



Project no. FP6-004265

**CoreGRID**

European Research Network on Foundations, Software Infrastructures and Applications for large-scale distributed GRID and Peer-to-Peer Technologies

Network of Excellence

GRID-based Systems for solving complex problems

**D.KDM.01 – Roadmap version 1 on Knowledge and Data Management**

Due date of deliverable: February 28, 2005  
Actual submission date: 15 April 2005  
Revised version submitted on March 21, 2006

Start date of project: 1 September 2004

Duration: 48 months

**Organisation name of lead contractor for this deliverable: FORTH**

**Revision draft**

<b>Project co-funded by the European Commission within the Sixth Framework Programme (2002-2006)</b>		
<b>Dissemination Level</b>		
<b>PU</b>	Public	<b>PU</b>

**Keyword List:** Distributed Data Management, Information and Knowledge Management, Data Mining and Knowledge Discovery

**Revised version of**  
**D.KDM.01 – Roadmap version 1 on Knowledge and Data Management**  
**Responses to Reviewers**

***Recommendation 1:***

*A short overview of the planned scientific GRID activities, versus state of the art in GRID and non-GRID technology per section (Distributed Data Management, Information and Knowledge Management, Data Mining and Knowledge Discovery) should be provided in the document.*

***Action:***

To address this recommendation we added three sub-sections in the revised version of the D.KDM.01 deliverable to describe the planned scientific activities for each one of the three tasks that compose the workpackage 2. In particular, the three new sub-sections are 4.2.1, 4.2.2, 4.2.3. They respectively describe scientific plans in the areas of Distributed Storage Management, Information and Knowledge Management, and Data Mining and Knowledge Discovery.

***Recommendation 2:***

*The results of the DataMiningGrid project should be taken into account in the Data Mining and Knowledge Discovery section.*

***Action:***

We introduced the DataMiningGrid project goals in section 3.1.3 and in section 3.2.3. Then in section 4.3.2 we described the potential synergy with that project and the common research and development interests. Moreover, it is necessary to mention that, as presented in the review meeting, the leader of the DataMiningGrid project gave a lecture at the second meeting of the KDM Institute in Barcelona and we discussed there some potential future joint activities. Finally, we are involved together in the EC Technical Group 5 on Data Management.

**Table of content**

<b>1 EXECUTIVE SUMMARY .....</b>	<b>4</b>
<b>2 INTRODUCTION.....</b>	<b>5</b>
2.1 CONTEXT .....	5
2.2 PROBLEM(S).....	6
2.3 OBJECTIVES .....	6
2.4 TASKS.....	7
2.5 DRIVERS.....	7
<b>3 POSITIONING.....</b>	<b>8</b>
3.1 STATE OF THE ART (EXISTING APPROACHES) .....	8
3.1.1 <i>Distributed Data Management</i> .....	8
3.1.2 <i>Information and Knowledge Management</i> .....	9
3.2 EXTENDED CONTEXT .....	11
3.2.3 <i>Data Mining and Knowledge Discovery</i> .....	13
<b>4 VISION, STRATEGY AND ROADMAP .....</b>	<b>15</b>
4.1 VISION AND SCENARIOS .....	15
4.2 STRATEGY .....	15
4.2.1 <i>Planned Scientific Activities in Distributed Storage Management</i> .....	15
4.2.2 <i>Planned Scientific Activities in Information and Knowledge Management</i> .....	16
4.2.3 <i>Planned Scientific Activities in Data Mining and Knowledge Discovery</i> .....	17
4.3 ROADMAP .....	17
4.3.1 <i>Distributed Storage Management</i> .....	17
4.3.2 <i>Information and Knowledge Management</i> .....	22
4.3.3 <i>Data Mining and Knowledge Discovery</i> .....	25
4.3.4 <i>Phases of the roadmap</i> .....	30
4.3.5 <i>Mechanisms</i> .....	30
4.3.6 <i>Future steps</i> .....	30
<b>5 LINK WITH OTHER COREGRID SCIENTIFIC WORKPACKAGES .....</b>	<b>32</b>
<b>6 REFERENCES .....</b>	<b>33</b>
<b>7 PARTICIPANTS.....</b>	<b>37</b>

# 1 Executive Summary

This document is the first version of the roadmap of the CoreGRID Virtual Institute on Knowledge and Data Management (Workpackage 2, WP2). Knowledge and data management is a key topic in Grid computing. As mentioned in the NGG expert group report, “*The Grids environment must behave intelligently using knowledge-based techniques, including semantic-rich operation and interoperation.*” and later “*there is a need for semantically-rich knowledge-based services in both Grids Foundations Middleware and Grids Services Middleware both to improve the functionality but also to support applications on the Grids surface which require semantic and / or knowledge-based support. It is also apparent in Grids scheduling, management, durability and dynamic reconfiguration.*” From those statements, we can figure conclude that data, information, and knowledge are critical elements in the application of Grids in several sectors of our society.

The overall objective of this workpackage is to further the integration of data management and knowledge discovery with Grid technologies for providing knowledge-based Grid services commonly known as the Semantic Grid through a strong coordination of European researchers active in those areas. The workpackage will provide a collaborative setting for European research teams working on: distributed storage management on Grids; knowledge techniques and tools for supporting data intensive applications; and the integration of data and computation Grids with information and knowledge Grids. The goal is to strengthen the combined activity of research groups that today have sporadic and partial collaboration thus promoting larger leading teams and supporting efforts towards standard models and tools for data and knowledge management on Grids and P2P systems.

The Grid should be effectively exploited for deploying data-driven and knowledge-based applications. To support this class of applications, tools and services for data and knowledge management are vital. The planned joint activities of the partners involved in this workpackage will result in new research collaboration and work in the area. In the next years the Grid will be used as a platform for implementing and deploying geographically data intensive applications, distributed knowledge discovery systems and knowledge management platforms. We aim to contribute to this technological development.

## 2 Introduction

### 2.1 Context

The primary objective of the CoreGRID Network of Excellence is to build solid methodological and technological foundations for Grid computing services that will remain at the forefront of the technology and based on excellence. This will be achieved by structuring and integrating the research activities of the experts in distributed systems, middleware, programming models, algorithms, tools and environments. In particular, the objective of the WP2 is to provide a collaborative infrastructure for European research teams working on the distributed storage management of Grids, the programming techniques and tools to support data intensive applications, and the integration of data and computational Grids with information and knowledge Grids. This document describes the roadmap of the activities for Workpackage 2. We present the main research activities of the partners involved in the workpackage and how they will be integrated to perform joint research activities in the areas of data and knowledge management.

- Current work at ICS-FORTH focuses on building future scalable storage systems based on commodity components that can automatically handle resource management tasks and can offer quality-of-service guarantees. Such systems and services form the foundations of large-scale data-Grid installations distributed across geographically dispersed regions.
- The technical background of PSNC concerns storage infrastructure and results from 10 years of providing storage-related services and the involvement in storage-related R&D projects. PSNC exploits centralized storage systems (high-end SAN infrastructure, tape libraries etc.) but has also experience with building and exploiting geographically dispersed data Grids consisting of many low-end components. Both kinds of storage systems can be included in Grids and used for implementing different classes of services in Grids.
- CETIC is a R&D centre "serving the industry" - a connecting agent between academic research and companies -providing technology transfer in software engineering and electronic systems. The research areas of CETIC include software processes and products quality, requirements engineering, distributed and Grid Systems, web mining and reverse-engineering, electronic systems, free and open-source software.
- The Grid research team at UNICAL is working on two main topics: Grid-aware data mining and knowledge discovery for designing Grid services and programming tools for distributed knowledge discovery and P2P data integration models for supporting data-intensive applications on Grids. The UNICAL team has developed prototypes and research models that can be used by other research teams and can be extended through joint research work.
- CNR-ISTI has expertise in the development of state-of-the-art data mining services. In particular, the development of distributed algorithms that are adaptive w.r.t. the executing platforms and the dataset features. Such tools are available for this research community. CNR-ISTI has studied the issues concerning performance prediction in the context of the scheduling of data mining tasks execution, together with associated data transfers.
- The GRID research group at UCY has extensive experience in Grid Computing; prior work includes the development of tools for the performance evaluation of Grid resources using benchmarking and the exploration of the Grid information space using navigational tools and search engines. Additional experiences include Computer Architecture, Performance Evaluation, and Benchmarking.
- The School of Computer Science and the e-Science North West Regional Centre at the University of Manchester have strong expertise in, and an international reputation for their work on, a broad range of scientific areas especially related to data and knowledge management and the Grid. Manchester is a shaper in the UK e-Science programme and plays a leading role in the Semantic Grid, Semantic Web Services and Web ontology languages. A team of researchers also works on extending core database functionality and widens its applicability, by developing Grid-enabled

distributed query processing systems, investigating adaptive query processing, and conducting research on schedulers.

- The research areas of the Business and Information Technology Department (BITD) at CCLRC include distributed systems and Grid, knowledge management and data mining, security and trust management, and Web technologies in general. In the Grid area, BITD is currently working on the application of Grid technologies into the business domain by developing middleware to exploit new business models for ASP and Utility computing for dynamic networks of service providers as well as constructing frameworks for trust, security and contract management in virtual organizations.

## 2.2 Problem(s)

The main scientific challenges in distributed data management are the unified access of numerous heterogeneous distributed storage devices, the cost-efficient administration of the hardware and software resources involved, and the offering of quality-of-service guarantees to the system users. Here we identify some basic challenges that must be faced in the development of data and knowledge management systems and applications in Grids.

- *Data heterogeneity and large data sets.* Systems must be able to cope with very large and high dimensional data sets that are geographically distributed and stored in different types of repositories as structured data in DBMS, text in files or semi-structured data in Web sites.
- *Algorithm integration and independence.* The architectural solutions must allow the integration of different algorithms and tools and must be as independent as possible from the data sources.
- *Compatibility with Grid infrastructure and Grid awareness.* The designed services must be interfaced with the lower levels of the Grid infrastructure. The interface must be aware of the Grid services and access them for supporting applications.
- *Openness.* Solutions must be open in order to be integrated with data management tools and knowledge oriented systems.
- *Scalability.* Designed systems must be scalable both in terms of number of nodes used for performing distributed tasks and in terms of performance achieved by using large Grid configurations.
- *Security and data privacy.* Security and privacy issues are vital features in wide area distributed systems. Grid services must offer valid support to systems to cope with user authentication, security and privacy of data. Basic Grid functionality (e.g., Globus security infrastructure - GSI) must be exploited to support secure client-server interactions without impacting on the usability of the Grid infrastructure and services.

## 2.3 Objectives

Data and knowledge is going to play a more important role in current and future Grids. The issues surrounding the representation, discovery, and integration of data and knowledge in a dynamic distributed environment have to be addressed. The objective is to further the unification of data management and knowledge discovery and management with Grid technologies to provide knowledge-based Grid services for the Semantic Grid and the Knowledge Grid. Examples of this approach are techniques for managing storage resources and providing “high-quality” storage at low cost to Grid users and tools for supporting data intensive and knowledge-based applications on Grids.

Our main objectives in distributed data management are to provide commodity-based connectivity among heterogeneous distributed storage devices, management automation of administration tasks traditionally handled manually, and storage virtualization for serving well-defined requirements from multiple users. Other objectives are the development of knowledge techniques and tools for supporting data intensive applications and the integration of data and computation Grids with information and knowledge Grids. Those objectives contribute to the goal to strengthen the joint activity of research groups that today have sporadic and partial collaboration promoting larger leading teams and

supporting efforts towards standard models and tools for data and knowledge management on Grids and P2P systems.

## 2.4 Tasks

To achieve the objectives described above, WP2 defines three main tasks:

**Distributed Data Management:** Providing infrastructures, techniques, and policies for managing storage resources in the Grid.

- **Storage Infrastructure:** Replacing existing high-end scalable storage systems with commodity physical storage devices, controllers, and interconnects within Grids and examining how current storage systems can migrate to this new architecture.
- **Providing Management Mechanisms:** Providing techniques for automatically managing storage resources in Grid and providing “high-quality” storage at low cost to users.
- **Specifying Management Policies:** Examining the different classes of storage services that could/should be offered to users and description methods and techniques for specifying service classes and management policies.

**Information and Knowledge Management:** Developing metadata, semantic representation, and protocols for Grid service discovery, information management and design of designing knowledge-oriented Grid services.

- **Semantic Modelling:** Developing metadata for Grid service discovery and information management and the design of designing knowledge-oriented Grid services
- **Semantic Representation:** Exploiting Semantic Web technologies for sharing machine-readable Semantic Grid models and techniques for knowledge intensive applications.
- **Agent Infrastructure:** Development and use of agent technologies to exploit semantic representation of users and resources to support workflow and knowledge management across distributed virtual organizations in science and business.
- **Standardisation and Integration:** Extending and standardizing the existing OGSA middleware for knowledge-based Grid services.

**Data Mining and Knowledge Discovery:** Design of Grid resource semantic mapping, database querying on Grids, and services and for distributed data mining and knowledge discovery on Grids.

- **Semantic Mapping:** Representation and mining of relationships between different Grid entities and resources.
- **Intelligent queries:** Query mechanism and intelligent agents for query formation,
- **Distributed Grid Services:** Design of services, and tools for distributed data mining and knowledge discovery on Grids, with Grid-aware highly adaptive data mining algorithms, considering data integrity and privacy.
- **Monitoring services:** Services providing accurate estimates of the cost of data mining tasks on Grids.

## 2.5 Drivers

The main motivation for the work performed in this WP is the pressing need to store, manage, and access more and more digital information. The increasing dependence of our society on information requires building infrastructures for storing large amounts of online information. Moreover, the complexity of managing stored information requires new techniques for automating the management process. Finally, new applications require high level abstractions for accessing storage (e.g. content-based search and resource discovery), as opposed to existing block- and file-level methods. This WP aims at consolidating research in these issues by structuring work in the tasks (layers) mentioned above.

## 3 Positioning

### 3.1 State of the art (existing approaches)

This section outlines existing approaches and their problems, tradeoffs, and limitations.

#### 3.1.1 Distributed Data Management

Existing storage systems are usually custom-built and custom-tuned to offer both scalability and good performance at high cost. Building large installations from low-cost commodity components would make data Grids less expensive and enable them to closely track the latest technological advancements of inexpensive mass-produced storage devices. Recent advancements in commodity interconnection network technologies and the continuing reduction in disk prices present new opportunities to address these issues. Various projects [REG03] currently strive to address the following issues:

- Build scalable storage systems that can hold petabytes of storage in a cost-effective manner.
- Make the storage infrastructure location-independent and client-agnostic in an efficient manner.
- Provide solid benchmarking methodologies.

#### *Building large-scale storage systems*

Building large-scale distributed data storage systems faces the problem storage resource virtualization incompatibility. The incompatibility results from the different levels of abstraction in the resources virtualization and the lack of a single, common framework for describing the storage services offered by the virtualized resources. Such a unified method of describing the services would make it possible to interface the different types of storage services in an effective way and would offer the starting point for building large-scale storage infrastructures. Many initiatives develop techniques for virtualizing the resources. For instance, initiatives like Lustre, GPFS, Frangipani and Petal etc. try to virtualize the blocks of filesystem by distributing them onto many storage nodes. The projects like Storage Resource Broker aim at the virtualization of files, file repositories and similar structures. P2P systems like Gnutella, Kazaa etc. use yet another level of abstraction of the virtualized data storage resources, providing access to files on the basis of metadata e.g. the song title or artist name.

Unfortunately, these approaches remain incompatible and different scientific data Grids remain isolated. While the methodologies for building the local, metropolitan, country-wide or continent-wide storage installations are quite well-developed, connecting them with other installations remains in the area of pioneer works. The situation in the commercial world is similar. Big market players like IBM, EMC, Cisco are able to interconnect their storage devices even by using 500-km long optical links while providing the possibility of updating the data synchronously from one device to the distant one. However, collaboration between the products of various vendors is still very difficult. Another problem is that the large data networks based on IP protocols and Grid middleware built within the confines of projects like DataGRID, European DataGRID, EGEE are tuned for specific applications. They allow the distribution of huge amounts of data produced by scientific equipment and their design is reasonable from the point of view of their developers, owners and founders/sponsors. But even if they are connected locally, there is no common approach to joining them into collaborating infrastructures.

Another issue is that the commodity components (inexpensive disks connected to PC machines), which seem to be the future of data Grids, rarely collaborate with the high-end storage components (like specialised disk matrices, tape libraries, hierarchical systems etc.) in order to realise common or complementary services. In many cases, high-end devices supply the services for the low-end ones, but it is uncommon to have both kinds/classes of devices available to end-users. Finally, the virtualization of the resources in Grids often loses the interesting features of these resources. The virtualization is needed to include the resources into Grids, but it often hides the interesting features of the resources like high-performance which are exploitable only using the native resource interface.

In summary, there are no common methods for describing storage services realised at various levels of abstraction by different elements of the storage infrastructure, and there is no common approach to the problem of virtualizing the storage resources without loosing their specific, interesting features. Without a means to effectively exchange the storage resources and services, the different data Grids remain incompatible, isolated and ineffective.

### **3.1.2 Information and Knowledge Management**

#### *Information Modeling*

The need to have common, platform-independent standards for representing Grid-related information has been recognized and is currently the subject of a number of projects and working groups. These efforts have been triggered primarily by the need to enable the interoperability between large, heterogeneous infrastructures and by the emergence of Open Grid Services. One of the earliest efforts in that direction comes from the DataTAG, iVDGL, Globus, and the DataGRID projects, which collaborated to agree upon a uniform description of Grid resources. This effort resulted in the Grid Laboratory Uniform Environment (GLUE) schema, which comprises a set of information specifications for Grid resources that are expected to be discoverable and subject to monitoring [GLUE04,ASV03]. GLUE represents an ontology that captures key aspects of the Grid architecture adopted by large Grid infrastructures deployed by projects like DataGRID, CrossGRID, the Large Hadron Collider Computing Grid (LCG), and EGEE. The GLUE ontology distinguishes two classes of entities: system resources and services that give access to system resources.

Going beyond the standardization of resources and services, a number of recent efforts are trying to devise common information representations for the structure and the status of jobs running on Grids. For example, the Job Submission Description Language Workgroup of the GGF (JSDL-WG) develops the specification of the Job Submission Description Language, an XML Schema for describing computational batch jobs and their required execution environments. Another effort, led by the CIM Grid Schema Workgroup of the GGF, seeks to standardize the information that could be published by Grid schedulers about the characteristics and status of Grid jobs submitted for execution. This workgroup has adopted the Common Information Schema (CIM) of the Distributed Management Task Force's (DTMF) [DTMF03]; based on CIM v.2.8, the GGF CIM workgroup of GGF has proposed a Job Submission Interface Model (JSIM) to describe the structure and attributes of batch jobs that run on Grid infrastructures. Finally, the need to provide basic Grid-job accounting and resource usage information in a common format is addressed by the Usage Record (UR-WG) and the Resource Usage Service (RUS-WG) workgroups of the GGF. These workgroups have started working towards the proposal of XML schemas that will describe accounting information in a general, platform-independent, way.

#### *Semantic Modeling*

Because of the lack of a global schema for Grid information, several researchers are investigating the application of semantic Web technologies as an alternative for bridging the gap that exists between infrastructures with incompatible information schemas. One of the earlier efforts came from the Grid Interoperability Project (GRIP) [GRIP04]; GRIP introduces two ontologies representing the structure and attributes of UNICORE and GLUE resources, respectively. These ontologies are described in XML and fed into a tool that supports the semi-automatic association between the two ontologies; this association is used for the mapping of resource requests to hardware resources that belong to Globus and UNICORE infrastructures [BFGC04]. A similar approach for the development of an ontology-based resource matchmaker was proposed by Tangmunarunkit, Decker and Kesselman in [TDK03]; their system comprised a matchmaker, which consisted of three components: (i) an ontologies component, which represents the domain model and the vocabulary for expressing resource advertisements and resource requests; (ii) a domain background knowledge component containing rules that express axioms, which cannot be expressed with an ontology language; (iii) a set of

matchmaking rules, which define the matching constraints between requests and resources and are expressed in a rule-based language. An ontology editor is used for the development of three domain ontologies for resources, requests, and applicable policies; these ontologies are described with the RDF-Schema specification of W3C. Matchmaking is conducted with the help of a deductive database [TDK03].

### 3.1.3 Data Mining and Knowledge Discovery

The Grid can be effectively exploited for deploying data-driven and knowledge discovery applications. It is a well-suited infrastructure for managing very large data sources and providing high-level mechanisms for extracting valuable knowledge from them. To perform this class of tasks, advanced tools and services for knowledge discovery are vital. Today research teams are devising implementations of knowledge Grids in terms of the OGSA model. According to this approach, knowledge Grid services are exposed as a persistent service, using the OGSA conventions and mechanisms. In future years the Grid will be used as a platform for implementing and deploying geographically distributed knowledge discovery and knowledge management platforms and applications. Some ongoing efforts in this direction have been recently started.

1. OGSA-DQP: In collaboration with the University of Newcastle, a system for distributed query processing on the Grid has been developed [[www.ogsadai.org.uk/dqp](http://www.ogsadai.org.uk/dqp)]
2. Polar\* : investigation and prototype development of techniques for adaptive query processing on the Grid and resource scheduling on the Grid [[www.ncl.ac.uk/polarstar/index.htm](http://www.ncl.ac.uk/polarstar/index.htm)]
3. KNOWLEDGE Grid: a prototype is available from University of Calabria running on Globus 3.0 and a WSRF-based version is under development to offer OGSA-based services for knowledge discovery in Grids.

#### *Grid-enabled database query system*

Structured and unstructured query systems are key components in data Grids. To date, Polar\* [SGW+02] and OGSA-DQP [AMP+03] are the only fully fledged generic Grid-enabled query processors, but there is an increasingly growing interest in Grid databases. For example, SkyQuery [MSBT03] applies DQP over Grid-connected databases that contain astronomical data. The execution approach that it follows has similarities with OGSADQP/ Polar\*, e.g., calls to Ws are regarded as typed UDFs. The main differences is that OGSA-DQP (i) accesses Grid rather than Web services, (ii) supports partitioned parallelism, (iii) can employ Grid machines that may not hold data in order to evaluate parts of the query plan, and (iv) is generic with respect to the underlying databases supported and not tailored to a specific scientific scenario.

GridDB-lite [NCK+03] is a project motivated by data-intensive scientific applications on the Grid, built upon DataCutter [BFK+00], in which the users express their retrieval tasks as SQL-like queries. However, the query is not evaluated using database technologies. Overall, GridDB-lite is benefited from the declarative manner of expressing potentially complex tasks in query processing, but develops its own execution mechanisms, thus not exploiting the full potential of a DQP system. Another project that supports database table interfaces for data processed in a workflow is GridDB [LFP03]. As in GridDB-Lite, this feature enables declarative task expression. However, GridDB takes one step further, and employs techniques devised for adaptive query processors to prioritize partial results. Generic interfaces to Grid databases have been developed in two European projects, OGSA-DAI [[www.ogsadai.org.uk](http://www.ogsadai.org.uk)] and European Datagrid's Spitfire [BBH+02].

#### *Design of Grid-aware distributed tools for data mining and knowledge discovery*

Currently large amounts of data are continuously collected in distributed sites, and data mining (DM) is emerging as a new discipline that is able to furnish tools and techniques to support knowledge extraction and decision making. This knowledge extraction process is both computationally intensive, and collaborative and distributed in nature. So in the last years many distributed data mining (DDM)

algorithms have been proposed [KC00, PK02]. DDM algorithms, when employed to devise Grid services and tools, must deal not only with distributed sources of huge datasets and multiple compute nodes, but also with distributed user community and privacy concerns. A further emerging challenge regards the updating of mined knowledge when the databases are also dynamic and evolving. For example, consider a warehouse continuously updated by streams of information [W02, MM02]. Grid-aware DDM services and tools must be adaptive with respect to data and platform features [OPP02], able to decide whether a tightly-coupled [OPP02a] or a loosely-coupled approximate solution [SO05] must be adopted, and optimize the use of resources [OPP02b, POP04].

The data and information patrimony today available can be effectively exploited if used as a source to produce the knowledge necessary to support decision making. Examples of large and distributed datasets available today include gene and protein databases, network access and intrusion data, drug features and effects data repositories, astronomy data files, and data about web usage, content, and structure. The knowledge extraction process is both computationally intensive, and collaborative and distributed in nature. Unfortunately, the number of high-level instruments to support the knowledge discovery and management in distributed environments is very low. This is particularly true in Grid-based knowledge discovery [Berman, 2001], although some research and development projects and activities in this area are going to be activated mainly in Europe and USA, such as the Knowledge Grid, the Discovery Net, the AdAM, and the DataMiningGrid projects. In particular, the Knowledge Grid [Cannataro and Talia, 2003] provides a middleware for knowledge discovery services targeted to a wide range of high-performance distributed applications. The DataMiningGrid project [DMG05] is a EU STREP project that is developing tools and services for deploying data mining applications on the grid.

## 3.2 Extended context

### 3.2.1 Distributed Data Management

Hippodrome [AHK02] at HP Labs automates storage device configuration in large data centers. After identifying disparities between workload measurements and user specifications, Hippodrome finds appropriate data placements and device configurations to better match the performance requirements of the system users. Similar approaches in larger scale could make data Grid management more cost-effective and less dependent on tedious manual setups amenable to human error. Work in full-system virtualization tries to address resource sharing and management at the system level. For instance, Vmware [GTH00] allows the sharing the resources of large servers among multiple users. One challenge is to keep the sharing cost low, while offering performance guarantees to different users coexisting on a server. Another challenge is to be able to virtualize the resources of server farms with specified performance goals per user.

#### *Storage backend Gridification*

The idea of pervasive Grid services remains a dream rather than the close reality. This is caused by the lack of standardized, well-understood description methods for data-related services, “deep” enough to present the abilities of a given system to the external world. For instance, the Storage Resource Broker [AGC05] from the San Diego Supercomputer Center has been deployed broadly during the last few years to integrate distributed heterogeneous data storage servers. It combines relational database technology with custom-built middleware software to offer unified view of disparate data sources. Even though SRB can access successfully a variety of data sources including ftp servers, database systems, file systems, and web servers, its integration with legacy applications manifests several performance problems which makes its use inflexible.

PSNC's everyday operation as the storage-service provider teaches us about the difficulty of interconnecting storage systems on the one hand and offering the storage services to the clients on the other. In fact, it is hard to bring together the easiness of using the services and the possibility of exploiting the storage infrastructure while not losing their advanced features, e.g. predictability, data

access optimisation techniques, guaranteed level of service, security, confidentiality and safety of data. If one decides to have an easy access method to the service, one loses the awareness of its way of operation, state and many specific features of the resource which often decide its use. However, considering all the complexities of a given resource while implementing services on top of it, is so difficult that the advantages of resource-specific-aware access methods are not worthwhile.

### *Storage System Benchmarking*

Benchmarking is an essential tool in evaluating new techniques. In order to test a wide variety of features and be more realistic, benchmark programs have become more complex. As a consequence, evaluating new solutions is becoming a very time consuming task. Therefore, there is an increasing interest in techniques that focus on keeping the benchmarking results accurate while reducing their execution time. Such techniques are usually based on statistical methods. Two important examples of such work in this area are: benchmark application sub-setting techniques using Principal Component Analysis (PCA) [VD04b] or Cluster analysis [VD04a], and execution phase analysis in order to reduce the execution time of each application [PHV03, SPH03].

One of the goals of WP2 is to study and propose new storage solutions for the Grid platform. Benchmarking will be used in order to evaluate these solutions. For this setup, a realistic benchmark is the set of queries from the standard decision support system benchmark known as TPC-H [TPC03]. As the objective for the experiments is to test the I/O system, the size of the database must be considerable large. This results in large execution times. As such, it is common practice to select a set of queries instead of executing the complete set from the benchmark. WP2's contribution to this task is to apply the techniques above described (application sub-setting and execution phase analysis) to this benchmark. Previously, these techniques have only been applied to the scientific SPEC benchmark suite and the metrics used were execution time and memory system metrics such as cache misses. The metrics and the criteria used for I/O benchmarking are significantly different from the ones used for architecture benchmarking. Therefore we expect that the application of the previously developed techniques to the new setup is a considerable challenge.

### **3.2.2 Information and Knowledge Management**

The means used for representing and publishing resource information, in typical Grid middleware like Globus or UNICORE, do not aim to support sophisticated, user-customized queries or allow the user to decide from a number of different options. Instead, they are tied to the job submission needs within the particular environment. As we move towards a fully-deployed Grid - with a massive and ever-expanding base of computing and storage nodes, network resources, and a huge corpus of available programs, services, data, and logs - providing an effective service related to the availability, the characteristics, and the usage of Grid resources can be expected to be a challenging and complex task. As discussed earlier, efforts to address this problem are focusing on the development and standardization of information schemas (mainly defined in XML or RDF) for the description of Grid-related information. Such schemas, however, often overlap in scope and there is a clear need to re-use existing or emerging standards. Most standardization efforts, however, are still at a very early stage of development and are not adopted by new middleware systems that emerge with an increasing pace. Therefore, it is practically impossible to materialize the vision of a widely established collection of mutually compatible schemas for encoding Grid-related information. On the other hand, the use of Semantic web technologies (ontologies, rule-based reasoning and semantic matching) faces known scalability limitations, although it enables the resolution of complex queries upon information bases spanning across syntactically incompatible infrastructures. Finally, if we draw from the WWW experience, the identification of interesting resources has proven to be very hard in the presence of too many dynamically changing resources without well-defined rules for classifying the degree of relevance and interest of a given resource for a particular user.

Searching for information and services on the Web typically involves navigation from already known resources, browsing through Web directories that classify a part of the Web (like Yahoo!), or

submitting a query to search engines. In the context of the Grid, one can easily envisage scenarios where users may have to “shop around” for solutions that satisfy their requirements best, use simultaneously different middleware systems (which employ different ways to publish resource information), or consider additional information (such as, historical or statistical information) in choosing an option. The integration of data discovered in and retrieved by those sources can help in the establishment and maintenance of knowledge bases for the Grid that could provide answers to various end-user queries [DSI05].

The main objective of Knowledge Grids is to allow heterogeneous Grid applications that are created by different independent organizations to interact for the transfer of data, service operations and knowledge. This objective relates the classic ones of heterogeneous database interaction and heterogeneous service interaction as well the issues key to service composition and agent interaction. Currently there is a gap between Grid computing endeavours and the vision of Grid computing. This vision sees a Grid in which there is a high degree of easy-to-use and seamless automation and in which there are flexible collaboration and computation on a global scale [GEL04]. To support the full richness of the Grid computing vision requires Semantic Web technologies for Grid middleware and applications, i.e. the Knowledge Grid. From a technical point of view, there are several challenges for Knowledge Grids. From the Semantic Web point of view, it is necessary to develop semantics for privacy, security and access-rights as well as dealing with dynamic information, state, QoS and states. From the Grid point of view, it is important to move from fixed-pipeline processes to dynamic compositions.

### 3.2.3 Data Mining and Knowledge Discovery

#### *Semantic Web*

Semantic web technologies have been proposed as a platform for the discovery of information about software and services deployed on the Grid. An early approach comes from the ICENI project in UK, which focuses on the semantic matching between Grid services and service requests in an autonomic computing context, even when requests and resources are syntactically incompatible [HLN03]. To achieve this goal, the ICENI project proposed the concept of a metadata space. This is an environment distinguished from the space of Grid services and resource requests. The metadata space hosts Grid-related semantic metadata, published and discovered through standard protocols. The operation of the metadata space is supported by meta-services providing semantic matching and service adaptation capabilities. Service adaptation refers to the automatic adaptation of a Grid service's output to the requirements of a semantically matched but syntactically incompatible resource request. The ICENI approach was demonstrated in the case of a very simple adaptation scenario [HLN03].

#### *Workflows*

Knowledge discovery procedures typically require the creation and management of complex, dynamic, multi-step work-flows. At each step, data from various sources can be moved, filtered, integrated and fed into a data mining tool. Examining the output results, the analyst chooses which other data sets and mining components can be integrated in the workflow or how to iterate the process to get a knowledge model. Workflows are mapped on a Grid by assigning abstract computing nodes to Grid hosts and exploiting communication facilities to ensure information/data exchange among the workflow stages. One of the aims is to consider the problem of Grid services description (semantic representation) that we address with an ontology-based approach. We put the focus on non-functional aspects of services. We plan to investigate the usage of ontology instances (individuals) in order to enable inferences and matchmaking. One more goal, is to define Inferential Monitoring and Management algorithms that are the bases for designing Grid Agents system.

Currently, several Grid-middleware components collect, store, and publish collections of information that can be useful to Grid systems and users. These collections include:

- Information that describes the capabilities, the operation, the status, the pricing, and the usage of hardware resources available on the Grid.
- Metadata about services deployed on the Grid, such as descriptions of functionality and interface, guidelines for invocation, and policies of use.
- Metadata regarding data and software repositories deployed on the Grid, describing their organization, contents, semantics, and relevant policies of access and use.
- Job management information regarding jobs deployed on Grids: their composition in terms of software or service components, their mapping to hardware and networking resources, their cost, etc.

Information on the capability and status of Grids is typically collected and maintained by a variety of Grid-middleware sub-systems or Grid-application components, which are characterized as Grid information and/or monitoring services, although the boundaries between these two categories are not clearly defined. Grid-related information is also collected and maintained by other components: job management information is typically maintained by resource brokers, workflow engines, logging servers, etc; information about data repositories can be found in data-Grid services, such as replica catalogues, virtual file systems, and application-specific data archives.

Notably, different information sources employ diverse data models, formats, and encodings for the representation and storage of information. Some sources make their data available to third-parties (i.e., to other services, administrators or end-users) by providing support for binding, discovery, and lookup through a variety of protocols and query models. Because of the lack of a standard model or a common schema for organizing and representing information, it is difficult to establish the interoperation between different Grid platforms. Moreover, the lack of common information models and standards makes it practically impossible to achieve the automated retrieval of resources, services, software, and data, and the orchestration thereof into Grid work-flows that lead to the solution of complex problems.

The discovery and matching of bioinformatics workflows deployed on the Grid is the goal of the myGrid project [myGrid04], which provides mechanisms for the search and discovery of pre-existing workflows based on their functionality (“task-oriented” or “construction-time” discovery), on the kind and format of their input data (“data-driven” discovery), or on the type and format of their output data (“result-driven” discovery). To make workflows discoverable, myGrid introduces the workflow executive summary, a workflow-specific collection of metadata represented in an XML Schema [MPWLCM04]. Metadata belonging to the workflow executive summary include: (i) mandatory descriptions of the workflow's definition (e.g. its URI address, its script, its invocation interface, the types of its input and output data); (ii) optional syntactic descriptions about the format encoding of the workflow's input and output data, and (iii) optional conceptual descriptions of the workflow's characteristics. Workflow executive summary information is encoded in RDF with additional pointers to semantic descriptions described in OWL. Two key modules in the myGrid system architecture are the registry and the semantic find component. MyGrid's registry is designed to accept and store workflow descriptions, in accordance to the UDDI specification. Furthermore, it supports the annotation of stored workflows with conceptual metadata [MPPDM04]. MyGrid's semantic find component is responsible for executing OWL queries upon the conceptual metadata attached to the workflow descriptions stored in myGrid's registry. Each time the semantic-find component receives notifications about metadata newly added to the registry, it updates accordingly an index with metadata descriptions. This index is used for fast replies to semantic queries. Alternatively, it can invoke a description-logic reasoner to answer semantic queries. Also the DataMiningGrid project is working in this area by developing grid-enabled data-mining data interfaces and services, grid-enabled data-mining workflow management tools to provide a workflow editor that facilitates the composition, execution and management of data-mining workflows in grid computing environments.

## 4 Vision, Strategy and Roadmap

### 4.1 Vision and Scenarios

Research activities in the area of data Grids and knowledge Grids are being pursued in Europe, in USA and in ASIA by several research teams and, at the same time, companies such as IBM, HP and SUN Microsystems, are very active in the area. This demonstrates the key role of data management in Grids and the importance of developing knowledge-based applications that exploit the Grid features to achieve high performance and high availability.

The research tasks that compose Workpackage 2 give a unified vision of the data and knowledge management in Grids through a layered approach that starts from efficient data storage techniques up to information management and knowledge representation and discovery. The main vision of this Workpackage is based on a common model and framework that can integrate the research results of the involved partners and will result in future common activities that will advance the present results and systems.

According to this vision, we identified the main phases for this roadmap that drive the research activities of the partners in implementing the WP2 activities:

- **Phase 1:** Exchanging partner information, experiences, and knowledge about techniques, tools and systems for Data and Knowledge Grids.
- **Phase 2:** Sharing and integration of common goals, research results, projects and system prototypes of Environments and Services for Data and Knowledge-based Grids.
- **Phase 3:** Use of the results of the previous phases for envisioning a unified framework for handling data, information and knowledge on Grids. Definition of joint proposals, research activities and projects in the WP2 research area.

### 4.2 Strategy

The strategic approach of this Workpackage is give an integrated view of data and knowledge management in Grids through the three subtasks of the Workpackage that investigate the three main layers of data and knowledge-based systems.

WP2 follows a layered approach in addressing data and knowledge related issues in Grid systems. The lowest layer (Task 2.1) deals with systems-level, distributed storage management issues. The middle layer (Task 2.2) explores techniques that will turn storage systems into knowledge representation systems. Finally, the top-most layer (Task 2.3) addresses issues in automatic mining and resource discovery techniques. Here we outline the main scientific activities planned in the three tasks that compose our workpackage. In section 4.3, an extensive description of the roadmap for each research activity is given. In that section, for each planned scientific activity in GRID computing, the state of the art is given by discussing GRID and non-GRID research and development activities.

#### 4.2.1 Planned Scientific Activities in Distributed Storage Management

The planned scientific activities of Task 2.1 are centered around the goal of advancing scientific results and technologies in the areas of storage management infrastructures, mechanisms and policies.

- **Storage Infrastructure:** Studying the replacement of existing high-end scalable storage systems with commodity physical storage devices, controllers, and interconnects within Grids and examining how current storage systems can migrate to this new architecture. The activities undertaken by project partners in this direction are related to research performed currently both in leading industrial and academic organizations for reducing the cost of ownership of storage facilities. All efforts in this direction, essentially aim at replacing existing centralized solutions (e.g. large, centralized storage controllers) with more distributed solutions, at varying degrees of

distribution and heterogeneity, depending on the target applications. Thus, our work in the context of CoreGRID is inline with these efforts and has the additional advantage that it interfaces with current efforts to build general-purpose GRID environments and infrastructures.

- **Providing Management Mechanisms:** Providing techniques for automatically managing storage resources in the Grid and providing “high-quality” storage at low cost to users. Besides the ability to built storage infrastructures in a cost-effective manner, it is also important to reduce cost of ownership (maintenance). Task activities in this direction aim at building in the infrastructure mechanisms that are not available today (due to the limitations of existing storage architectures) that will lead to lowering the cost of managing storage. This work is related to issues also important to autonomic computing initiatives, which align very well with goals in Grid research.
- **Specifying Management Policies:** Examining the different classes of storage services that could/should be offered to users and description methods and techniques for specifying service classes and management policies. Once new storage architectures and low-level mechanisms are available, it becomes important to examine how these can be used by high-level policies to facilitate building new services and satisfying application needs. This work in storage systems is to some extent related to parallel efforts in specifying policies at the network level for managing network infrastructure, however, they target storage systems. Performing this research in the context of CoreGRID allows us to be in touch and exchange ideas with CoreGRID researchers from the networks community.

#### **4.2.2 Planned Scientific Activities in Information and Knowledge Management**

The objectives of Task 2.2 can be grouped around the following themes:

- *Semantic Modeling and Representation:* Developing metadata models and systems for Grid service discovery and information management and the design of knowledge-oriented Grid services. Exploiting Semantic Web technologies for sharing machine-readable Semantic Grid models and techniques for knowledge intensive applications. Task activities in this direction aim at building a core Grid ontology, which can provide a common basis for representing Grid knowledge about grid resource, grid middleware, services, and applications.
- *Distributed Information Management:* Studying the mechanisms of applying the semantic technique in the Grid information management. Scalability is the key issue to be addressed of this research area. Task activities in this direction are included: (1) designing and developing a semantic information system to collect, manage, and share the distributed semantic Grid information. In particular, an ontology-based Grid information integration sub-system will be developed to merge the different data/information about the same Grid entity of a Grid system from different information sources; (2) providing a semantic-based infrastructure to share the distributed information efficiently. Examining the distributed and centralized mechanisms for sharing semantic Grid information regarding to the scalability.
- *Distributed Query Processing over Grid Data Resources & discovery services:* Ontology-driven Query mechanisms over Grid data resources are investigated. Research activities undertaken by project partners in this area are to design an ontology-driven query algorithm for querying the distributed semantic metadata. Moreover, research activities includes designing/developing scalable information discovery services to locate data and generic Grid resources, like physical resources or Grid services identified and described by suitable metadata, by coping with aspects of dynamicity, at both system configuration and information level.
- *Agent Infrastructure :* Analyzing the use of agent technologies to exploit semantic representation of users and resources to support workflow and knowledge management across distributed virtual organizations in science and business.
- *Service-based data integration on Grids:* design of scalable P2P models and architectures for distributed data integration. In a distributed setting, data GRIDs must manage data coming from different source and having different structure. In this scenario data integration is a key issue. Schema integration models and algorithms will be studied and developed.

### 4.2.3 Planned Scientific Activities in Data Mining and Knowledge Discovery

The scientific activities of Task 2.3 on data mining and knowledge discovery in GRIDs can be grouped around the following main themes:

- *Semantic Mapping, services discovery and negotiation*: Studies about the representation and mining of relationships between different Grid entities and resources. This research activity relates to semantic web technologies and investigates semantic description for service discovery, management and security. Technologies such as OWL-S, XACML and P3P are used.
- *Intelligent queries*: Query mechanisms and intelligent agents for query formation are investigated in the context of the OGSA-DAI and OGSA-DQP systems and services. Research teams involved in those key projects in the area of data access and querying in GRIDs are designing a framework for defining and constructing adaptive query processing (AQP) systems both for Grids and for more traditional environments. Moreover, a generic mechanism for extracting monitoring information from the query execution to support multiple AQP techniques will be considered.
- *Distributed Grid Services*: This research area deals with the design of services, and tools for distributed data mining and knowledge discovery on Grids, with Grid-aware highly adaptive data mining algorithms, considering data integrity and privacy. The two main scientific activities in this them concern the design and implementation of high-performance data mining algorithms and the development of knowledge discovery services for supporting the implementation of distributed data mining application on GRIDs. These scientific activities represent an advancement of the state of the art in distributed data mining since they will leverage OGSA and WSRF technologies for enabling the exploitation of data mining technologies in GRID environments. This will result in knowledge discovery WSRF web services.
- *Monitoring and scheduling services for data mining*: Services providing accurate estimates of the cost of data mining tasks on Grids can improve the performance of distributed data mining applications and services on GRIDs. Research activities in this area will find synergy with activities carried out in the EU NextGrid project for predicting the performance of data mining tasks on the basis of statistical analysis on the source dataset. Minimum Completion Time heuristic strategies are also investigating for scheduling high performance data mining tasks on a Knowledge Grid.

## 4.3 Roadmap

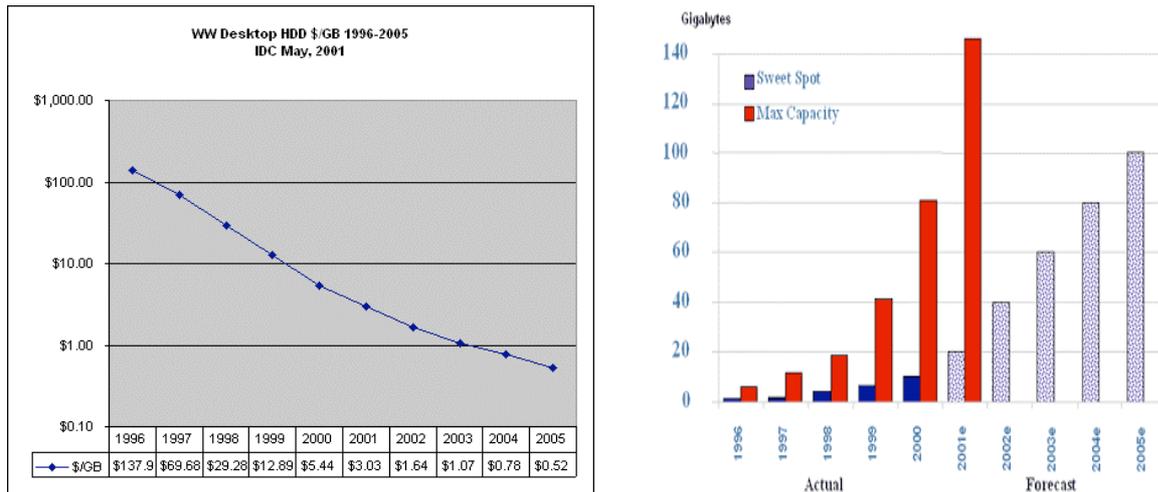
### 4.3.1 Distributed Storage Management

#### *Storage Infrastructure*

Over the last few years scalable computing infrastructure has undergone significant architectural changes. New technologies and trends have led to commoditization of components and systems that previously relied on expensive custom designs. Today, scalable computing infrastructure is based on loosely- or tightly-coupled servers (PCs or workstations) interconnected in larger Grids. Loosely-coupled servers usually take the form of clusters or desktops interconnected with local area networks, whereas tightly-coupled systems make use of low-latency high-bandwidth interconnects and provide single system image to applications. Over the last few years, a lot of research has also been conducted in building larger, tightly-coupled parallel systems based on commodity PCs and blade servers. We observe similar needs and trends in the scalable storage infrastructure.

First, demand for on-line storage has been growing dramatically over the last few years. For example, the Internet is currently used not only by people, but also by agents, bots, and spiders that search, retrieve and process information on behalf of their owners. It is speculated that within a few years, most of the information on the network will never be seen by a human; instead it will be entirely created and consumed by computer applications, which will skyrocket the amount of stored

information [GRA99]. Furthermore, a growing number of thin clients access the Internet and demand novel services, including handling storage-hungry rich-media data. Rich-media production, editing, and distribution have increased storage needs dramatically both in various industry sectors as well as end users. Secondly, magnetic disks cost per Gbyte of stored information are dropping rapidly (Figure 1), leading many application domains to all-on-line storage systems, where all available information is stored on disks. With disk capacities approaching the TByte-level in 2005, it becomes feasible to store most of the information produced [LVA03] online and to provide continuous access to users and applications. Thus, both application needs and the underlying technologies push towards the creation of even larger disk-based information repositories. A central challenge for the information infrastructure is to build the scalable disk-based storage systems that will be able to accommodate future storage needs.



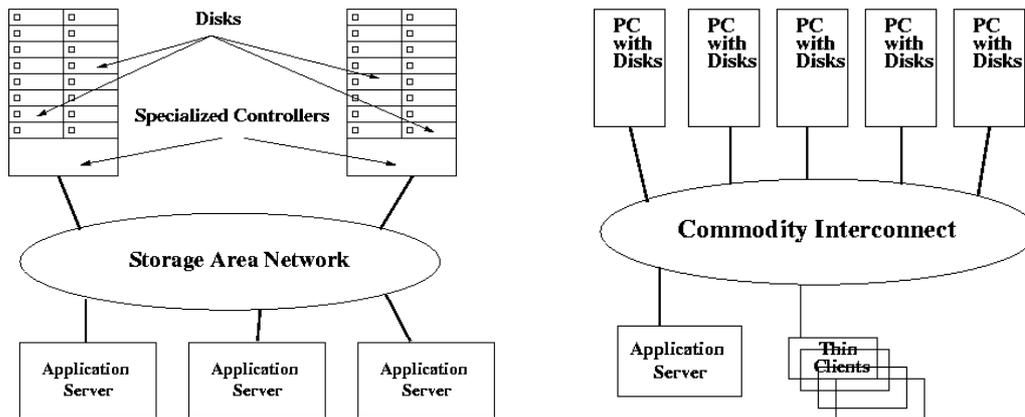
**Figure 1:** Magnetic disk cost per GByte of stored information and magnetic disk capacity trends. Note that (left) capacity about doubles every year and that (right) storage cost drops by about 1.3-1.8 every year.

To understand the issues involved, examine the requirements imposed on storage systems today by existing and new applications:

1. **Cost-effective scalable storage (storage scalability):** Improvements in CPU speeds allow cost-effective generation and processing of large amounts of data. Improvements in intra- and inter-network capacities result in the ability to rapidly transfer large quantities of data between different locations. In general, our ability today to generate, process, and transfer large amounts of data leads to increased demands for storage that are projected to double every year [GRA99]. However, existing services and applications do not become cost-effective, despite the lowering CPU processing and memory costs.
2. **High-quality storage (storage virtualization):** Applications require not only large amounts of storage but high-quality storage as well. *High-quality* is a term we use to describe storage that provides advanced functionality. Scalable storage systems of the future need to be intelligent and capable of adapting to application needs [GRA99]. For instance, they need to shrink and expand on demand, be robust, redistribute data based on application needs, tolerate catastrophic failures, etc. Intel shows that 55% of server downtime is related to problems with the storage subsystem.
3. **Easily-accessible storage (storage Gridification):** The increasing number of portable devices or thin clients and the general desire for seamless access to information results in requirements for storage systems to facilitate such applications. Moreover, due to the limited processing capabilities of thin clients, there is a need for pushing storage-related functionality towards lower layers and into the storage systems. Scalable storage systems require extensive and expensive IT expertise to deploy and maintain, making it practically impossible to deploy on demand and in arbitrary locations.

Our goal is to explore important problems that have not been addressed yet and may lead to solutions that will satisfy application needs. In order to address these issues, research is required that will provide the enabling technologies for building future storage systems:

- Use generic commodity components and develop the required technology to build future storage infrastructure.
- Provide location-independent and client-agnostic storage to reduce access costs and enable new applications.
- Reduce management complexity and cost by identifying the key contributing factors and providing mechanisms for policy-based storage management.



**Figure 2:** Current (left) and future (right) architecture of scalable storage systems.

### *Storage Management Mechanisms*

Ongoing research effort at FORTH-ICS focuses on the development of systems software that unifies off-the-shelf storage devices over commodity local networking equipment. Multiple server nodes contribute storage and computation resources towards providing the view of a single storage device with large space capacity and data transfer bandwidth that is both reliable and scalable. The software is structured in multiple layers each adding particular functions and abstracting the details of the layers lying underneath. Each system client should be able to connect to an arbitrary storage node and receive the illusion of accessing the entire collection of resources that the client is authorized to use. Ultimately, each client should be able to mount a storage system partition as a local file system and access it through the traditional file system interface. This implies that legacy applications can take advantage of the storage system resources without any modification requirements.

Particular features such as data reliability are enabled through specific software layers that add the necessary redundancy to the stored data and distribute them across distinct physical devices. When a physical device fails, the data available on the remaining devices can be used to reconstruct the missing data. Due to the potential scale of the storage system, estimations about the expected failure rate and the required reliability can be mapped to particular levels of redundancy and scale of data distribution. Interaction with task 4.5 (dependability) of WP4 will further enhance the reliable operation of the entire system. Having multiple users sharing the same resources can lead to contention and reduced performance. The plan is to continuously monitor the utilization of the resources and identify hot spots in data links or storage devices. Subsequently, we can apply dynamic data migration techniques to redistribute data and hide from the users the heterogeneous features of the storage system components. Such functionality has been investigated for several years in the storage systems community, and remains an open problem to a large extent. The plan is to further study issues of data administration automation, and develop solutions applicable to data Grid environments. Additionally, there will be interactions with task 5.1 (information and monitoring services) of WP5

(Grid information and monitoring services) on the problem of monitoring resource usage required for management automation.

The plan is to achieve scalability by combining hierarchical and peer-to-peer organization techniques in the structure of the system wherever necessary, and avoid single points of failure. Searching for data resources is another problem that shows up in distributed data management. Unlike small file systems where hierarchical organization of the name space usually provides means for finding data, in large data collections data a search becomes a challenging issue. Note this issue presents synergies with task 4.3 (scalable services) of WP4 (scalable architectures).

### *Storage Management Policies*

In the Internet age, the needs for exchanging large quantities of high quality numerical data are increasing. This trend has evolved from exchanging text, to image and photography today, and video tomorrow. However, software for exchanging such kind of data has remained basic and not very reliable. More reliable and efficient ways of exchanging this kind of data is required. A CETIC project called Filestamp answers this need, using a modern service-provider-independent solution. It will allow users to exchange files based on a peer to peer approach for the file exchange functions. The problem is that these large amounts of high quality numerical data need to be stored and retrieved. Current file systems can be used to store the data, but offer little support for searching and retrieving the data. The availability of storage services with guaranteed quality of service is needed. The availability of Grid storage services which offer guaranteed quality of service would improve the overall efficiency of the Filestamp approach.

At PSNC, we believe that the large data Grid structures of tomorrow need to encompass very different kinds of resources (disks, filesystems, virtual filesystems, repositories like P2P systems etc.). They should also satisfy the needs of larger groups of users, which have more diverse profiles (educational, scientific, government, business) and expectations resulting from these profiles. Therefore, future data Grids should implement very different kinds of services, from database front-ends to "virtual filesystems", coexisting in the same environment. The access methods ought to be described using a coherent framework, which allows for high level of abstraction on the one hand, and the very detailed description of the service features on the other. We see the place for both commodity components and high-end components as elements of the storage infrastructure in data Grids. In our vision, they can collaborate, exchange the virtualized resources, but they can also provide separate, disjoint services on the different guaranteed levels of service. The "good" services, that have the parameters guaranteed with a given probability, could be based on the commodity components. The "high-end" services for mission-critical applications could be implemented based on the high-end infrastructure elements. However, both kinds of the services should be available using the same, standardized framework.

Storage resource virtualization is a very attractive feature that should lead to the possibility of using the storage services seamlessly, without worrying about underlying components. "Mounting" the "general Grid filesystem" on the workstation is possible like plugging the cable into the socket would be a good marketing is a very attractive perspective. However, some of the potential users will not only worry about "give me the service", but they will be interested in knowing the values of the service parameters such as stability, availability, MTBF, max available "bandwidth", max amount of operations/transactions per second etc. Possibly they will also be interested not only in the chance of learning the behaviour of the infrastructure, but also in tuning them to their needs or at least the means to find in the Grid the access point that will satisfy their needs. In order to bring the pervasive data Grids into life without losing the native high-end features of the Grid resources we need to build the framework for describing the data Grid services, interfaces and access points in a common well-defined easily-understood and expandable way. This should allow the various data Grid infrastructures to interface with each other in an effective way and provide the potential clients (end-users, systems, applications) the means to automatically find the access points to the services they need.

To reconcile the contradictory requirements of the particular clients of the storage services we propose the expandable, stackable model of the storage infrastructure i.e. storage resource and storage services. We see this model as the stack similar to the stack of the network protocols. The lowest level relates to the actions like block data access, transmission, synchronization, etc. The higher services relate to operations like implementing the filesystems, transferring the files or making the replicas of them etc. Each intermediate layer uses the services of one of the lower layers (not necessarily the close neighbour) in order to implement its services. Each layer can also implement the services for the higher layers. We assume the possibility of crossing/skipping the intermediate layers of service in order to assure the possibility of using the native interfaces of the resources at the lower layers. However, the levels of abstraction of the particular layers should be standardized to facilitate the exchange of services between the different Grids. We hope that this approach, although not new in computing, can be useful in structuring the abstraction levels for the storage services and the storage infrastructure descriptions. In this approach, clients, end-users and subsystems could access the layer of the service stack that is most suitable for them, having more or less awareness about what is behind the stack layer they interface with.

In CoreGRID, we are going to examine the existing methods for defining the storage resources in Grids and attempting to unify them. Similarly, we will examine the models for describing the client's (end-users and Grid services) needs and try to find a common, universal model for them. The framework we try to develop can have different forms. This will be the subject of further investigations in CoreGRID. It can be somehow similar to the approach used in the World Wide Web. In the Web, having even the simplest web browser (like lynx) you can at least download and see the html-based text pages. Wanting to see the images, you have to upgrade to the graphical browser like MS IE or Mozilla. Going further, while you explore the web, from time to time, you find the web pages your browser cannot interpret, e.g. containing flash animations or VRML objects. Then you need the plug-ins, your browser can download automatically or you install them manually. Moreover, when you connect e.g. a special project's or corporation's internal page, your machine should be equipped in the smart-card reader in order to be able to authenticate you to the serves. Still, having even the simplest browser, you can use the basic functionality (getting the pages using HTTP protocol and interpreting them using the HTML engine). Such approach could be accepted for building large data Grids. Accessing Globus pools, you should just install the Globus plug-in to your "mount" command. Accessing a well-protected, ciphered filesystem you could use the smart-card reader plug-in. The framework could also make putting to the Grid both the not-guaranteed-level-of-service services and resources and guaranteed ones possible. Possibly, it provides the way to get money in Grids, since offering the storage services at the agreed level is possible. This could be a good starting point to make the industry more interested in Grids.

The comparison presented above is perhaps exaggerated, since storage services are more complicated than just get/post operations in HTTP protocols. However, like www made the common framework for providing different kinds of services in the Internet (from text documents to video streaming), the model we plan to work on should provide the common approach to provide storage-related services in data Grids. The layered structure of the model should help to join the different storage infrastructures nevertheless on what abstraction level they virtualize the resources. Along with the Grid information services, this model should be useful for the potential Grid clients looking for the access point to the storage service.

The PSNC's contribution to this work can be related with the experience connected with the services we provide to projects we have in the centre. From 2004, we provide storage services for the Police in our city. It is the backup/archive service for video streams collected by the video-cameras located on the streets. In this project we had to reconcile the features of the application controlling the data streams and the features of our storage devices and software. The application makes constant data transfers and the random accesses to the stored video files assuming the unlimited storage space of the filesystem, while the only service we could provide in this setup (because of software licensing policy, limited project budget etc.) was the backup/archive service. Introducing some intermediate steps in the

data processing allowed the implementation of the backup/archive service without the modification of the original application. The application got the impression of having the unlimited storage space.

The filesystem used by the application was virtualized using the put/get mechanisms provided by the storage software. In the R&D project funded by the Polish government we are implementing the secured storage service for the database system. The database is used for a PKI application (signing the local administration documents) so its backup/archive copy must be secured (need of confidentiality and integrity of data) with a very high level of protection. A huge amount of the backed-up data requires the ciphering to be hardware-aided. Again, we have to develop the intermediate layer allowing the database management service to use the secured backup/archive service as if it has used the normal data storage. Modification of neither the commercial database management system nor the storage software was possible. Therefore, the virtualization of the backup storage space was needed and performed by the application we develop.

We have prepared the R&D project entitled "National Data Store". Although it has been accepted by the Polish government, currently we are looking for a commercial partner to launch the work. In this project we plan the virtualization of the filesystem at the block level. Preliminary work on classifying the groups of the potential users of the solution was performed. In addition, the analysis of the different aspects of services offered by the solution was done during the preparations of the projects. We have analysed the performance, ease of use and configuration of the end-user systems, consistency of the data stored as well as the security feature at the various levels (starting from the code security, secure network technologies, finishing at the safe rooms for the storage devices). We believe that the knowledge gained by PSNC during the aforementioned projects and the experience in employing the real, productive storage services can be the important contribution to the universal model of the storage infrastructure and services.

### **4.3.2 Information and Knowledge Management**

#### *Semantic Modeling*

CETIC is involved in the HPC4U (Highly Predictable Clusters for Internet-Grids) Grid FP6 project. It is based on the fact that commodity-based Grids are spreading across the industry, supported by their excellent cost/performance ratio. However, they lack reliability, usability and manageability. The HPC4U project will address those issues and provide as a solution a generic and modular Grid middleware software covering multiple administrative domains to enable increased fault-tolerant level. The objective of the HPC4U project is to expand the potential of the Grid approach to Complex Problems Solving through the development of software components for a dependable and reliable Grid environments and combining this with Service Level Agreements (SLA) and commodity-based clusters providing Quality of Service (QoS). Development of HPC4U will take place in a Grid context following standards of the Global Grid Forum (GGF). The HPC4U results will provide Next Generation Grids with the possibility to guarantee the completion of Grid jobs and leverage the larger uptake of Grid environments.

In order for an end-user to access such a commodity based cluster through the Grid, he has to describe the job he wishes to submit, and the required qualities of service (QoS). Based on this service request a suitable service provider then needs to be identified, and the QoS needs to be negotiated with him. The approach that is being taken in the project is to use an ontology:

- to describe the service request and the required QoS
- to describe the services that are offered by different service providers
- to match the service request with the service descriptions
- to negotiate the QoS between the user and the service requestor and the service provider

An ontology for describing resources, and the qualities of service that can be guaranteed for these resources by the resource management system.

Data integration on Grids has to deal with unpredictable, highly dynamic data volumes provided by unpredictable membership of nodes that happen to be participating at any given time. Then, traditional approaches to data integration, such as the federation of database management systems (FDBMSs) and the use of mediator/wrapper middleware, are not suitable in Grid settings. Recently, several works on data management in peer-to-peer systems are adopting an integration approach not based on a global schema. According to those works each peer manages an autonomous information system, and data integration is achieved by establishing mappings directly among the various peers. All these systems allow for a decentralized, wide-scale sharing of data preserving semantics. The Grid community is devoting great attention at managing structured and semi-structured data such as databases and XML data. The most significant examples of such efforts are the OGSA Data Access and Integration (OGSA-DAI) and the OGSA Distributed Query Processor (OGSA-DQP) projects. However, none of those projects actually meets schema-integration issues necessary for establishing semantic connections among heterogeneous data sources.

Grid Data Integration System (GDIS) [COT04] is a decentralized service-based data integration architecture for Grid databases that we refer to as. The design of the GDIS framework at UNICAL has been guided by the goal of developing a decentralized network of semantically-related schemas that enables the formulation of distributed queries over heterogeneous, different located data sources. In order to effectively exploit the available Grid resources and their dynamic allocation, GDIS adopts a scalable P2P model allowing to map data in the most convenient manner. GDIS implements a P2P-based integration formalism whose key feature is the query reformulation algorithm: when a query is posed over the schema of a peer, the system will utilize data from any peer that is transitively connected by semantic mappings, and reformulate the given query expanding and translating it into appropriate queries over semantically related peers. The GDIS infrastructure exploits the middleware provided by OGSA-DQP, OGSA-DAI, and Globus Toolkit 3, building on top of them schema-integration services. All the nodes in the system expose their resources as Grid services (GSs) except data resources and data integration facility that are exposed as Grid Data Services (GDSs).

The effective use of a Grid requires the definition of a model to manage the heterogeneity of the involved resources. The management of such resources requires the use of metadata that, through an accurate categorization of resources, provide useful information about the features of resources and their effective use. In particular, metadata information is essential in the publishing and discovery of resources. The adoption of the service-oriented model in novel Grid systems (in particular the Open Grid Services Architecture, and the Web Services Resource Framework) will have an impact on the management of metadata and on the architecture of information services, since such systems allow to expose all services and resources as Grid services (also called WS-Resources in WSRF). The information model of service-oriented Grid frameworks is essentially based on two features: (i) metadata about Grid service instances is stored into XML-encoded documents; (ii) information is collected and indexed by means of hierarchical information services that subscribe to the information stored in Grid services, aggregate it and provide it to high level browsing and querying services.

Metadata are used to classify and manage a resource, but classification parameters, i.e. the structure of metadata information, depend on the type of the resource (i.e. software, hardware, data etc.) and on the application domain in which it is used. We designed both a metadata model and an information system that offer a uniform and at the same time flexible approach to the management of heterogeneous metadata structure. The proposed metadata model permits to classify and describe resources needed for different domains. A metadata document, associated to each resource, includes an ontological metadata section that identifies the category of the resource, a semantic metadata section that characterizes the resource and is used to assist discovery services, and a resource metadata section that gives details about the access mechanisms. The information system allows for a uniform and flexible management of metadata by exploiting an ontology system to semantically describe application domains and resources, and the basic information services of a service-oriented Grid framework, namely the Globus Toolkit 3, to aggregate and index metadata.

INFN-CNAF research teams are trying to find the best formalism inside Logic and Functional Programming Language as Prolog, LISP or Evolving Algebra, in order to develop an ontology-based information modeling for Grids. We are investigating a family of languages known as Description Logic that is a powerful paradigm for specifying conceptualization and classification. It offers reasoning support on concepts, defining axioms and deductive rules. This would possibly enable semantic matchmaking over Grid. Moreover, the SoA of Ontology Modeling over Grid is investigating a special binding between DL and a form of Semantic Web language known as OWL-DL. We address semantic representation by describing ontologies with Ontology Web Language (OWL). Our basis is OWL-S: an OWL based ontology to describe services. We also try to represent transactional support as a QoS feature in the ontology-based service description (OWL-S). We provide a taxonomy of transaction types, embody them in OWL-S as a QoS feature, and extend the constructs for service process definition with the notion of transaction. Our goal is an OWL-S extension that takes into account the provided taxonomy and includes notions of Grid/Web services integration.

### *Semantic Representation*

Designed to be accessed and navigated interactively by humans, web sites are often unstructured and hardly accessible automatically by computers. In that context, wrappers are software that extract structured and interpreted data from web content. For the Semantic Grid to become a reality, tools developed by CETIC help finding, indexing and semantically interpreting web pages. Most of the content of the web is unstructured. For the semantic Grid to become a reality some ways of accessing this data are required. We use a semantic browser to interactively structure unstructured web content with a user, thus allowing it to be accessed by structured content tools, such as semantic searching. The approach being followed is to build a wrapper to build a XML file based on meta-level information, and extract a schema used to validate it. This XML file is considered as the structured version of the unstructured content.

### *Semantic Agents*

The current implementation of information system in computing Grids include just a static model based on a tree distributed structure. Such a model cannot be used to effectively organize highly dynamic data and the relationships among Grid concepts. The languages used are not expressive enough both syntactically and semantically to model the complex relationships among concepts relevant to the Grid. A typical monitoring and Information System should provide the following three functionalities: raw data collection, analysis and management.

INFN-CNAF proposed a new architecture for a Grid inferential monitoring and discovery system. The main characteristic of the system is that it supports reasoning activities on a knowledge base (KB) that formalizes the concepts, and relationships relevant to a Grid computing environment. This KB models a new powerful Information Data Model for the Grid nodes and services. Moreover, our solution adds inferential capabilities, based on the DL paradigm, to all the components GIS (IP, GRIS, MDS). The solution focuses on the creation of an autonomous and adaptive management of Grid environment. This goal is achieved by a "conceptualization process" aimed to formally define each component of a Grid node environment. Moreover we want to use an agent component based for creating "assertion on concepts and individuals (TBOX and ABOX)". This agent has reasoning capabilities for enriching and querying the KB through forward and backward rules. A Service Oriented Architecture is defined as a set of network addressable services that can be invoked by other services. Services provide a well defined interface that can be published and discovered. At present, the most known instances of SOA are the Web Service Architecture (WSA) and the Open Grid Service Architecture (OGSA). In these environments services may need to dynamically discover non-functional properties of possible other services to cooperate with.

A topic of focus in the above mentioned loosely-coupled scenario is the automated matchmaking of services in a system composed by dynamic, heterogeneous and distributed resources. We propose an

algorithm for matchmaking and ranking of services based on the OWL ontologies. Our objective is to address the matchmaking of single services through investigating the notion of closeness between ontology instances, and use coordination aspects for matchmaking of multiple services. An algorithm addressing the matchmaking of a single service and possibly a workflow schedule of coordinated services. In many of the present Grid architectures the matchmaking phase is performed by entities called Resource Brokers (RBs). For scalability and efficiency reasons, each of these RB could be aware of only a fraction of the resources and services available on the Grid. In order to allow every single RB to satisfy the users' requests, we need a mechanism that exploits collaboration between RBs, enabling a RB to discover and match against resources that it did not previously know of. P2P systems and taxonomy and ontologies of resources and services, coupled with Service Oriented Architectures like the Open Grid Service Architecture (OGSA), seem to be the most promising technologies for this purpose.

#### *Standardization and Evolution*

The use of computers is changing our way to make discoveries and is improving both speed and quality of the discovery processes. Several scientific problems and industrial processes require the analysis of large and distributed repositories of data. In this scenario future Grids might provide an effective environment for distributed data mining and knowledge discovery from large data sets. To support users in solving that class of problems UNICAL recently designed a system for distributed data mining applications on Grids called Knowledge Grid. The Knowledge Grid architecture uses basic Grid mechanisms to build specific knowledge discovery services. These services can be developed in different ways using the available Grid environments. This approach benefits from "standard" Grid services that are more and more utilized and offers an open distributed knowledge discovery architecture that can be configured on top of Grid middleware in a simple way.

UNICAL is working on the design and composition of distributed knowledge-discovery services, according to the OGSA model, by using the Knowledge Grid environment starting from searching Grid resources, composing software and data elements, and executing the resulting application on a Grid. The expected result of this research activity is an implementation of the Knowledge Grid in terms of the OGSA model. In this implementation, each Knowledge Grid service, called K-Grid service, is exposed as a Web Service that exports one or more operations by using the WSRF conventions and mechanisms. The Knowledge Grid services can be used to orchestrate data mining applications on very large data sets available over Grids as workflow of services, to make scientific discoveries, improve industrial processes and organization models, and uncover business valuable information. Scientific activities and results of this research activities can find synergy and cooperation with research activities carried out in the DataMiningGrid project. Actually, a joint meeting has been planned in the June 2005 to discuss interaction and cooperation. The main common areas of research are in designing Grid services for knowledge discovery and workflow composition for distributed data mining.

### **4.3.3 Data Mining and Knowledge Discovery**

#### *Semantic Mapping*

The main area of related research at CCLRC is the definition of semantic description for service discovery and information management, emphasizing aspects of trust and security management. In Virtual Organisations, security management must become autonomic and adaptation must occur automatically in real-time, rather than through human intervention. Furthermore, autonomic security management will have to be complemented by extensible and machine processable standards for negotiating, validating and amending collaboration agreements, encoded by means of electronic contracts, which can be autonomically enacted by the platform. Such extensible and machine processable standards require the development of common vocabularies and negotiation protocols.

#### *Policy Publication and Enforcement*

Service providers will publish policies for their use, detailing the obligations, privileges and expected levels of service, which a user should accept before using the service. Some initial efforts in the use of Semantic Web representations for basic security applications (authentication, access control, data integrity, encryption) have begun to bear fruit. For example, Denker et al. [DEN03] have integrated a set of ontologies (credentials, security mechanisms) and security extensions for Web Service profiles with the CMU Semantic Matchmaker. Kagal et al. [KAG03] are also developing Rei, a Semantic Web based policy language. Furthermore, KAoS services and tools allow for the specification, management, conflict resolution, and enforcement of policies within the specific contexts established by complex organizational structures represented as domains [BRA03]. A comparison of KAoS, Rei, and more traditional policy approaches such as Ponder can be found in [TON03]. KAoS provides a powerful tool-set that appears to be capable to address publication and deployment of complex policies for Semantic Web Services. However the incorporation of trust metrics and a distributed enforcement and performance assessment schemes remain the main challenges, in addition to the production of a critical mass of domain/application-specific ontologies to allow its uptake and validation in large scale systems. With respect to the latter there is an ongoing effort to adapt KAoS for use in Grid Computing environments in conjunction to OGSA [JOH03].

### *Service Discovery*

In order for a new service to be used it needs to be discovered and a mapping needs to be established between the requirements of the client and the capabilities of the service. On the service side, discovery is facilitated in the presence of a set of semantic descriptions. In Web Service architectures, WSDL descriptions can be used to support this, but they fall short in providing any unambiguous semantic content for the service interface description they provide. Thus there have been approaches to describing the functionality of web services using Semantic Web technologies such OWL-S where in addition to publishing their interfaces, Web Services publish statements describing their intended or normative behaviour. These statements should be given common, machine processable, extensible semantics that support judgment of whether a service can perform a given task; the relative ranking of a set of services with respect to basic QoS criteria; and to then using reasoning to match service descriptions against requirements. On the client side, the client objectives must also be given the semantics in order to enable achieving a "sufficiently good" similarity between objectives of requestor and the capabilities of the service, advertised by its provider. Generally, a match can be determined by heuristic algorithms, aided by domain-specific ontologies that define the terms used for service description as well as the objectives of the requestor. Again, there is a need to extend this work to non-functional requirements. P3P [McB02] adds policies and requirement of the client with respect to Privacy; this would need to be extended to express the wider quality-of-service expectations of the client.

### *Service Negotiation*

Once a service has been selected, there needs to be a negotiation between service and user to establish a relationship. As part of this process, the policies of both parties have to be interrogated and a contract of use established, and a conversation needs to take place between the parties, establishing a mutually intelligible vocabulary of terms for data and process descriptions. This negotiation may involve third parties (brokers, guarantors, service framework providers etc), which may facilitate the relationship and foster trust between the parties. In this process, there is a step of trust evaluation, either from previous experience of one another, as recorded in a "trustbase" of trust valuations, or an evaluation of the trust value from recommendations from third parties, or a calculation of trust across the network via intermediate trust valuations. Preliminary work in calculating trust values across trust networks in the semantic web have been studied by Goldbeck, Hendler and Parsia [GOL03] and Richardson, Agrawal, and Domingos [RIC03] which use a relatively straightforward model of trust which does not take into account context or uncertainty.

### *Monitoring and Policy Enforcement*

During the execution of the service, which may be over a long period, its progress is monitored. The experience of the quality of the service may modify the relationship between the parties. For example, if the experience so far is good, then the parties may relax restrictions for the remainder of the service. Policy statements need to be interpreted into lower-level rules which are then enforced at each network end-point. Web Services standards for SOAP-based message security and XML-based languages for access control (e.g. XACML) are emerging. The use of XML as a basis for expression specification has the advantage of extensibility. Its semantics however are mostly implicit as meaning depends on a shared understanding derived from human consensus, and allow incompatible representation variations. Semantic Web-based policy representations could be mapped to lower level XML representations if required by an implementation. Once an agreement has been established, then the client can start using the service. This usage may be long-lived, and the experience of the parties during the interaction may modify their behaviour for its remainder. For example, good experience may result in the loosening of restrictions and a higher-level of trust, changing the valuations in internal “trustbases”, and reducing the policy enforcement overhead.

#### *Intelligent Queries*

OGSA-DAI components are either data-access components or data-integration components. A Distributed Query Processing (DQP) system is an example of a data integration component and can potentially provide effective declarative support for service orchestration as well as data integration. The service-based DQP framework described in [AMP+03], termed as OGSA-DQP, provides an approach that:

- supports queries over GDSs (Grid Data Services) and over other services available on the Grid, thereby combining data access with analysis;
- uses the facilities of the OGSA to dynamically obtain the resources necessary for efficient evaluation of a distributed query;
- adapts techniques from parallel databases to provide implicit parallelism for complex data-intensive requests; and
- uses the emerging standard for GDSs to provide consistent access to database metadata and to interact with databases on the Grid.
- The service-based DQP framework extends the OGSA-DAI with two new services (and their corresponding factories):
- Grid Distributed Query Service (GDQS). The GDQS is the main interaction point for the clients. When a GDQS is set up, it interacts with the appropriate registries to obtain the metadata and computational resource information that it needs to compile, optimize, partition and schedule distributed query execution plans over multiple execution nodes in the Grid. The implementation of the GDQS builds on a previous work on the Polar\* distributed query processor for the Grid [SGW+02] by encapsulating its compilation and optimisation functionality.
- Grid Query Evaluation Service (GQES). GQES instances are created by the GDQS based on the query plans generated by the query compiler, optimiser and scheduler. Each GQES instance evaluates a partition of the query execution plan assigned to it by a GDQS. The set of GQES instances form a tree through which the data flows from leaf GQESs which interact with GDSs, up the tree to reach its destination.

As well as using the services provided by Grid Data Services (GDSs), the GDQS and GQES both implement GDS port type, and thus can be discovered and invoked in the same way as other GDSs. Consequently, the Grid stands to benefit from DQP, through the provision of facilities for declarative request formulation that complement existing approaches to service orchestration, via uniform interfaces and interaction semantics. Expected results: OGSA-DQP can be used as a generic query processor over commercial databases connected to the Grid. In particular, it can be used in the context of the Information and Knowledge Management task for supporting data-intensive applications on Grids, and in the context of the Data Mining and Knowledge Discovery - Intelligent Queries task to provide a query mechanism.

UoM is currently working on the construction of a generic framework for describing and constructing AQP systems in a systematic way, allowing for component reuse. AQP thus far suffers from two major shortcomings:

- current techniques are designed in an isolated way, which does not permit them to be combined; and,
- most AQP proposals have focused either on completely centralised query processing, or on centralised processing of data retrieved from remote sources and data streams. By following such an approach, the resources used for query execution are predefined, and thus the focus is mostly on adapting to changing properties of the data processed (such as cardinalities of intermediate results and operator selectivities). This is of paramount importance for query processing on the Grid as crucial information about the data may be missing at compile time. However, of equal significance are adaptations to changing properties of the arbitrary set of resources that query processing on the Grid uses both for data retrieval and other types of data manipulation, such as joins. Currently, AQP with respect to changing resources is not addressed as satisfactorily, as with respect to changing data properties.

The framework being developed tackles the first of the two above limitations. It is based on the decomposition into and the separate investigation of three distinct phases that are inherently present in any AQP system:

- monitoring of query execution and environment to collect the necessary feedback;
- assessment of the feedback collected to establish issues with the current query execution or opportunities for improvement; and
- responding to the identified monitoring events based upon the results of the assessment process.

The framework has been instantiated in the context of OGSA-DQP to tackle the second limitation of existing AQP techniques identified previously, i.e., adapting to changing resources; and a prototype has been produced. In particular, two important cases are being investigated:

- adaptive workload balancing of parallel query processing; and
- adaptive resource scheduling.

Future work includes the application of the framework in other domains except query processing, e.g., workflows, data mining, etc. Expected results: As Grid is an inherently volatile environment, adaptivity is likely to be crucial for executing tasks on the Grid in an efficient way. Thus, it is felt that the adaptive extensions to OGSA-DQP will improve the performance of any system or toll that uses DQP significantly.

### *Knowledge Grid Services*

Knowledge discovery processes and data mining applications generally need to handle large data sets and, at the same time, are compute intensive tasks that in many cases involve distribution of data and computations. This is a scenario where the use of parallel and distributed computers can be effective for solving complex data mining tasks. The Knowledge Grid system is an environment for distributed data mining on Grids that was developed at UNICAL. A Knowledge Grid prototype has been used for implementing data mining applications on Globus-based Grids. Currently under development are the system mechanisms and operations as Grid Services. In the Knowledge Grid framework, data mining tasks and knowledge discovery processes will be made available as OGSA-WSRF services that will export data and tool access, data and knowledge transmission, and mining services. These services will make possible the design and orchestration of distributed data mining applications running on large-scale Grids. This approach will support the integration of algorithms, tools, and data sources for implementing knowledge discovery application on Grids.

CNR-ISTI has developed state-of-the-art data mining algorithms to extract frequent patterns from transactional databases. Frequent pattern extraction is the most computational expensive phase of Association Rule Mining (ARM), one of the most popular topics in the KDD field. The main contribution of these algorithms, that makes them suitable for Grids, is their adaptivity to the features

of the executing platform, i.e. the possibility of running on clusters of SMPs when distributed data can be moved to a centralized site, or using a loosely-coupled distributed approximate approach when data cannot be moved and the final results must be inferred by merging partial results/models independently computed by each distributed site.

They have developed DCI (Direct Count & Intersect), a scalable algorithm for discovering frequent itemsets in large databases. The main contribution of DCI relies on the multiple heuristics strategies employed. DCI adapts its behaviour to the features of the specific dataset mined and of the computing platform used as well, so that it turns out to be effective in mining both short and long patterns from sparse and dense datasets. We also developed ParDCI, a parallel implementation of DCI. ParDCI explicitly targets clusters of SMP nodes, so that shared memory and message passing paradigms are exploited at the intra- and inter-node levels, respectively.

One of the main problems addressed by ParDCI is load imbalance. When data cannot be moved to be centralized on a single site, and network communications are slow, a more distributed approach must be adopted. They have thus, developed a novel communication-efficient distributed algorithm, AP<sub>Interp</sub>, for approximate the mining of frequent itemsets. The proposed algorithm consists of the distributed exact computation of locally frequent itemsets, and an effective method for merging local results and inferring the local support of locally infrequent itemsets. The combination of the two give a good approximation of the set of the globally frequent patterns and their supports.

We plan to focus our future activities on devising Grid-aware distributed/approximate algorithms to extract frequent patterns and/or condensed lossless representation of them. Moreover, we want to consider the streamed nature of data sources, which can also referenced in space and time like trajectories collected from mobile devices. This research activity is continuing within GeoPKDD, a national project. Ref: [OPP02, OPP02a, SO04, SO05]

#### *Monitoring and predicting performance of data-mining tools*

In order to optimize Grid resources and reduce the response time of submitted tasks, we need to predict complexity and performance of tasks. Unfortunately, most of the complexity of common data mining tasks is due to the unknown amount of information contained in the data being mined. The more patterns and correlations are contained in data sources, the more resources are needed to extract them. This is confirmed by the fact that, in general, there is not a single best algorithm for a given data mining task on any possible kind of input dataset. On the contrary, in order to achieve good performances, strategies and optimizations have to be adopted according to the dataset specific characteristics. For example, one typical distinction in transactional databases is between sparse and dense datasets. The Frequent Itemset Mining (FIM) problem, used as a case study for performance prediction, is usual considered very expensive when used to extract knowledge from dense datasets. Moreover, the user parameters supplied by analysts, like the minimum support threshold, strongly influence resource usage and make performance prediction difficult. Small support thresholds can cause a combinatorial explosion in the complexity of the problem. We have thus proposed a statistical analysis of the properties of transactional datasets that allows us to give a characterization of the complexity of the mining process, a very useful feature when DM tools have to be exploited in a Grid environment, where resource usage must be optimized. The prediction requires to monitor and profile the execution of a FIM tool, when executed to extract pattern from a dataset given a user supplied support threshold. We will continue this research activity in the EU NextGrid project, during which we aim to build a model trained with a multitude of statistical features, extracted from a dataset before starting the actual mining process [POP04].

#### *Scheduling high-performance data-mining tasks on data Grids*

Increasingly the datasets used for data mining are huge and physically distributed. Since the distributed knowledge discovery process is both data and computational intensive, the Grid is a natural platform for deploying a high performance data mining service. The focus of this activity is on the

core services of such a Grid infrastructure. In particular we concentrate on the design and implementation of specialized Resource Allocation and Execution Management services, aware of the location of data sources and of the resource needs of the data mining tasks. Scheduling decisions are taken on the basis of cost metrics and models that exploit knowledge about previous executions, and use sampling to acquire estimate about execution behavior. We explored the limit of such performance estimation, but also put in evidence its benefits in reducing the whole completion time of all tasks. In particular, we explored the use of an on-line MCT (Minimum Completion Time) heuristic strategy for scheduling high performance DM tasks on a Knowledge Grid [OPP02b].

#### 4.3.4 Phases of the roadmap

We identify three main phases for this roadmap that should drive the activities of the partners in implementing the WP2 activities:

- **Phase 1:** Exchanging partner information, experiences, and knowledge about techniques, tools and systems for Data and Knowledge Grids.
- **Phase 2:** Sharing and integration of common goals, research results, pursued projects and system prototypes of Environments and Services for Data and Knowledge-based Grids.
- **Phase 3:** Use of the results of the previous phases for envisioning a unified framework for handling data, information and knowledge on Grids. Definition of joint proposals, research activities and projects in the WP2 research area.

#### 4.3.5 Mechanisms

- *Short visits:* The partners planned a series of short visits with duration from 2 days to 2 weeks that will be used to exchange research experiences and set up collaborations on common research activities.
- *Workshops:* Workshops are periodically organized to exchange knowledge achieved by the Workpackage partners through presentation of results and foster collaboration on common research goals. A first workshop was held on the 14<sup>th</sup> and 15<sup>th</sup> of January in Iraklion (Crete) to present research activities and organizes the contributions to the first deliverable.
- *Documents:* Sharing of documents, reports and other material among the partners also by using the CoreGRID portal.
- *E-meetings & tele-conference meetings*
- *New projects:* Definition of provisional new projects that can involve partners involved in the Workpackage.
- *Invite external people and involve other institutions:* Invitation of external researchers and institutions active in the areas of the Workpackage.
- *WP2 member participation in committees and forums:* Participation in standards activities or Grid community where topics dealt in the Workpackage are discussed
- *Exchange of tools, data, logs, and metadata:* Sharing of software tools, prototypes, data sources, metadata and other digital content among the partners also by using the CoreGRID portal.
- *Common testbed:* if need arises and is possible.

#### 4.3.6 Future steps

Basically partners want to foster collaboration in common research activities, integration of research teams and research results. Here are listed the main future steps that will be done to implement the Workpackage roadmap:

- Promote partner collaboration to fill gaps in the ‘bigger’ picture.
- Promote participation of partners into new research proposals, assuming the results of this WP as starting point.
- Ensure stable and durable cooperation between WP partners.
- Try to establish a common testbed for future work.
- Participation in standardization activities.

- Result dissemination to industry.

**Table 1: Partner contributions in joint activities.**

INFN	Visit to UNICAL
UNICAL	Visit to UoM
UCY	Visit to UNICAL
FORTH	Visit to UCY
FORTH	Visit to PSNC
UCY	Visit to FORTH
UCY	Visit to UoM
CETIC	Visit to CCLRC
UOM	Visit to FORTH
UOM	Visit to UCY
UOM	Visit to UNICAL

## 5 Link with other CoreGRID scientific workpackages

The activities performed by the partners involved in the WP2 are related to the workpackage goals and tasks and at the same time related to the activities of the other CoreGRID workpackages. The following figure shows the main interactions between WP2 and the other WPs.

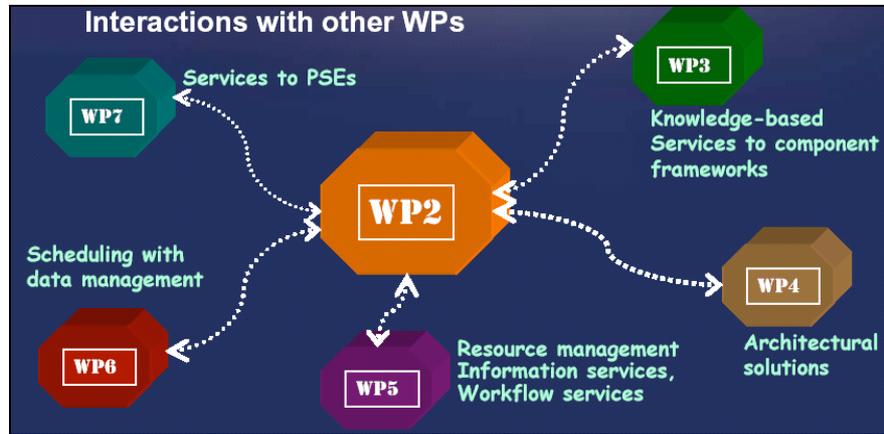


Figure 3: WP relationships

WP3 aims at defining a common Component Model for the Grid, including the notion of parallelism and distribution. The interaction between WP2 and WP3 should be concerned with investigating how the component programming models and frameworks proposed in WP3 could be used to make it easier to design and deploy parallel and distributed Grid-aware data mining tools.

Interactions between WP2 and WP4 concern the use of peer-to-peer techniques, protocols and architectures for implementing scalable data and knowledge-based applications on Grids. P2P solutions can be used both to implement non-hierarchical services for data and knowledge-based systems and to implement data-intensive systems and applications on Grids such as P2P DBMS, P2P data integration systems, and p2p data mining on Grids.

Topics and issues that link WP2 and WP5 are resource management, information services and workflow services. In particular, WP5 can provide solutions concerning resource management and workflow services that may be used in research activities and results of WP2 partners. At the same time, WP2 solutions in the area of data and knowledge management can be exploited in systems developed in WP5.

Data management and movement can be a critical part of a Grid job in terms of performance. Job scheduling in Grids should take into account data movement by properly considering the execution of data movement tasks in the evaluation of job execution times, and adopting estimates of intermediate result sizes and dynamic network performance. At the same time, data management should be an integrated component of a Grid scheduling architecture and of Grid scheduling algorithms. With respect to WP6, UoM is interested in contributing to the area of resource scheduling for distributed query processing on the Grid. Those topics are common areas of interaction between WP2 and WP6.

A main objective of WP7 is the development of a component-based infrastructure for integrating applications, tools and system facilities. This involves the design of application-level metadata storage and caching services to facilitate the scheduling tasks in an efficient way. In that direction, WP2 could contribute its data management and integration expertise for building the metadata repository and cache, while WP2 could benefit from scheduling information made available from WP7 towards optimizing the distributed data management activities.

## 6 References

- [AGC05] S. V. Anastasiadis, S. Gadde, J. S. Chase. Scale and Performance in Semantic Storage Management of Data Grids, *International Journal on Digital Libraries*. 2005 (to appear).
- [AHK02] E. Anderson, M. Hobbs, K. Keeton, S. Spence, M. Uysal, A. Veitch. Hippodrome: running around storage administration. Usenix FAST Conference, Monterey, CA, January 2002.
- [AMP+03] M. N. Alpdemir, A. Mukherjee, N.W. Paton, P. Watson, A. A. Fernandes, A. Gounaris, and J. Smith. Service-based distributed querying on the grid. In the Proceedings of the First International Conference on Service Oriented Computing, pages 467-482. Springer, 15-18 December 2003.
- [ANR04] Knowledge Base for Adaptive Decision Making in Autonomous Grid Monitoring Middleware. Ashiq Anjum, Fawad Nazir, Nasir Rasul, Arshad Ali. European Organization for Nuclear Research (CERN) Geneva. In Posters Session of Middleware2004. Toronto, Canada. Oct 2004.
- [ASV03] S. Andreozzi, M. Sgaravatto and C. Vistoli. Sharing a conceptual model of Grid resources and services. In *Proceedings of the 2003 Conference for Computing in High Energy and Nuclear Physics*, March 2003. <http://www.slac.stanford.edu/econf/C0303241/>.
- [BBH+ 02] W. H. Bell, D. Bosio, W. Hoschek, P. Kunszt, G. McCance, and M. Silander. Project Spitfire - Towards Grid Web Service Databases. In *Global Grid Forum 5*, 2002.
- [BER01] F. Berman, "From TeraGrid to Knowledge Grid," *Comm. ACM*, vol. 44, no. 11, Nov. 2001, pp. 27–28.
- [BFGC04] J. Brooke, D. Fellows, K. Garwood and C. Coble. Semantic matching of Grid Resource Descriptions. In M.D. Dikaiakos, editor, *Grid Computing. Second European Across Grids Conference, AxGrids 2004, Nicosia, Cyprus, January 2004, Revised Papers*, volume 3165 of *Lecture Notes in Computer Science*, pages 240-249. Springer, 2004.
- [BFK+00] M. Beynon, R. Ferreira, T. Kurc, A. Sussman, and J. Saltz. Datacutter: Middleware for filtering very large scientific datasets on archival storage systems. In *IEEE Symposium on Mass Storage Systems*, pages 119–134, 2000.
- [BRA03] Bradshaw, J., Uszok, A., Jeffers, R., Suri, N., Hayes, P., Burstein, M., Acquisti, A., Benyo, B., Breedy, M., Carvalho, M., Diller, D., Johnson, M., Kulkarni, S., Lott, J., Sierhuis, M. and Van Hoof, R. Representation and reasoning about DAML-based policy and domain services in KAoS. In *Proc. of The 2nd Int. Joint Conf. on Autonomous Agents and Multi Agent Systems (AAMAS2003)*.
- [CAT03] M. Cannataro, D. Talia, "The Knowledge Grid", *Communications of the ACM*, vol. 46, no. 1, pp. 89-93, 2003.
- [CNT+04] Computer Network Technology Corporation. UltraNet® Replication Appliance for IBM. Product description. <http://www.cnt.com/documents/?ext=pdf&filename=PL821>.
- [COT04] C. Comito, D. Talia, "GDIS: A Service-based Architecture for Data Integration on Grids", *Proc. Conference OTM 2004, Cyprus, 2004*, Springer Verlag, LNCS, pp. 88-98, 2004.
- [CS+04] Cisco Systems, Inc. White Paper. EMC and Cisco: Building disaster recovery and business continuity solutions. [http://www.cisco.com/warp/public/cc/so/neso/datactr/emcwp\\_wp.pdf](http://www.cisco.com/warp/public/cc/so/neso/datactr/emcwp_wp.pdf).
- [DEN03] Denker, G., Kagal, L., Finin, T., Paolucci, M. and Sycara, K. Security for DAML Web Services: Annotation and Matchmaking. In D. Fensel, K. Sycara, & J. Mylopoulos (Ed.), *The Semantic Web—ISWC 2003. Proceedings of the 2nd International Semantic Web Conference*, Sanibel Island, Florida, USA, October 2003, LNCS 2870.
- [DMG05] <http://www.datamininggrid.org>.
- [DSI05] M. Dikaiakos, R. Sakellariou, Y. Ioannidis, "Information Services for Large-Scale Grids: A Case for a Grid Search Engine." In *Engineering the Grid: status and perspective*, Jack Dongarra, Hans Zima, Adolfo Hoisie, Laurence Yang, Beniamino DiMartino (Editors), Nova publishers (to appear, 2005).
- [DTMF03] Distributed Management Task Force. CIM Concepts White Paper. CIM Versions 2.4+. June 2003. <http://www.dmtf.org/standards/documents/CIM/DSP0110.pdf> (accessed Oct. 2004).

- [EGE+04] EGEE. Enabling Grids for E-science in Europe. Fact Sheet. <http://public.eu-egee.org/files/EGEEfact-sheet2.pdf>.
- [GEL04] Geldof, M. The Semantic Grid: Will Semantic Web and Grid go hand in hand? European Commission, DG Information Society Unit “Grid Technologies”, 2004. Available at <http://www.semanticgrid.org/documents/Semantic%20Grid%20report%20public.pdf>.
- [GLUE04] Glue Schema Official Documents, <http://www.cnaf.infn.it/~sergio/datatag/glue> (last accessed Sept. 2004).
- [GOL03] Golbeck J., Parsia B., and Hendler J.: Trust networks on the semantic web. In Proceedings of Cooperative Intelligent Agents 2003, Helsinki, Finland, August 2003.
- [GRA99] J. Gray. What Next? A Few Remaining Problems in Information Technology (Turing Lecture). ACM Federated Computer Research Conferences (FCRC). May 1999. Atlanta, Georgia
- [GRIP04] Grid Interoperability Project. <http://www.grid-interoperability.org> (last accessed Jan. 2005).
- [GTH00] K. Govil, D. Teodosiu, Y. Huang, M. Rosenblum. Cellular disco: resource management using virtual clusters on shared-memory multiprocessors. ACM Trans. Comput. Syst. 18(3). 229-262 (2000).
- [HDL+03] H. Stockinger, F. Donno, E. Laure, S. Muzaffar, P. Kunszt, Grid Data Management in Action: Experience in Running and Supporting Data Management Services in the EU DataGrid Project. 2003 Conference for Computing in High-Energy and Nuclear Physics (CHEP 03), La Jolla, California, 24-28 Mar 2003.
- [HLN03] J. Hau, W. Lee and S. Newhouse. Autonomic Service Adaptation in ICENI using Ontological Annotation. In *Proceedings of the 4<sup>th</sup> International Workshop on Grid Computing*, pages 10-17. IEEE Computer Society, April 2003.
- [HMS+00] W. Hoschek, JJ. Martinez, A. Samar, H. Stockinger, K. Stockinger. Data Management in an International Data Grid Project, EEE/ACM International Workshop on Grid Computing Grid'2000 - 17-20 December 2000 Bangalore.
- [HST00] I. Horrocks, U. Sattler, and S. Tobies. Reasoning with individuals for the description logic SHIQ. In David MacAllester, editor, Proc. of the 17th Int. Conf. on Automated Deduction (CADE 2000), number 1831 in Lecture Notes in Artificial Intelligence, pages 482-496. Springer-Verlag, 2000.
- [JOH03] Johnson, M., Chang, P., Jeffers, R., Bradshaw, J. M., Soo, V.-W., Breedy, M. R., Bunch, L., Kulkarni, S., Lott, J., Suri, N., & Uszok, A. KAoS semantic policy and domain services: An application of DAML to Web services-based grid architectures. Proceedings of the AAMAS 03 Workshop on Web Services and Agent-Based Engineering. Melbourne, Australia, 2003.
- [KAG03] Kagal, K., Finin, T. and Anupam, J. A Logical Policy Language for a Pervasive Computing Environment., 4th IEEE Int. Workshop on Policies for Distributed Systems and Networks, Lake Como, 4-6 June, 2003.
- [KC00] H. Kargupta, P. Chan (Eds.). *Advances in Distributed and Parallel Knowledge Discovery*. AAAI Press/The MIT Press, 2000.
- [LFP03] David T. Liu, Michael J. Franklin, and Devesh Parekh. GridDB: a relational interface for the grid. In Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data 2003, pages 660–660. ACM Press, 2003.
- [LVA03] Peter Lyman and Hal R. Varian. *How Much Storage is Enough?* ACM Queue vol. 1, no. 4 - June 2003. <http://www.acmqueue.org/modules.php?name=Content&pa=showpage&pid=45>
- [LVS03] Peter Lyman, Hal R. Varian, Kirsten Swearingen, Peter Charles, Nathan Good, Laheem Lamar Jordan, Joyojeet Pal. How Much Information? 2003. <http://www.sims.berkeley.edu/research/projects/how-much-info-2003>.
- [McB02] McBride, B., Wenning, R. and Cranor, L.: An RDF Schema for P3P. W3C Note 25 January 2002: <http://www.w3.org/TR/p3p-rdfschema>.
- [MM02] G. S. Manku, R. Motwani. “Approximate Frequency Counts over Data Streams”. VLDB 2002.
- [MPPDM04] S. Miles, J. Papay, K. Decker, and L. Moreau. Towards a Protocol for the Attachment of Semantic Descriptions to Grid Services. In M. D. Dikaiakos, editor, *Grid Computing. Second European Across Grids Conference, AxGrids 2004, Nicosia, Cyprus, January*

- 2004, *Revised Papers*, volume 3165 of *Lecture Notes in Computer Science*, pages 240-249, Springer, 2004.
- [MPWL04] S. Miles, J. Papay, C. Wroe, P. Lord, C. Goble, and L. Moreau. Semantic Description, Publication and Discovery of Workflows in myGrid. Technical Report ECSTR-IAM04-001, Electronics and Computer Science, University of Southampton, 2004.
  - [MSBT03] Tanu Malik, Alex S. Szalay, Tamas Budavari, and Ani R. Thakar. SkyQuery: A Web Service Approach to Federate Databases. In Proc. CIDR, 2003.
  - [myGrid04] myGrid UK e-Science Project. <http://www.myGrid.org> (accessed Nov. 2004).
  - [NCK+03] Sivaramakrishnan Narayanan, Umit V. Catalyrek, , Tahsin M. Kurc, Xi Zhang, and Joel H. Saltz. Applying database support for large scale data driven science in distributed environemnts. In Proc. of the 4th Workshop on Grid Computing, GRID'03, 2003.
  - [OPP02] S. Orlando, P. Palmerini, R. Perego, F. Silvestri. "Adaptive and Resource-Aware Mining of Frequent Sets". IEEE ICDM, pp. 338-345, 2002.
  - [OPP02a] S. Orlando, P. Palmerini, R. Perego, F. Silvestri. "An Efficient Parallel and Distributed Algorithm for Counting Frequent Sets". Int. Conf. VECPAR 2002, 2002. (Appeared as Selected Paper in LNCS 2565, Springer, pp. 197-204).
  - [OPP02b] S. Orlando, P. Palmerini, R. Perego, F. Silvestri. "Scheduling High Performance Data Mining Tasks on a Data Grid Environment". Euro-Par 2002, LNCS 2400, Springer, pp. 375-384, 2002.
  - [PHV03] Perelman, E.; Hamerly, G.; Van Biesbrouck, M.; Sherwood, T., and Calder, B., Using SimPoint for Accurate and Efficient Simulation, *Proceedings of the ACM SIGMETRICS the International Conference on Measurement and Modeling of Computer Systems*, June 2003.
  - [PK02] B. Park, H. Kargupta. "Distributed Data Mining: Algorithms, Systems, and Applications". In *Data Mining Handbook*, IEA, pp. 341-358, 2002.
  - [POP04] P. Palmerini, S. Orlando, R. Perego. "On Statistical Properties of Transactional Datasets". ACM - SAC 2004, Special track on Data Mining, pp. 515-519, 2004.
  - [REG03] S. Rhea, P. Eaton, D. Geels, H. Weatherspoon, B. Zhao, J. Kubiatiowicz. Pond: The OceanStore Prototype. Usenix FAST Conference, San Francisco, CA, March 2003
  - [RIC03] Richardson, M., Agrawal, R. and Domingos, P. Trust Management and the Semantic Web In D. Fensel, K. Sycara, & J. Mylopoulos (Eds.), *The Semantic Web—ISWC 2003*. Proc. of the 2nd Int. Semantic Web Conf., Sanibel Island, Florida, USA, October 2003, LNCS 2870.
  - [RWM+03] A. Rajasekar, M. Wan, R. Moore, W. Schroeder, G. Kremenek, A. Jagatheesan, C. Cowart, B. Zhu, S-Y. Chen, R. Olschanowsky, *Computer Society of India Journal, Special Issue on SAN*, Vol. 33, No. 4, pp. 42-54 Oct 2003.
  - [SGW+02] J. Smith, A. Gounaris, P. Watson, N. W. Paton, A. A. A. Fernandes, and R. Sakellariou. Distributed Query Processing on the Grid. In Proc. Grid Computing 2002, pages 279-290. Springer, LNCS 2536, 2002.
  - [SO04] C. Silvestri, S. Orlando. "A new algorithm for gap constrained sequence mining". ACM - SAC 2004, Special track on Data Mining, 2004.
  - [SO05] C. Silvestri, S. Orlando. "Distributed Approximate Mining of Frequent Patterns". ACM - SAC 2005, Special track on Data Mining, 2005.
  - [SPH03] Sherwood, T.; Perelman, E.; Hamerly, G.; Sair, S., and Calder, B., Discovering and Exploiting Program Phases, *IEEE Micro: Micro's Top Picks from Computer Architecture Conferences*, December 2003.
  - [TDK03] H. Tangmunarunkit, S. Decker, and C. Kesselman. Ontology-Based Resource Matching in the Grid – The Grid Meets the Semantic Web. In D. Fensel, K.P. Sycara, and J. Mylopoulos, editors, *The Semantic Web –ISWC 2003, Second International Semantic Web Conference, Sanibel Island, FL, USA, Oct. 20-23, 2003, Proceedings*, volume 2870 of *Lecture Notes in Computer Science*, pages 706-721. 2003.
  - [TON03] Tonti, G., Bradshaw, J. M., Jeffers, R., Montanari, R., Suri, N., & Uszok, A. (2003). Semantic Web languages for policy representation and reasoning: A comparison of KAoS, Rei, and Ponder. In D. Fensel, K. Sycara, & J. Mylopoulos (Eds.), *The Semantic Web—ISWC 2003*. Proc. of the 2nd Int. Semantic Web Conf., Sanibel Island, Florida, USA, October 2003, LNCS 2870

- [TPC03] Transaction Processing Performance Council. TPC Benchmark H (Decision Support), Standard Specification, Revision 2.1.0, August 2003.
- [VD04a] Vandierendonck, H.; De Bosschere, K. Experiments with Subsetting Benchmark Suites. *Proceedings of the Seventh Annual IEEE International Workshop on Workload Characterization*. 2004.
- [VD04b] Vandierendonck, H.; De Bosschere, K. Many Benchmarks Stress the Same Bottlenecks. *Workshop on Computer Architecture Evaluation Using Commercial Workloads*. 2004.
- [W02] O. Wolfson. "Moving Objects Information Management: The Database Challenge". 5th Workshop on Next Generation Information Technologies and Systems (NGITS'2002), Israel, 2002.
- [YAN03] The Yankee Group. DSL to Dominate Europe's Broadband Growth. [http://www.yankeegroup.com/-public/products/research\\\_note.jsp?ID=9523](http://www.yankeegroup.com/-public/products/research\_note.jsp?ID=9523). February 2003.

## 7 Participants

The partners that are involved in WP2 and have contributed to this roadmap are:

<b>Task</b>	<b>Name</b>	<b>Involved Partners</b>
<b>T2.1</b>	<b>Distributed Data Management</b>	<b><u>FORTH</u>, CETIC, PSNC, UCY</b>
<b>T2.2</b>	<b>Resource Discovery and Data Management</b>	<b><u>CCLRC</u>, CETIC, INFN, UCY, UNICAL, UoM</b>
<b>T2.3</b>	<b>Data Mining and Knowledge Discovery</b>	<b><u>UNICAL</u>, CETIC, CNR-ISTI, RAL, UCY, UoM</b>

The individual researchers that have participated in preparing the roadmap are:

Alvaro Arenas	CCLRC
Anne Falien	CETIC
Antonia Ghiselli	INFN
Angelos Bilas	FORTH
Carmela Comito	UNICAL
Christophe Noel	CETIC
Christiana Christophi	UCY
Pedro Trancoso	UCY
Michail Flouris	FORTH
Stergios Anastasiadis	FORTH
Donatien Grolaux	CETIC
Edgardo Ambrosi	INFN
Fabrice Estievenart	CETIC
Francesco Bonchi	ISTI
Giorgos Tsouloupas	UCY
Tasos Gounaris	UoM
Laura Bocchi	INFN
Maciej Brzezniak	PSNC
Matteo Mordacchini	INFN
Mario Dikaiakos	UCY
Norbert Meyer	PSNC
Salvatore Orlando	ISTI
Philippe Massonet	CETIC
Raffaele Perego	ISTI
Rizos Sakellariou	UoM
Talia Domenico	UNICAL
Wei Xing	UCY