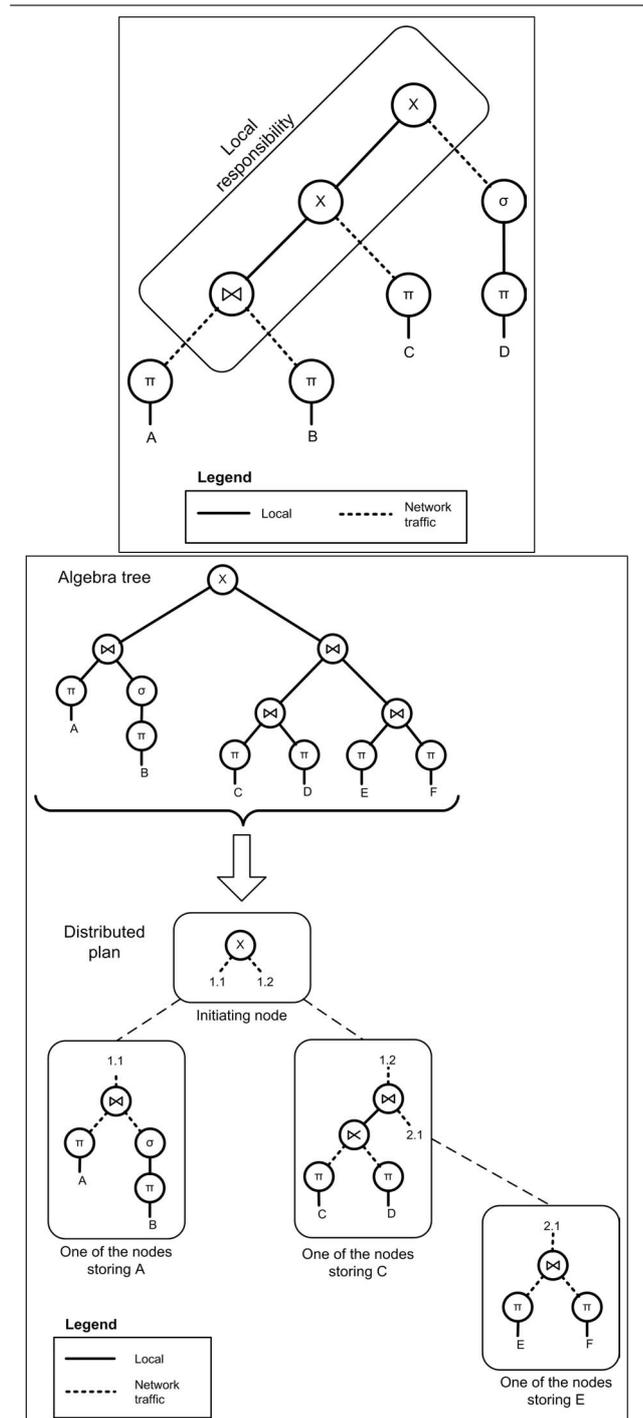


Figure 3. Centralized plan (top) and transformation from an algebra tree to a distributed plan (bottom).

choice is arbitrary, it could have been any of the owners. The chosen owner will receive the entire join subtree. If this subtree has any more joins in it, these joins are delegated in the same manner as well. In the figure this can be seen as the second join under the cartesian product is delegated to two separate nodes.

Execution: The execution phase begins with the initiating node requesting data from all base relations in the query. If the query plan is a distributed query plan subtrees are delegated to other nodes. The execution is pipelined, which means that processed tuples are sent up the tree as soon as possible. This avoids having to store temporary caches with intermediate results. To know where a tuple belongs at the receiving end, all tuples sent across the network are tagged with an identifier. The identifier uniquely specifies the query and the point in the query tree where the tuple is expected. This can be seen in the plan in Figure 3. For instance, the children of the cartesian product are tagged as 1.1 and 1.2. At the sites responsible for those parts of the query, every resulting tuple is tagged with 1.1 or 1.2, respectively.

5. Experimental results

In this section study the different query planning variants. The experiments were conducted on a cluster of 36 computers, each having a 3 GHz Pentium 4 and 1 GB RAM. The experiments were performed by a test application that simulated the activity of a regular PORDaS node. The parameters of the test application are summarized in Table 1, and are chosen to simulate a probable usage pattern under medium load. The time between requests is Poisson distributed with a given mean value. The queries will be delayed if there are more concurrent queries than the maximum number of concurrent queries allowed. The actual databases have been kept small in order to be sure PORDaS is tested and not the local database system (i.e., Derby).

Due to space constraint we limit the discussion in this paper to two of the experiments that were conducted. In the *first experiment*, all nodes were sharing and active. The purpose was to compare every permutation of type of algebra tree and type of planner, which means that 4 simulations were run. The interesting data in this case are the number of started queries versus the number of finished queries, and the response times for each permutation. Figure 4 (left) shows a comparison of all permutations of type of query planner and type of algebra tree with respect to the number of started and finished queries. In both cases using centralized plans the loss of queries is minimal, while in the distributed case, the loss is noticeable. With a linear query tree, the loss is over 13%. The circumstances indicate that the reason for the loss is that queries timed out before results were received. The time-out threshold was set to two minutes. By looking at the maximum response times for

the queries that did finish, these are close to this limit. Figure 4 (left) also shows a distinct difference between the centralized and distributed cases. The number of started queries is much lower in the distributed case. The reason for not starting more queries is because of the limit on the number of concurrent queries. Waiting for a query to be executed has the effect of lowering the throughput.

The purpose of the *second experiment* was to study and isolate the effect of executing queries in parallel. This was achieved by having one active, non-sharing node query the rest of the nodes, which were sharing and inactive. The only task of the non-active sites in the system was to resolve the queries the one active node gave them. Two tests were run. The first test had the query processor make linear trees and used a centralized execution strategy. The second test used bushy trees and a distributed execution strategy. Figure 4 (right) shows a comparison of the minimum, average and maximum response times observed for queries with a result size between 25.000 and 50.000 tuples. The distributed execution strategy performs consistently worse than the centralized strategy.

Both experiments reported above show that the centralized execution strategy is better than the distributed. The distributed execution strategy has the advantage of executing operators in parallel, but failed since the joins always had a huge selection rate. This caused the distributed execution strategy to generate a lot more network traffic than the centralized strategy. The distributed strategy could have worked better in comparison to the centralized strategy if the joins had lower selection rates.

6. Conclusions and further work

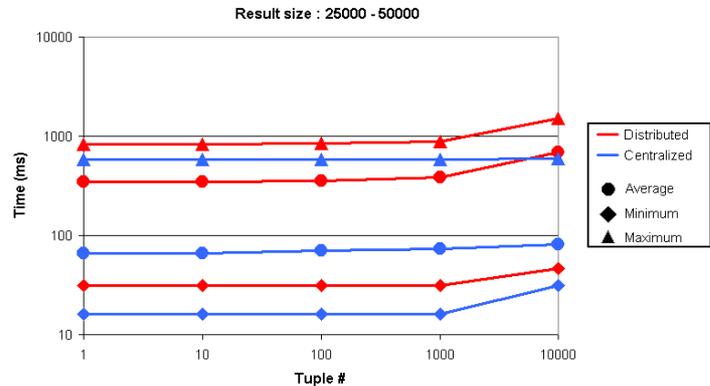
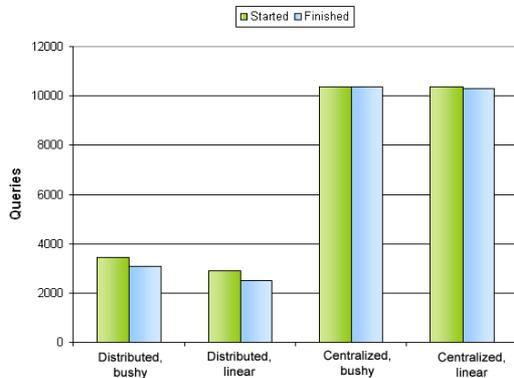
In this paper, we have given an overview of the P2P DBMS PORDaS, and described some aspects of query processing and query planning in this system. We have also presented some results from an experimental evaluation of different query planning variants in PORDaS.

Future work will be in two directions: 1) large-scale experiments and 2) extended functionality. Experiments are planned to be performed both on larger clusters and Grids in Norway, as well as using PlanetLab [15]. Extended functionality that will be implemented includes replication, enhanced indexing, and a further development of aspects of query optimization in the context of P2P DBMSs.

References

- [1] Apache Derby, <http://db.apache.org/derby/>.
- [2] D. Bauer, P. Hurley, R. Pletka, and M. Waldvogel. Bringing efficient advanced queries to distributed hash tables. In *Proceedings of the 29th Annual IEEE International Conference on Local Computer Networks (LCN'04)*, 2004.

Parameter	Value	Parameter	Value
Tables per node	8	Tuples per node	2000
# of nodes	36	Simulation length	10 minutes
Request intensity	2 seconds	# of concurrent queries	5
Joins in query	3 (deviation 2)	# of tables in query	4 (deviation 2)
Maximum result size	100 000 (deviation 90 000)		



- [3] P. Boncz and C. Treijtel. AmbientDB: relational query processing in a P2P network. In *Proceedings of DBISP2P'2003*, 2003.
- [4] R. Braumandl, M. Keidl, A. Kemper, D. Kossmann, A. Kreuz, S. Seltzsam, and K. Stocker. ObjectGlobe: ubiquitous query processing on the Internet. *VLDB Journal*, 10(1):48–71, 2001.
- [5] FreePastry, [http://freepastry.org/FreePastry/](http://freepastry.org/).
- [6] M. Harren, J. M. Hellerstein, R. Huebsch, B. T. Loo, S. Shenker, and I. Stoica. Complex queries in DHT-based peer-to-peer networks. In *Proceedings of IPTPS 2002*, 2002.
- [7] R. Huebsch, J. M. Hellerstein, N. Lanham, B. T. Loo, S. Shenker, and I. Stoica. Querying the internet with PIER. In *Proceedings of VLDB'2003*, 2003.
- [8] D. Kossmann. The state of the art in distributed query processing. *ACM Computing Surveys*, 32(4):422–469, 2000.
- [9] P. Maniatis, M. Roussopoulos, T. Giuli, D. S. H. Rosenthal, M. Baker, and Y. Muliadi. Preserving peer replicas by rate-limited sampled voting. In *Proceedings of the 19th ACM SOSP*, 2003.
- [10] W. S. Ng, B. C. Ooi, K.-L. Tan, and A. Zhou. PeerDB: A P2P-based system for distributed data sharing. In *Proceedings of the 19th International Conference on Data Engineering*, 2003.
- [11] K. Nørnvåg. DASCOSA: database support for computational science applications. In *Proceedings of GLOBE'06*, 2006.
- [12] N. Ntarmos and P. Triantafillou. SeAI: managing accesses and data in peer-to-peer data sharing networks. In *Proceedings of HDMS*, 2004.
- [13] N. Ntarmos, P. Triantafillou, and G. Weikum. Counting at large: Efficient cardinality estimation in internet-scale data networks. In *Proceedings of the 22nd International Conference on Data Engineering (ICDE'06)*, 2006.
- [14] OGSA-DAI. Open grid services architecture data access and integration, <http://www.ogsadai.org.uk/>.
- [15] PlanetLab, <http://www.planet-lab.org/>.
- [16] R. van Renesse, K. P. Birman, and W. Vogels. Astrolabe: A robust and scalable technology for distributed system monitoring, management, and data mining. *ACM Trans. Comput. Syst.*, 21(2):164–206, 2003.