



IN PARTNERSHIP WITH:
CNRS

**Université des sciences et
techniques du Languedoc
(Montpellier 2)**

Activity Report 2012

Project-Team ZENITH

Scientific Data Management

IN COLLABORATION WITH: Laboratoire d'informatique, de robotique et de microélectronique de Montpellier (LIRMM)

RESEARCH CENTER
Sophia Antipolis - Méditerranée

THEME
**Knowledge and Data Representation
and Management**

Table of contents

1. Members	1
2. Overall Objectives	2
2.1. Introduction	2
2.2. Highlights of the Year	3
3. Scientific Foundations	3
3.1. Data Management	3
3.2. Distributed Data Management	3
3.3. Cloud Data Management	5
3.4. Uncertain Data Management	6
3.5. Metadata Integration	6
3.6. Data Mining	7
3.7. Content-based Information Retrieval	8
4. Application Domains	9
5. Software	10
5.1. WebSmatch (Web Schema Matching)	10
5.2. YAM++ ((not) Yet Another Matcher)	10
5.3. SON (Shared-data Overlay Network)	10
5.4. P2Prec (P2P recommendation service)	11
5.5. ProbDB (Probabilistic Database)	11
5.6. Pl@ntNet-Identify	11
5.7. Pl@ntNet-DataManager	11
5.8. SnoopIm	12
6. New Results	12
6.1. Data and Metadata Management	12
6.1.1. Uncertain Data Management	12
6.1.2. Metadata Integration	13
6.1.3. High-dimensional data management	13
6.2. Data and Process Sharing	14
6.2.1. Hybrid P2P/cloud Architecture	14
6.2.2. Social-based P2P Data Sharing	14
6.2.3. View Selection in Distributed Data Warehousing	15
6.2.4. Scientific Workflow Management	15
6.2.5. Plants identification and classification from social image data	15
6.3. Scalable Data Analysis	16
6.3.1. StreamCloud	16
6.3.2. Mining Uncertain Data Streams	16
6.3.3. Detecting Rare Events in Massive Datasets	17
6.3.4. Highly Informative Feature Set Mining	17
6.3.5. Clustering Users with Evolving Profiles in Usage Streams	17
6.3.6. Scalable Mining of Small Visual Objects	17
7. Bilateral Contracts and Grants with Industry	18
8. Partnerships and Cooperations	18
8.1. Regional Initiatives	18
8.1.1. Labex NUMEV, Montpellier	18
8.1.2. Institut de Biologie Computationnelle (IBC), Montpellier	18
8.1.3. ModSiCS2020 Working Group, Montpellier	19
8.2. National Initiatives	19
8.2.1. ANR	19
8.2.1.1. VERSO DataRing(2008-2012, 300Keuros)	19

8.2.1.2.	OTMedia (2011-2013), 150Keuros	19
8.2.2.	Others	19
8.2.2.1.	RTRA PI@ntNet (2009-2013), 1Meuros	19
8.2.2.2.	CIFRE INA/Inria (2011-2013), 100Keuros	20
8.2.2.3.	CNRS INS2I Mastodons (2012), 30Keuros	20
8.3.	European Initiatives	20
8.4.	International Initiatives	20
8.4.1.	Inria International Partners	20
8.4.2.	Participation In International Programs	21
8.5.	International Research Visitors	21
8.5.1.	Visits of International Scientists	21
8.5.2.	Visits to International Teams	21
9.	Dissemination	21
9.1.	Scientific Animation	21
9.2.	Teaching - Supervision - Juries	23
9.2.1.	Teaching	23
9.2.2.	Supervision	24
9.2.3.	Juries	24
9.3.	Popularization	24
10.	Bibliography	25

Project-Team ZENITH

Keywords: Data Management, Scientific Data, Information Indexing And Retrieval, Workflow, Parallelism

Zenith is a joint team with University Montpellier 2 (UM2) and is located at LIRMM (CNRS and UM2), Montpellier.

Creation of the Project-Team: January 01, 2011 , Updated into Project-Team: January 01, 2012 .

1. Members

Research Scientists

Reza Akbarinia [Junior Researcher, Inria]
Alexis Joly [Junior Researcher, Inria]
Florent Masegla [Junior Researcher, Inria]
Didier Parigot [Junior Researcher, Inria]
Patrick Valduriez [Team Leader, Senior Researcher, Inria, HdR]

Faculty Members

Zohra Bellahsene [Professor, UM2, HdR]
Hinde Bouziane [Associate Professor, UM2]
Esther Pacitti [Associate Team Leader, Professor, UM2, HdR]
Konstantin Todorov [Associate Professor, UM2, since September]

Engineers

Emmanuel Castanier [Engineer, ADT WebSmatch]
Julien Champ [Engineer, ANR Otmedia]
Mathias Chouet [Engineer, RTRA PI@ntNet]
Guillaume Verger [Engineer, ANR DataRing]

PhD Students

Ayoub Ait Lahcen [ANR STAMP fellowship]
Yoann Couillec [Inria CORDIS fellowship (with the INDES team), since October]
Vincenzo Gulisano [Universidad Politecnica de Madrid, Spain]
Pierre Letessier [CIFRE fellowship, Telecom ParisTech]
Miguel Liroz [Inria CORDIS fellowship]
Saloua Litayem [Telecom ParisTech]
Imen Mami [UM2 fellowship]
Duy Hoa Ngo [ANR DataRing fellowship]
Jonas Dias [Universidade Federal de Rio de Janeiro, Brazil]
Mohamed Riadh Trad [Telecom ParisTech]
Saber Salah [Inria Hemera fellowship, since November]
Toufik Sarni [CNRS fellowship, LINA]
Maximilien Servajean [NUMEV-UM2 fellowship]

Post-Doctoral Fellow

Hervé Goëau [RTRA PI@ntNet]

Administrative Assistant

Annie Aliaga

2. Overall Objectives

2.1. Introduction

Modern science such as agronomy, bio-informatics, astronomy and environmental science must deal with overwhelming amounts of experimental data produced through empirical observation and simulation (<http://www.computational-sustainability.org>). Such data must be processed (cleaned, transformed, analyzed) in all kinds of ways in order to draw new conclusions, prove scientific theories and produce knowledge. However, constant progress in scientific observational instruments (e.g. satellites, sensors, large hadron collider) and simulation tools (that foster *in silico* experimentation, as opposed to traditional *in situ* or *in vivo* experimentation) creates a huge data overload. For example, climate modeling data are growing so fast that they will lead to collections of hundreds of exabytes (10^{18} bytes) expected by 2020.

Scientific data is also very complex, in particular because of heterogeneous methods used for producing data, the uncertainty of captured data, the inherently multi-scale nature (spatial scale, temporal scale) of many sciences and the growing use of imaging (e.g. satellite images), resulting in data with hundreds of attributes, dimensions or descriptors. Processing and analyzing such massive sets of complex scientific data is therefore a major challenge since solutions must combine new data management techniques with large-scale parallelism in cluster, grid or cloud environments.

Furthermore, modern science research is a highly collaborative process, involving scientists from different disciplines (e.g. biologists, soil scientists, and geologists working on an environmental project), in some cases from different organizations distributed over different countries. Since each discipline or organization tends to produce and manage its own data, in specific formats, with its own processes, integrating distributed data and processes gets difficult as the amounts of heterogeneous data grow.

Despite their variety, we can identify common features of scientific data: massive scale; manipulated through complex, distributed workflows; typically complex, e.g. multidimensional or graph-based; with uncertainty in the data values, e.g., to reflect data capture or observation; important metadata about experiments and their provenance; and mostly append-only (with rare updates).

Generic data management solutions (e.g. relational DBMS) which have proved effective in many application domains (e.g. business transactions) are not efficient for dealing with scientific data, thereby forcing scientists to build ad-hoc solutions which are labor-intensive and cannot scale. In particular, relational DBMSs have been lately criticized for their “one size fits all” approach. Although they have been able to integrate support for all kinds of data (e.g., multimedia objects, XML documents and new functions), this has resulted in a loss of performance and flexibility for applications with specific requirements because they provide both “too much” and “too little”. Therefore, it has been argued that more specialized DBMS engines are needed. For instance, column-oriented DBMSs, which store column data together rather than rows in traditional row-oriented relational DBMSs, have been shown to perform more than an order of magnitude better on decision-support workloads. The “one size does not fit all” counter-argument generally applies to cloud data management as well. Cloud data can be very large, unstructured (e.g. text-based) or semi-structured, and typically append-only (with rare updates). And cloud users and application developers may be in high numbers, but not DBMS experts. Therefore, current cloud data management solutions have traded consistency for scalability, simplicity and flexibility. As alternative to relational DBMS (which use the standard SQL language), these alternative solutions have been quoted as Not Only SQL (NOSQL) by the database research community.

The three main challenges of scientific data management can be summarized by: (1) scale (big data, big applications); (2) complexity (uncertain, multi-scale data with lots of dimensions), (3) heterogeneity (in particular, data semantics heterogeneity). The overall goal of Zenith is to address these challenges, by proposing innovative solutions with significant advantages in terms of scalability, functionality, ease of use, and performance. To produce generic results, these solutions will be in terms of architectures, models and algorithms that can be implemented in terms of components or services in specific computing environments, e.g. grid, cloud. To maximize impact, a good balance between conceptual aspects (e.g. algorithms) and practical aspects (e.g. software development) is necessary. We plan to design and validate our solutions by

working closely with scientific application partners (CIRAD, INRA, CEMAGREF, etc.). To further validate our solutions and extend the scope of our results, we also want to foster industrial collaborations, even in non scientific applications, provided that they exhibit similar challenges.

2.2. Highlights of the Year

Patrick Valduriez has been elected ACM Fellow (2013).

At the 2012 competition of the Ontology Alignment Evaluation Initiative (<http://oei.ontologymatching.org>), our YAM++ ontology matching tool ranked first at the Large Biomedical Ontologies (largebio) track.

Members of the team have published the first textbook on P2P data management [9], in the series Synthesis Lectures on Data Management by Morgan & Claypool Publishers.

3. Scientific Foundations

3.1. Data Management

Data management is concerned with the storage, organization, retrieval and manipulation of data of all kinds, from small and simple to very large and complex. It has become a major domain of computer science, with a large international research community and a strong industry. Continuous technology transfer from research to industry has led to the development of powerful DBMSs, now at the heart of any information system, and of advanced data management capabilities in many kinds of software products (application servers, document systems, search engines, directories, etc.).

The fundamental principle behind data management is *data independence*, which enables applications and users to deal with the data at a high conceptual level while ignoring implementation details. The relational model, by resting on a strong theory (set theory and first-order logic) to provide data independence, has revolutionized data management. The major innovation of relational DBMS has been to allow data manipulation through queries expressed in a high-level (declarative) language such as SQL. Queries can then be automatically translated into optimized query plans that take advantage of underlying access methods and indices. Many other advanced capabilities have been made possible by data independence : data and metadata modeling, schema management, consistency through integrity rules and triggers, transaction support, etc.

This data independence principle has also enabled DBMS to continuously integrate new advanced capabilities such as object and XML support and to adapt to all kinds of hardware/software platforms from very small smart devices (smart phone, PDA, smart card, etc.) to very large computers (multiprocessor, cluster, etc.) in distributed environments.

Following the invention of the relational model, research in data management has continued with the elaboration of strong database theory (query languages, schema normalization, complexity of data management algorithms, transaction theory, etc.) and the design and implementation of DBMS. For a long time, the focus was on providing advanced database capabilities with good performance, for both transaction processing and decision support applications. And the main objective was to support all these capabilities within a single DBMS.

The problems of scientific data management (massive scale, complexity and heterogeneity) go well beyond the traditional context of DBMS. To address them, we capitalize on scientific foundations in closely related domains: distributed data management, cloud data management, uncertain data management, metadata integration, data mining and content-based information retrieval.

3.2. Distributed Data Management

To deal with the massive scale of scientific data, we exploit large-scale distributed systems, with the objective of making distribution transparent to the users and applications. Thus, we capitalize on the principles of large-scale distributed systems such as clusters, peer-to-peer (P2P) and cloud, to address issues in data integration, scientific workflows, recommendation, query processing and data analysis.

Data management in distributed systems has been traditionally achieved by distributed database systems which enable users to transparently access and update several databases in a network using a high-level query language (e.g. SQL) [13]. Transparency is achieved through a global schema which hides the local databases' heterogeneity. In its simplest form, a distributed database system is a centralized server that supports a global schema and implements distributed database techniques (query processing, transaction management, consistency management, etc.). This approach has proved effective for applications that can benefit from centralized control and full-fledge database capabilities, e.g. information systems. However, it cannot scale up to more than tens of databases. Data integration systems, e.g. price comparators such as KelKoo, extend the distributed database approach to access data sources on the Internet with a simpler query language in read-only mode.

Parallel database systems extend the distributed database approach to improve performance (transaction throughput or query response time) by exploiting database partitioning using a multiprocessor or cluster system. Although data integration systems and parallel database systems can scale up to hundreds of data sources or database partitions, they still rely on a centralized global schema and strong assumptions about the network.

Scientific workflow management systems (SWfMS) such as Kepler (<http://kepler-project.org>) and Taverna (<http://www.taverna.org.uk>) allow scientists to describe and execute complex scientific procedures and activities, by automating data derivation processes, and supporting various functions such as provenance management, queries, reuse, etc. Some workflow activities may access or produce huge amounts of distributed data and demand high performance computing (HPC) environments with highly distributed data sources and computing resources. However, combining SWfMS with HPC to improve throughput and performance remains a difficult challenge. In particular, existing workflow development and computing environments have limited support for data parallelism patterns. Such limitation makes complex the automation and ability to perform efficient parallel execution on large sets of data, which may significantly slow down the execution of a workflow.

In contrast, peer-to-peer (P2P) systems [9] adopt a completely decentralized approach to data sharing. By distributing data storage and processing across autonomous peers in the network, they can scale without the need for powerful servers. Popular examples of P2P systems such as Gnutella and BitTorrent have millions of users sharing petabytes of data over the Internet. Although very useful, these systems are quite simple (e.g. file sharing), support limited functions (e.g. keyword search) and use simple techniques (e.g. resource location by flooding) which have performance problems. To deal with the dynamic behavior of peers that can join and leave the system at any time, they rely on the fact that popular data get massively duplicated.

Initial research on P2P systems has focused on improving the performance of query routing in the unstructured systems which rely on flooding, whereby peers forward messages to their neighbors. This work led to structured solutions based on Distributed Hash Tables (DHT), e.g. CHORD and Pastry, or hybrid solutions with super-peers that index subsets of peers. Another approach is to exploit gossiping protocols, also known as epidemic protocols. Gossiping has been initially proposed to maintain the mutual consistency of replicated data by spreading replica updates to all nodes over the network. It has since been successfully used in P2P networks for data dissemination. Basic gossiping is simple. Each peer has a complete view of the network (i.e. a list of all peers' addresses) and chooses a node at random to spread the request. The main advantage of gossiping is robustness over node failures since, with very high probability, the request is eventually propagated to all nodes in the network. In large P2P networks, however, the basic gossiping model does not scale as maintaining the complete view of the network at each node would generate very heavy communication traffic. A solution to scalable gossiping is by having each peer with only a partial view of the network, e.g. a list of tens of neighbor peers. To gossip a request, a peer chooses at random a peer in its partial view to send it the request. In addition, the peers involved in a gossip exchange their partial views to reflect network changes in their own views. Thus, by continuously refreshing their partial views, nodes can self-organize into randomized overlays which scale up very well.

We claim that a P2P solution is the right solution to support the collaborative nature of scientific applications as it provides scalability, dynamicity, autonomy and decentralized control. Peers can be the participants or

organizations involved in collaboration and may share data and applications while keeping full control over their (local) data sources.

But for very-large scale scientific data analysis or to execute very large data-intensive workflow activities (activities that manipulate huge amounts of data), we believe cloud computing (see next section), is the right approach as it can provide virtually infinite computing, storage and networking resources. However, current cloud architectures are proprietary, ad-hoc, and may deprive users of the control of their own data. Thus, we postulate that a hybrid P2P/cloud architecture is more appropriate for scientific data management, by combining the bests of both. In particular, it will enable the clean integration of the users' own computational resources with different clouds.

3.3. Cloud Data Management

Cloud computing encompasses on demand, reliable services provided over the Internet (typically represented as a cloud) with easy access to virtually infinite computing, storage and networking resources. Through very simple Web interfaces and at small incremental cost, users can outsource complex tasks, such as data storage, system administration, or application deployment, to very large data centers operated by cloud providers. Thus, the complexity of managing the software/hardware infrastructure gets shifted from the users' organization to the cloud provider. From a technical point of view, the grand challenge is to support in a cost-effective way the very large scale of the infrastructure which has to manage lots of users and resources with high quality of service.

Cloud customers could move all or part of their information technology (IT) services to the cloud, with the following main benefits:

- **Cost.** The cost for the customer can be greatly reduced since the IT infrastructure does not need to be owned and managed; billing is only based on resource consumption. For the cloud provider, using a consolidated infrastructure and sharing costs for multiple customers reduces the cost of ownership and operation.
- **Ease of access and use.** The cloud hides the complexity of the IT infrastructure and makes location and distribution transparent. Thus, customers can have access to IT services anytime, and from anywhere with an Internet connection.
- **Quality of Service (QoS).** The operation of the IT infrastructure by a specialized provider that has extensive experience in running very large infrastructures (including its own infrastructure) increases QoS.
- **Elasticity.** The ability to scale resources out, up and down dynamically to accommodate changing conditions is a major advantage. In particular, it makes it easy for customers to deal with sudden increases in loads by simply creating more virtual machines.

However, cloud computing has some drawbacks and not all applications are good candidates for being "cloudified". The major concern is wrt. data security and privacy, and trust in the provider (which may use not so trustful providers to operate). One earlier criticism of cloud computing was that customers get locked in proprietary clouds. It is true that most clouds are proprietary and there are no standards for cloud interoperability. But this is changing with open source cloud software such as Hadoop, an Apache project implementing Google's major cloud services such as Google File System and MapReduce, and Eucalyptus, an open source cloud software infrastructure, which are attracting much interest from research and industry.

There is much more variety in cloud data than in scientific data since there are many different kinds of customers (individuals, SME, large corporations, etc.). However, we can identify common features. Cloud data can be very large, unstructured (e.g. text-based) or semi-structured, and typically append-only (with rare updates). And cloud users and application developers may be in high numbers, but not DBMS experts.

Current cloud data management (NOSQL) solutions typically trade consistency for scalability, simplicity and flexibility. They use a radically different architecture than RDBMS, by exploiting (rather than embedding) a distributed file system such as Google File System (GFS) or Hadoop Distributed File System (HDFS), to store and manage data in a highly fault-tolerant manner. They tend to rely on a more specific data model, e.g. key-value store such as Google Bigtable, Hadoop Hbase or Apache CouchDB) with a simple set of operators easy to use from a programming language. For instance, to address the requirements of social network applications, new solutions rely on a graph data model and graph-based operators. User-defined functions also allow for more specific data processing. MapReduce is a good example of generic parallel data processing framework, on top of a distributed file system (GFS or HDFS). It supports a simple data model (sets of (key, value) pairs), which allows user-defined functions (map and reduce). Although quite successful among developers, it is relatively low-level and rigid, leading to custom user code that is hard to maintain and reuse. In Zenith, we exploit or extend these NOSQL technologies to fit our needs for scientific workflow management and scalable data analysis.

3.4. Uncertain Data Management

Data uncertainty is present in many scientific applications. For instance, in the monitoring of plant contamination by INRA teams, sensors generate periodically data which may be uncertain. Instead of ignoring (or correcting) uncertainty, which may generate major errors, we need to manage it rigorously and provide support for querying.

To deal with uncertainty, there are several approaches, e.g. probabilistic, possibilistic, fuzzy logic, etc. The *probabilistic approach* is often used by scientists to model the behavior of their underlying environments. However, in many scientific applications, data management and uncertain query processing are not integrated, i.e. the queries are usually answered using ad-hoc methods after doing manual or semi-automatic statistical treatment on the data which are retrieved from a database. In Zenith, we aim at integrating scientific data management and query processing within one system. This should allow scientists to issue their queries in a query language without thinking about the probabilistic treatment which should be done in background in order to answer the queries. There are two important issues which any PDBMS should address: 1) how to represent a probabilistic database, i.e. data model; 2) how to answer queries using the chosen representation, i.e. query evaluation.

One of the problems on which we focus is *scalable query processing* over uncertain data. A naive solution for evaluating probabilistic queries is to enumerate all possible worlds, i.e. all possible instances of the database, execute the query in each world, and return the possible answers together with their cumulative probabilities. However, this solution can not scale up due to the exponential number of possible worlds which a probabilistic database may have. Thus, the problem is quite challenging, particularly due to exponential number of possibilities that should be considered for evaluating queries. In addition, most of our underlying scientific applications are not centralized; the scientists share part of their data in a *P2P* manner. This distribution of data makes very complicated the processing of probabilistic queries. To develop efficient query processing techniques for distributed scientific applications, we can take advantage of two main distributed technologies: *P2P* and *Cloud*. Our research experience in P2P systems has proved us that we can propose scalable solutions for many data management problems. In addition, we can use the cloud parallel solutions, e.g. MapReduce, to parallelize the task of query processing, when possible, and answer queries of scientists in reasonable execution times. Another challenge for supporting scientific applications is uncertain data integration. In addition to managing the uncertain data for each user, we need to integrate uncertain data from different sources. This requires revisiting traditional data integration in major ways and dealing with the problems of uncertain mediated schema generation and uncertain schema mapping.

3.5. Metadata Integration

Nowdays, scientists can rely on web 2.0 tools to quickly share their data and/or knowledge (e.g. ontologies of the domain knowledge). Therefore, when performing a given study, a scientist would typically need to access and integrate data from many data sources (including public databases). To make high numbers

of scientific data sources easily accessible to community members, it is necessary to identifying semantic correspondences between metadata structures or models of the related data sources. The main underlying task is called matching, which is the process of discovering semantic correspondences between metadata structures such as database schema and ontologies. Ontology is a formal and explicit description of a shared conceptualization in term of concepts (i.e., classes, properties and relations). For example, the matching may be used to align gene ontologies or anatomical metadata structures.

To understand a data source content, metadata (data that describe the data) is crucial. Metadata can be initially provided by the data publisher to describe the data structure (e.g. schema), data semantics based on ontologies (that provide a formal representation of the domain knowledge) and other useful information about data provenance (publisher, tools, methods, etc.). Scientific metadata is very heterogeneous, in particular because of the great autonomy of the underlying data sources, which leads to a large variety of models and formats. The high heterogeneity makes the matching problem very challenging. Furthermore, the number of ontologies and their size grow fastly, so does their diversity and heterogeneity. As a result, schema/ontology matching has become a prominent and challenging topic [4].

3.6. Data Mining

Data mining provides methods to discover new and useful patterns from very large sets of data. These patterns may take different forms, depending on the end-user's request, such as:

- **Frequent itemsets and association rules [1].** In this case, the data is usually a table with a high number of rows and the algorithm extracts correlations between column values. This problem was first motivated by commercial and marketing purposes (e.g. discovering frequent correlations between items bought in a shop, which could help selling more). A typical example of frequent itemset from a sensor network in a smart building would say that “in 20% rooms, the door is closed, the room is empty, and lights are on.”
- **Frequent sequential pattern extraction.** This problem is very similar to frequent itemset mining, but in this case, the order between events has to be considered. Let us consider the smart-building example again. A frequent sequence, in this case, could say that “in 40% rooms, lights are on at time i , the room is empty at time $i+j$ and the door is closed at time $i+j+k$ ”. Discovering frequent sequences has become a crucial need in marketing, but also in security (detecting network intrusions for instance) in usage analysis (web usage is one of the main applications) and any domain where data arrive in a specific order (usually given by timestamps).
- **Clustering [12].** The goal of clustering algorithms is to group together data that have similar characteristics, while ensuring that dissimilar data will not be in the same cluster. In our example of smart buildings, we would find clusters of rooms, where offices will be in one category and copy machine rooms in another one because of their characteristics (hours of people presence, number of times lights are turned on and off, etc.).

One of the main problems for data mining methods recently was to deal with data streams. Actually, data mining methods have first been designed for very large data sets where complex algorithms of artificial intelligence were not able to complete within reasonable time responses because of data size. The problem was thus to find a good trade-off between time response and results relevance. The patterns described above well match this trade-off since they both provide interesting knowledge for data analysts and allow algorithm having good time complexity on the number of records. Itemset mining algorithms, for instance, depend more on the number of columns (for a sensor it would be the number of possible items such as temperature, presence, status of lights, etc.) than the number of lines (number of sensors in the network). However, with the ever growing size of data and their production rate, a new kind of data source has recently emerged as data streams. A data stream is a sequence of events arriving at high rate. By “high rate”, we usually admit that traditional data mining methods reach their limits and cannot complete in real-time, given the data size. In order to extract knowledge from such streams, a new trade-off had to be found and the data mining community has investigated approximation methods that could allow maintaining a good quality of results for the above patterns extraction.

For scientific data, data mining now has to deal with new and challenging characteristics. First, scientific data is often associated to a level of uncertainty (typically, sensed values have to be associated to the probability that this value is correct or not). Second, scientific data might be extremely large and need cloud computing solutions for their storage and analysis. Eventually, we will have to deal with high dimension and heterogeneous data.

3.7. Content-based Information Retrieval

Today's technologies for searching information in scientific data mainly rely on relational DBMS or text based indexing methods. However, content-based information retrieval has progressed much in the last decade and is now considered as one of the most promising for future search engines. Rather than restricting search to the use of metadata, content-based methods attempt to index, search and browse digital objects by means of signatures describing their actual content. Such methods have been intensively studied in the multimedia community to allow searching the massive amount of raw multimedia documents created every day (e.g. 99% of web data are audio-visual content with very sparse metadata). Successful and scalable content-based methods have been proposed for searching objects in large image collections or detecting copies in huge video archives. Besides multimedia contents, content-based information retrieval methods recently started to be studied on more diverse data such as medical images, 3D models or even molecular data. Potential applications in scientific data management are numerous. First of all, to allow searching the huge collections of scientific images (earth observation, medical images, botanical images, biology images, etc.) but also to browse large datasets of experimental data (e.g. multisensor data, molecular data or instrumental data). Despite recent progress, scalability remains a major issue, involving complex algorithms (such as similarity search, clustering or supervised retrieval), in high dimensional spaces (up to millions of dimensions) with complex metrics (Lp, Kernels, sets intersections, edit distances, etc.). Most of these algorithms have linear, quadratic or even cubic complexities so that their use at large scale is not affordable without consistent breakthrough. In Zenith, we plan to investigate the following challenges:

- **High-dimensional similarity search.** Whereas many indexing methods were designed in the last 20 years to retrieve efficiently multidimensional data with relatively small dimensions, high-dimensional data have been more challenging due to the well-known dimensionality curse. Only recently have some methods appeared that allow approximate Nearest Neighbors queries in sub-linear time, in particular, Locality Sensitive Hashing methods which offer new theoretical insights in high-dimensional Euclidean spaces and proved the interest of random projections. But there are still some challenging issues that need to be solved including efficient similarity search in any kernel or metric spaces, efficient construction of knn-graphs or relational similarity queries.
- **Large-scale supervised retrieval.** Supervised retrieval aims at retrieving relevant objects in a dataset by providing some positive and/or negative training samples. To solve such task, there has been a focused interest on using Support Vector Machines (SVM) that offer the possibility to construct generalized, non-linear predictors in high-dimensional spaces using small training sets. The prediction time complexity of these methods is usually linear in dataset size. Allowing hyperplane similarity queries in sub-linear time is for example a challenging research issue. A symmetric problem in supervised retrieval consists in retrieving the most relevant object categories that might contain a given query object, providing huge labeled datasets (up to millions of classes and billions of objects) and very few objects per category (from 1 to 100 objects). SVM methods that are formulated as quadratic programming with cubic training time complexity and quadratic space complexity are clearly not usable. Promising solutions to such problems include hybrid supervised-unsupervised methods and supervised hashing methods.
- **P2P content-based retrieval.** Content-based P2P retrieval methods appeared recently as a promising solution to manage masses of data distributed over large social networks, particularly when the data cannot be centralized for privacy or cost reasons (which is often the case in scientific social networks, e.g. botanist social networks). However, current methods are limited to very simple similarity search paradigms. In Zenith, we will consider more advanced P2P content-based retrieval and

mining methods such as k-nn graphs construction, large-scale supervised retrieval or multi-source clustering.

4. Application Domains

4.1. Data-intensive Scientific Applications

The application domains covered by Zenith are very wide and diverse, as they concern data-intensive scientific applications, i.e. most scientific applications. Since the interaction with scientists is crucial to identify and tackle data management problems, we are dealing primarily with application domains for which Montpellier has an excellent track record, i.e. agronomy, environmental science, life science, with scientific partners like INRA, IRD, CIRAD and IRSTEA (prev. CEMAGREF). However, we are also addressing other scientific domains (e.g. astronomy, oil extraction) through our international collaborations (e.g. in Brazil).

Let us briefly illustrate some representative examples of scientific applications on which we have been working on.

- **Management of astronomical catalogs.** An example of data-intensive scientific applications is the management of astronomical catalogs generated by the Dark Energy Survey (DES) project on which we are collaborating with researchers from Brazil. In this project, huge tables with billions of tuples and hundreds of attributes (corresponding to dimensions, mainly double precision real numbers) store the collected sky data. Data are appended to the catalog database as new observations are performed and the resulting database size is estimated to reach 100TB very soon. Scientists around the globe can query the database with queries that may contain a considerable number of attributes. The volume of data that this application holds poses important challenges for data management. In particular, efficient solutions are needed to partition and distribute the data in several servers. An efficient partitioning scheme should try to minimize the number of fragments accessed in the execution of a query, thus reducing the overhead associated to handle the distributed execution.
- **Pesticide reduction.** In a pesticide reduction application, with CEMAGREF, we plan to work on sensor data for plant monitoring. Sensors are used to observe the development of diseases and insect attacks in the agricultural farms, aiming at using pesticides only when necessary. The sensors periodically send to a central system their data about different measures such as plants contamination, temperature or moisture level. A decision support system analyzes the sent data, and triggers a pesticide treatment only when needed. However, the data sent by sensors are not entirely certain. The main reasons for uncertainty are the effect of climate events on sensors, e.g. rain, unreliability of the data transmission media, fault in sensors, etc. This requires to deal with uncertain data in modeling and querying to be used for data analysis and data mining.
- **Botanical data sharing.** Botanical data is highly decentralized and heterogeneous. Each actor has its own expertise domain, hosts its own data, and describes them in a specific format. Furthermore, botanical data is complex. A single plant's observation might include many structured and unstructured tags, several images of different organs, some empirical measurements and a few other contextual data (time, location, author, etc.). A noticeable consequence is that simply identifying plant species is often a very difficult task; even for the botanists themselves (the so-called taxonomic gap). Botanical data sharing should thus speed up the integration of raw observation data, while providing users an easy and efficient access to integrated data. This requires to deal with social-based data integration and sharing, massive data analysis and scalable content-based information retrieval. We address this application in the context of the French initiative PI@ntNet, with CIRAD and IRD.
- **Deepwater oil exploitation.** An important step in oil exploitation is pumping oil from ultra-deepwater from thousand meters up to the surface through long tubular structures, called risers. Maintaining and repairing risers under deep water is difficult, costly and critical for the environment. Thus, scientists must predict risers fatigue based on complex scientific models and observed data for the risers. Risers fatigue analysis requires a complex workflow of data-intensive activities which

may take a very long time to compute. A typical workflow takes as input files containing riser information, such as finite element meshes, winds, waves and sea currents, and produces result analysis files to be further studied by the scientists. It can have thousands of input and output files and tens of activities (e.g. dynamic analysis of risers movements, tension analysis, etc.). Some activities, e.g. dynamic analysis, are repeated for many different input files, and depending on the mesh refinements, each single execution may take hours to complete. To speed up risers fatigue analysis requires parallelizing workflow execution, which is hard to do with existing SWfMS. We address this application in collaboration with UFRJ, and Petrobras.

These application examples illustrate the diversity of requirements and issues which we are addressing with our scientific application partners (CIRAD, INRA, CEMAGREF, etc.). To further validate our solutions and extend the scope of our results, we also want to foster industrial collaborations, even in non scientific applications, provided that they exhibit similar challenges.

5. Software

5.1. WebSmatch (Web Schema Matching)

Participants: Zohra Bellahsène, Emmanuel Castanier, Rémi Coletta, Duy Hoa Ngo, Patrick Valduriez [contact].

URL: <http://websmatch.gforge.inria.fr/>

In the context of the Action de Développement Technologique (ADT) started in october 2010, WebSmatch is a flexible, open environment for discovering and matching complex schemas from many heterogeneous data sources over the Web. It provides three basic functions: (1) metadata extraction from data sources; (2) schema matching (both 2-way and n-way schema matching), (3) schema clustering to group similar schemas together. WebSmatch is being delivered through Web services, to be used directly by data integrators or other tools, with RIA clients. Implemented in Java, delivered as Open Source Software (under LGPL) and protected by a deposit at APP (Agence de Protection des Programmes). WebSmatch is being used by Datapublica and CIRAD to integrate public data sources.

5.2. YAM++ ((not) Yet Another Matcher)

Participants: Zohra Bellahsène [contact], Duy Hoa Ngo, Konstantin Todorov.

URL: <http://www2.lirmm.fr/~dngo/>

YAM++ is a tool for discovering semantic correspondences between ontologies. YAM++ supports several matching strategies: machine learning; generic methods when learning data are not available; discovering alignment of ontologies represented in different languages. Furthermore, since this year YAM++ is able to deal with large scale ontology matching.

5.3. SON (Shared-data Overlay Network)

Participants: Ayoub Ait Lahcen, Fady Draïdi, Esther Pacitti, Didier Parigot [contact], Patrick Valduriez, Guillaume Verger.

URL: <http://www-sop.inria.fr/teams/zenith/SON>

SON is an open source development platform for P2P networks using web services, JXTA and OSGi. SON combines three powerful paradigms: components, SOA and P2P. Components communicate by asynchronous message passing to provide weak coupling between system entities. To scale up and ease deployment, we rely on a decentralized organization based on a DHT for publishing and discovering services or data. In terms of communication, the infrastructure is based on JXTA virtual communication pipes, a technology that has been extensively used within the Grid community. Using SON, the development of a P2P application is done through the design and implementation of a set of components. Each component includes a technical code that provides the component services and a code component that provides the component logic (in Java). The complex aspects of asynchronous distributed programming (technical code) are separated from code components and automatically generated from an abstract description of services (provided or required) for each component by the component generator.

5.4. P2Prec (P2P recommendation service)

Participants: Fady Draidi, Esther Pacitti [contact], Didier Parigot, Guillaume Verger.

URL: <http://p2prec.gforge.inria.fr>

P2Prec is a recommendation service for P2P content sharing systems that exploits users social data. To manage users social data, we rely on Friend-Of-A-Friend (FOAF) descriptions. P2Prec has a hybrid P2P architecture to work on top of any P2P content sharing system. It combines efficient DHT indexing to manage the users FOAF files with gossip robustness to disseminate the topics of expertise between friends. P2Prec is implemented in java using the Data-Shared Overlay Network (SON) infrastructure which is the basis for the ANR DataRing project.

5.5. ProbDB (Probabilistic Database)

Participants: Reza Akbarinia [contact], Patrick Valduriez, Guillaume Verger.

URL: <http://probdb.gforge.inria.fr>

ProbDB is a probabilistic data management system to manage uncertain data on top of relational DBMSs. One of the main features of the prototype is its portability; that means with a minimum effort it can be implemented over any DBMS. In ProbDB, we take advantage of the functionalities provided by almost all DBMSs, particularly the query processing functions. It is implemented in Java on top of PostgreSQL.

5.6. Pl@ntNet-Identify

Participants: Mathias Chouet, Hervé Goëau, Alexis Joly [contact].

URL: <http://identify.plantnet-project.org>

Pl@ntNet-Identify is a web application dedicated to the image-based identification of plants. It has been developed jointly by ZENITH, the AMAP UMR team (CIRAD) and the Inria team IMEDIA. It allows submitting one or several query pictures of a plant and browse the matching species in a large collection of social image data, i.e. plant images collected by the members of a social network. It also allows users to enrich the knowledge of the application by uploading their own pictures in the reference collection. Nowadays, the dataset includes more than 17K images posted by about 100 members of Telabotanica¹ social network. In 2012, about 5000 identification sessions have been recorded. The client side of the application is implemented in Javascript whereas the server side (visual search engine) is mostly implemented in C++.

5.7. Pl@ntNet-DataManager

Participants: Mathias Chouet [contact], Alexis Joly.

¹<http://www.tela-botanica.org/>

PI@ntNet-DataManager is a software dedicated to managing and sharing distributed heterogeneous botanical data. It is developed jointly by ZENITH, the AMAP UMR team (CIRAD) and Tela Botanica non profit organization. It allows scientists to define data structures dedicated to their own datasets, and share parts of their structures and data with collaborators in a decentralized way. PI@ntNet DataManager offers innovative features like partial or complete P2P synchronization between distant databases (master-master), and a user friendly data structure editor. It also provides full text search, querying, CSV import/export, SQL export, image management, and geolocation. DataManager is built on NoSQL technology (CouchDB database), Javascript (Node.js), HTML5 and CSS3, and may be deployed on a server or run on a local machine (standalone version for Linux, Windows, Mac). It is being used by researchers and engineers of the PI@ntNet Project (CIRAD, INRA, Inria, IRD, Tela-Botanica) to manage taxonomical referentials, herbarium data and geolocated plant observations.

5.8. SnoopIm

Participants: Julien Champ [contact], Alexis Joly, Pierre Letessier.

URL: <http://otmedia.lirmm.fr/OTmedia/>

SnoopIm is a content-based search engine allowing to discover and retrieve small visual patterns or objects in large collections of pictures (such as logos on clothes, road signs in the background, paintings on walls, etc.) and to derive statistics from them (frequency, visual cover, size variations, etc.). Query objects to be searched can be either selected from the indexed collection of photos, or selected from an external picture (by simply providing its URL). The web application allows online search of multiple users and has a cache feature to speed-up the processing of seen queries. It is implemented in Javascript on top of a C++ library developed in collaboration with INA ². The software is used at INA by archivists and sociologists in the context of the Transmedia Observatory project.

6. New Results

6.1. Data and Metadata Management

6.1.1. Uncertain Data Management

Participants: Reza Akbarinia, Patrick Valduriez, Guillaume Verger.

Data uncertainty in scientific applications can be due to many different reasons: incomplete knowledge of the underlying system, inexact model parameters, inaccurate representation of initial boundary conditions, inaccuracy in equipments, error in data entry, etc.

One of the areas, in which uncertainty management is important, is the integration of heterogeneous data sources, in the sense where usually there may be an uncertainty in the possible mappings between the attributes of the sources. Usually the human interaction is demanded to help the system in choosing the correct mappings. In [30], we propose a pay-as-you-go data integration solution that aims at performing the data integration in a fully automated way. Our solution takes advantage of attribute correlations by using functional dependencies, and captures uncertainty in mediated schemas using a probabilistic data model. It allows integrating a given set of data sources, as well as incrementally integrating additional sources, without needing to restart the process from scratch. We implemented our solution, and compared it with a baseline approach. The performance evaluation results show significant performance gains of our solution in terms of recall and precision compared to the baseline approaches.

²<http://www.ina-sup.com/>

Another problem that arises in many applications such as data integration systems is that of Entity Resolution (ER). ER is the process of identifying tuples that represent the same real-world entity. It has been well studied in the literature for certain data, but it has not been deeply investigated for uncertain data. Existing proposals for the ER problem are not applicable to the above examples since they ignore probability values completely and return the most similar tuples as the solution. Furthermore, the semantics of the solution for the ERUD problem has not been clearly defined in the literature. In [31], we address the ERUD problem. We adopt the well-known possible worlds semantics for defining the semantics for the ERUD problem, and propose a PTIME algorithm for a large class of similarity functions, i.e. context-free. For the rest of similarity functions, i.e. context-sensitive, we use Monte-Carlo randomization for approximating the answer. We propose a parallel version of our Monte-Carlo algorithm using the MapReduce framework. To the best of our knowledge, this is the first study of the ERUD problem that adopts the possible world semantics and the first efficient algorithm for implementing it.

Another important problem in uncertain data management is the efficient processing of probabilistic queries. We have continued the development of our probabilistic database prototype, called ProbDB (Probabilistic Database) that deals with large-scale probabilistic data sharing. ProbDB divides each probabilistic query into two parts: probabilistic and deterministic (i.e. non probabilistic). The deterministic part is executed by the underlying RDBMS, and the rest of work is done by our probabilistic query processing algorithms that are executed over the data returned by the RDBMS.

6.1.2. Metadata Integration

Participants: Zohra Bellahsène, Emmanuel Castanier, Duy Hoa Ngo, Patrick Valduriez.

Our work on metadata integration encompassed ontology matching and open data source integration.

The major focus of our work in 2012 was to deal with large scale ontology matching and scalability. To improve the matching quality of YAM++, we designed a new IR-based measure to deal with terminological heterogeneity in real world ontologies. To deal with large ontology matching, we designed a method based on indexing concepts from their labels and comments. Our approach aims at reducing the search space when comparing the concepts of the input ontologies. For this purpose, we designed three filters: Description Filter, Context Filter and Label Filter. These methods make use of the Lucene search engine for indexing and searching the context of entities in the input ontologies. Another contribution lies on the Fast Semantic Filtering method, which refines the discovered mappings in the ontology matching task. The aim of the Semantic Filter is to detect and reject inconsistent mappings by exploring semantic information of entities in the input ontologies [45]. The originality of our method is to use a new structural indexing technique and a heuristic to generate relative disjointness axioms. At the 2012 competition of the Ontology Alignment Evaluation Initiative (<http://oaei.ontologymatching.org>), YAM++ was one of the best matchers, with very good results in all tracks. It obtained the first position in the Large BioMed Track [55].

Integrating open data sources can yield high value information but raises major problems in terms of metadata extraction, data source integration and visualization of integrated data. In [34], [33], we describe WebSmatch, a flexible environment for Web data integration, based on a real, end-to-end data integration scenario over public data from Data Publica. WebSmatch supports the full process of importing, refining and integrating data sources and uses third party tools for high quality visualization. We use a typical scenario of public data integration which involves problems not solved by current tools: poorly structured input data sources (XLS files) and rich visualization of integrated data.

6.1.3. High-dimensional data management

Participants: Mohamed Riadh Trad, Alexis Joly, Saloua Litayem.

High dimensional data hashing is essential for scaling up and distributing data analysis applications involving feature-rich objects, such as text documents, images or multi-modal entities (scientific observations, events, etc.). In this first research track, we first investigated the use of high dimensional hashing methods for efficiently approximating K-NN Graphs [47], particularly in distributed environments. We highlighted the importance of balancing issues on the performance of such approaches and show why the baseline approach

using Locality Sensitive Hashing does not perform well. Our new KNN-join method is based on RMMH, a hash function family based on randomly trained classifiers that we introduced in 2011. We show that the resulting hash tables are much more balanced and that the number of resulting collisions can be greatly reduced without degrading quality. We further improve the load balancing of our distributed approach by designing a parallelized local join algorithm, implemented within the MapReduce framework. In other work [43], we address the problem of speeding-up the prediction phase of linear Support Vector Machines via Locality Sensitive Hashing. Whereas the mainstream work in the field is focused on training classifiers on huge amount of data, less efforts are spent on the counterpart scalability issue: how to apply big trained models efficiently on huge non annotated collections ? In this work, we propose building efficient hash-based classifiers that are applied in a first stage in order to approximate the exact results and alter the hypothesis space. Experiments performed with millions of one-against-one classifiers show that the proposed hash-based classifier can be more than two orders of magnitude faster than the exact classifier with minor losses in quality.

6.2. Data and Process Sharing

6.2.1. Hybrid P2P/cloud Architecture

Participants: Esther Pacitti, Patrick Valduriez.

Zenith adopts a hybrid P2P/cloud architecture. P2P naturally supports the collaborative nature of scientific applications, with autonomy and decentralized control. Peers can be the participants or organizations involved in collaboration and may share data and applications while keeping full control over some of their data (a major requirement for our application partners). But for very-large scale data analysis or very large workflow activities, cloud computing is appropriate as it can provide virtually infinite computing, storage and networking resources. Such hybrid architecture also enables the clean integration of the users' own computational resources with different clouds.

In [24], we define Zenith's architecture with P2P data services and cloud data services. We model an online scientific community as a set of peers and relationships between them. The peers have their own data sources. The relationships are between any two or more peers and indicate how the peers and their data sources are related, e.g. friendship, same semantic domain, similar schema. The P2P data services include basic services (metadata and uncertain data management): recommendation, data analysis and workflow management through the Shared-data Overlay Network (SON) middleware. The cloud P2P services include data mining, content-based information retrieval and workflow execution. These services can be accessed through web services, and each peer can use the services of multiple clouds.

6.2.2. Social-based P2P Data Sharing

Participants: Reza Akbarinia, Emmanuel Castanier, Esther Pacitti, Didier Parigot, Patrick Valduriez, Guillaume Verger.

As a validation of the ANR DataRing project, we have developed P2PShare, a P2P system for large-scale probabilistic data sharing in scientific communities. P2PShare leverages content-based and expert-based recommendation. It is designed to manage probabilistic and deterministic data in P2P environments. It provides a flexible environment for integration of heterogeneous sources, and takes into account the social based aspects to discover high quality results for queries by privileging the data of friends (or friends of friends), who are expert on the topics related to the query.

Using the Shared-Data Overlay Network (SON), we have implemented a prototype of P2PShare that integrates three major DataRing services: ProbDB, a probabilistic database management service for relational data; WebSmatch, an environment for Web data integration; and P2Prec, a social-based P2P recommendation service for large-scale content sharing.

In [50], we describe the demo of P2PShare's main services, e.g., gossiping topics of interest among friends, key- word querying for contents, and probabilistic queries over datasets.

6.2.3. View Selection in Distributed Data Warehousing

Participants: Zohra Bellahsène, Imen Mami.

Scientific data generate large amounts of data which have to be collected and stored for analytical purpose. One way to help managing and analyzing large amounts of data is data warehousing, whereby views over data are materialized [23]. At large scale, a data warehouse can be distributed. We have examined the problem of choosing a set of views and a set of data warehouse nodes at which these views should be materialized so that the full query workload is answered with the lowest cost. To address this problem, we extended our view selection method that we proposed for the centralized case. Thus, we modelled the distributed view selection problem as a Constraint Satisfaction Problem (CSP). Furthermore, we introduced the distributed AND-OR view graph, which can be seen as an extensive form of the AND-OR view graph to reflect the relation between views and communication network within the distributed scenario. The experiment results show that our approach provides better performance compared with the genetic algorithm in term of the solution quality (i.e., the quality of the obtained set of materialized views). We demonstrated experimentally that our approach provides better results in term of cost savings when the view selection is decided under space and maintenance cost constraints [44].

6.2.4. Scientific Workflow Management

Participants: Ayoub Ait Lahcen, Jonas Dias, Didier Parigot, Patrick Valduriez.

Scientific experiments based on computer simulations can be defined, executed and monitored using Scientific Workflow Management Systems (SWfMS). Several SWfMS are available, each with a different goal and a different engine. Due to the exploratory analysis, scientists need to run parameter sweep (PS) workflows, which are workflows that are invoked repeatedly using different input data. These workflows generate a large amount of tasks that are submitted to High Performance Computing (HPC) environments. Different execution models for a workflow may have significant differences in performance in HPC. However, selecting the best execution model for a given workflow is difficult due to the existence of many characteristics of the workflow that may affect the parallel execution.

In [36], we develop a study to show performance impacts of using different execution models in running PS workflows in HPC. Our study contributes by presenting a characterization of PS workflow patterns (the basis for many existing scientific workflows) and its behavior under different execution models in HPC. We evaluated four execution models to run workflows in parallel. Our study measures the performance behavior of small, large and complex workflows among the evaluated execution models. The results can be used as a guideline to select the best model for a given scientific workflow execution in HPC. Our evaluation may also serve as a basis for workflow designers to analyze the expected behavior of an HPC workflow engine based on the characteristics of PS workflows.

This work was done in the context of the the CNPq-Inria project DatLuge and FAPERJ-Inria P2Pcloud project .

In the context of SON, we also proposed a declarative workflow language based on service/activity rules. In [27], [46], we present a formal approach that combines component-based development with well-understood methods and techniques from the field of Attribute Grammars and Data-Flow Analysis in order to specify the behavior of P2P applications, and then construct an abstract representation (i.e., Data-Dependency Graph) to perform analyzes on it. This formal approach makes it possible to infer a dependency graph for SON applications that provides for automatic parallelization.

6.2.5. Plants identification and classification from social image data

Participants: Hervé Goëau, Alexis Joly, Saloua Litayem.

This work is done in collaboration with the botanists of the AMAP UMR team (CIRAD) and with Inria team IMEDIA. Inspired by citizen sciences, the main goal of this trans-disciplinary work is to speed up the collection and integration of raw botanical observation data, while providing to potential users an easy and efficient access to this botanical knowledge. We therefore did continue working intensively on plants identification and classification [54], [37], [38], [26]. We first developed a new interactive method [37] for the visual identification of plants from social image data. Contrary to previous content-based identification methods and systems that mainly relied on leaves, or in few other cases on flowers, it makes use of five different organs and plant's views including habit, flowers, fruits, leaves and bark. Thanks to an interactive query widget, the tagging process of the different organs and views is as simple as drag-and-drop operations and does not require any expertise in botany. All training pictures used by the system were continuously collected during one year through a crowdsourcing application and more than 17K images are now integrated. System-oriented and human-centered evaluations of the application show that the results are already satisfactory and therefore very promising in the long term to identify a richer flora.

Besides, we did continue working on leaf-based identification notably through the organization of and participation to ImageCLEF plant identification evaluation campaign 2012 [54].

Finally we did apply one of our former work related to multi-source shared-nearest neighbors clustering to an original experiment aimed at evaluating if we were able to automatically recover morphological classifications built by the botanists themselves [38]. The results are very promising, since all clusters discovered automatically could be easily matched to one node of a morphological tree built by botanists.

6.3. Scalable Data Analysis

6.3.1. StreamCloud

Participants: Vincenzo Gulisano, Patrick Valduriez.

Recent years have witnessed the growth of a new class of data-intensive applications that do not fit the DBMS query paradigm. Instead, the data arrive at high speeds taking the form of an unbounded sequence of values (data streams) and queries run continuously returning new results as new data arrive. Examples of data streams are sensor data (e.g. in environmental applications) or IP packets (e.g. in a network monitoring application). The unbounded nature of data streams makes it impossible to store the data entirely in bounded memory. Current research efforts have mainly focused on scaling in the number of queries and/or query operators having overlooked the scalability with respect to the stream volume.

Current Stream Processing Engines do not scale with the input load due to single-node bottlenecks. Additionally, they are based on static configurations that lead to either under or over-provisioning. In [21], [22], we present StreamCloud, a scalable and elastic stream processing engine for processing large data stream volumes. StreamCloud uses a novel parallelization technique that splits queries into subqueries that are allocated to independent sets of nodes in a way that minimizes the distribution overhead. Its elastic protocols exhibit low intrusiveness, enabling effective adjustment of resources to the incoming load. Elasticity is combined with dynamic load balancing to minimize the computational resources used. We present the system design, implementation and a thorough evaluation of the scalability and elasticity of the fully implemented system.

6.3.2. Mining Uncertain Data Streams

Participants: Reza Akbarinia, Florent Masseglia.

Dealing with uncertainty by using probabilistic approaches has gained increasing attention these past few years. One of the main requirements for uncertain data mining is the ability to discover Probabilistic Frequent Itemsets (PFI). However, PFI mining, particularly in uncertain data streams, is very challenging and needs the development of new techniques, since approaches designed for deterministic data are not applicable in this context. In [29], we propose an efficient solution for exact PFI mining over data streams with sliding windows. Our proposal includes efficient solutions for updating frequentness probability of itemsets and thus fast extraction of PFI, whenever transactions are added or removed from the sliding window. To the best of our knowledge, this is the first efficient solution for data stream PFI mining. We have conducted an extensive experimental evaluation of our approach over synthetic and real-world data sets; the results illustrate its very good performance.

6.3.3. Detecting Rare Events in Massive Datasets

Participant: Florent Masseglia.

In this work, we consider that rare events are very small clusters typically representing less than 0.01% of the entire dataset. Finding these abnormal events allows to identify the emergence of possible anomalies in their very early stages. Such a scenario is generally difficult to handle as it lies at the frontier between outlier detection and clustering and is characterized by a clear challenge to avoid false negatives. To address this challenge, we take a backward approach and propose RARE, a framework that identifies and isolates the abnormal/rare regions. The dense regions are identified using a radius-limited density-driven variant of k-means and adjacent regions are merged to form new regions. These newly formed regions are gradually augmented as long as a density-driven condition is respected. When no more dense regions are observed, the remaining data is clustered and presented for further analysis to human experts. The framework is tested on a medical application and compared against human analysis. The experiments show that rare events that were missed during human analysis because of the multivariate character of the data can be discovered by our approach.

This work is funded by the labex NUMEV and a patent application involving Inria, CNRS, UM2 and INSERM has been filled.

6.3.4. Highly Informative Feature Set Mining

Participant: Florent Masseglia.

For many textual collections, the number of features is often overly large. As these features can be very redundant, it is desirable to have a small, succinct, yet highly informative collection of features that describes the key characteristics of a dataset. Information theory is one such tool for us to obtain this feature collection. In [48], we mainly contribute to the improvement of efficiency for the process of selecting the most informative feature set over high-dimensional unlabeled data. We propose a heuristic theory for informative feature set selection from high dimensional data. Moreover, we design data structures that enable us to compute the entropies of the candidate feature sets efficiently. We also develop a simple pruning strategy that eliminates the hopeless candidates at each forward selection step. We test our method through experiments on real-world data sets, showing that our proposal is very efficient.

6.3.5. Clustering Users with Evolving Profiles in Usage Streams

Participant: Florent Masseglia.

Existing data stream models commonly assume that users' records or profiles in data streams will not be updated once they arrive. In many applications such as web usage, however, the users' records/profiles may evolve along time. This kind of streaming transactions are referred to as bi-streaming data (*i.e.* the data evolves temporally in two dimensions, the flowing of transactions as with the traditional data streams, and the evolving of users' profiles inside the streams, which makes bi-streaming data different from traditional data streams). The two-dimensional evolving of bi-streaming data brings difficulties on modeling and clustering for exploring the users' behaviors. In [49], we propose three models to summarize bi-streaming data, which are the batch model, the Evolving Objects (EO) model and the Dynamic Data Stream (DDS) model. Through creating, updating and deleting user profiles, the models summarize the behaviors of each user as an object. Based on these models, clustering algorithms are employed to identify the user groups. The proposed models are tested on a real-world data set showing that the DDS model can summarize the bi-streaming data efficiently and effectively, providing better basis for clustering user profiles than the other two models.

6.3.6. Scalable Mining of Small Visual Objects

Participants: Pierre Letessier, Julien Champ, Alexis Joly.

Automatically linking multimedia documents that contain one or several instances of the same visual object has many applications including: salient events detection, relevant patterns discovery in scientific data or simply web browsing through hyper-visual links. Whereas efficient methods now exist for searching rigid objects in large collections, discovering them from scratch is still challenging in terms of scalability, particularly when the targeted objects are rather small. In this work [40], we formally revisit the problem of mining or discovering such objects, and then generalized two kinds of existing methods for probing candidate object seeds: weighted adaptive sampling and hashing based methods. We then introduce a new hashing strategy, working first at the visual level, and then at the geometric level. Experiments conducted on millions of images show that our method outperforms state-of-the-art.

This method was integrated within a visual-based media event detection system in the scope of a French project called the transmedia observatory. It allows the automatic discovery of the most circulated images across the main news media (news websites, press agencies, TV news and newspapers). The main originality of the detection is to rely on the transmedia contextual information to denoise the raw visual detections and consequently focus on the most salient trans-media events. This work was presented at ACM Multimedia Grand Challenge 2012 [39]. The movie presented during this event is available at <http://www.otmedia.fr/?p=217>.

7. Bilateral Contracts and Grants with Industry

7.1. Data Publica (2010-2013)

Participants: Emmanuel Castanier, Patrick Valduriez.

Data Publica (<http://www.data-publica.com>) is a startup providing a web portal for open data which can be public, private, free or charged. We collaborate with Data Publica through our WebSmatch tool on technologies for automatic schema extraction and matching from high numbers of data sources. A first contribution has been the development of an Excel extraction component based on machine learning techniques.

8. Partnerships and Cooperations

8.1. Regional Initiatives

8.1.1. Labex NUMEV, Montpellier

URL: <http://www2.lirmm.fr/numev>

We are participating in the Laboratory of Excellence (labex) NUMEV (Digital and Hardware Solutions, Modelling for the Environment and Life Sciences) headed by University of Montpellier 2 in partnership with CNRS, University of Montpellier 1, and Inria. NUMEV seeks to harmonize the approaches of hard sciences and life and environmental sciences in order to pave the way for an emerging interdisciplinary group with an international profile. The NUMEV project is decomposed in four complementary research themes: Modeling, Algorithms and computation, Scientific data (processing, integration, security), Model-Systems and measurements. Patrick Valduriez heads the theme on scientific data.

8.1.2. Institut de Biologie Computationnelle (IBC), Montpellier

URL: <http://www.ibc-montpellier.fr>

IBC is a 5 year project with a funding of 2Meuros by the MENRT (“Investissements d’Avenir” program) to develop innovative methods and software to integrate and analyze biological data at large scale in health, agronomy and environment. Patrick Valduriez heads the workpackage on integration of biological data and knowledge.

8.1.3. ModSiCS2020 Working Group, Montpellier

The ModSiCS2020 (Modeling and Simulation of Complex Systems in 2020) working group was set up by UM2 to analyze the local situation (forces and weaknesses, current projects), identify the critical research directions and propose concrete actions in terms of research projects, equipment facilities, human resources and training to be encouraged in Montpellier. The group was headed by Patrick Valduriez and gathered a small number of experts in different disciplines (agronomy, bioinformatics, computer science, environmental science, life science, etc.). The conclusions of the group [57] were presented at the ModSiCS2020 workshop on Data, Models and Theories for Complex Systems: new challenges and opportunities, organized by UM2 in march. Following the work of the group, a “Groupement d’Intérêt Scientifique (GIS)” is being proposed in Montpellier.

8.2. National Initiatives

8.2.1. ANR

8.2.1.1. VERSO DataRing(2008-2012, 300Keuros)

Participants: Reza Akbarinia, Zohra Bellahsène, Emmanuel Castanier, Duy Hoa Ngo, Esther Pacitti, Didier Parigot, Guillaume Verger, Patrick Valduriez [leader].

URL: <http://www-sop.inria.fr/teams/zenith/dataring>

The DataRing project, headed by P. Valduriez, involves the Leo project-team (Inria Saclay Ile de France), LIG, LIRMM and Telecom ParisTech. The objective is to address the problem of data sharing for online communities, such as social networks (e.g. sites like MySpace and Facebook) and professional communities (e.g. research communities, online technical support groups) which are becoming a major killer application of the web. The project addresses this problem by organizing community members in a peer-to-peer (P2P) network ring across distributed data source owners where each member can share data with the others through a P2P overlay network. In this project, we study the following problems: schema matching, query processing with data uncertainty, data indexing and caching, data privacy and trust. To validate our approach, we develop services based on our prototypes WebSmatch, SON, P2Prec and ProbDB.

8.2.1.2. OTMedia (2011-2013), 150Keuros

Participants: Alexis Joly, Julien Champ, Pierre Letessier.

The Transmedia Observatory project, launched in November 2010, aims to develop processes, tools and methods to better understand the challenges and changes in the media sphere. Studying and tracking media events on all media (web, press, radio and television) are the two prioritized research areas. OTMedia brings together six partners: Inria (ZENITH), AFP (French Press Agency), INA (French National Audiovisual Institute), Paris 3 Sorbonne Nouvelle (researchers in Information Science and Communication), Syllabs (a SME specialized in semantic analysis and automatic creation of text) and the Computer Science Laboratory of Avignon University. ZENITH addresses more specifically the research challenges related to the trans-media tracking of visual contents (images and videos) and the clustering of heterogeneous information sources.

8.2.2. Others

8.2.2.1. RTRA Pl@ntNet (2009-2013), 1Meuros

Participants: Alexis Joly, Hervé Goëau, Saloua Litayem, Mathias Chouet.

The PI@ntNet project <http://www.plantnet-project.org/> was launched in 2009 by a large international consortium headed by three groups with complementary skills (UMR AMAP ³, IMEDIA project team at Inria, and the French botanical network TelaBotanica ⁴), with financial support from the Agropolis Foundation. Due to the departure of Nozha Boujemaa from the head of IMEDIA and the mobility of Alexis Joly in 2011, ZENITH has been entrusted with the Inria's management and scientific coordination of the project in spring 2012. The objectives of the project are (i) to develop cutting-edge transdisciplinary research at the frontier between integrative botany and computational sciences, based on the use of large datasets and expertise in plant morphology, anatomy, agronomy, taxonomy, ecology, biogeography and practical uses (ii) provide free, easy-access software tools and methods for plant identification and for the aggregation, management, sharing and utilization of plant-related data (iii) promote citizen science as a powerful means to enrich databases with new information on plants and to meet the need for capacity building in agronomy, botany and ecology.

8.2.2.2. CIFRE INA/Inria (2011-2013), 100Keuros

Participants: Alexis Joly, Pierre Letessier.

This CIFRE contract with INA funds a 3-years PhD (Pierre Letessier) to address research challenges related to content-based mining of visual objects in large collections.

8.2.2.3. CNRS INS2I Mastodons (2012), 30Keuros

Participants: Florent Masseglia, Patrick Valduriez, Esther Pacitti [leader].

This project deals with the problems of big data in the context of life science, where masses of data are being produced, e.g. by Next Generation Sequencing technologies or plant phenotyping platforms. In this project, Zenith addresses the specific problems of large-scale data analysis and data sharing.

8.3. European Initiatives

8.3.1. FP7 Projects

Program: FP7

Project acronym: CHORUS+ (avmediasearch.eu)

Project title: European coordination action on Audio-Visual Media Search

Duration: 2010 - 2012

Coordinator: JCP consulting

Other partners: CERTH-ITI (Greece), University of Trento (Italy), HES-SO (Switzerland), Technicolor (France), Vienna University of Technology (Austria), Engineering Ingegneria Informatica SPA (Italy), JRC Institute for Prospective Technological Studies (EU)

Abstract: CHORUS+ <http://avmediasearch.eu/> objective is to coordinate national and international projects and initiatives in the Search-engine domain and to extend this Coordination in non-European countries. ZENITH actively participated to this action, Alexis Joly being member of the steering committee and leader of a work package. We particularly promoted scientific data as an essential challenge to be addressed by this community through the co-organization of international events (CBMI 2012 panel, ImageCLEF 2012, international workshop on search computing) and discussions with leaders of European projects belonging to the cluster of the coordination action. Besides, we did work on technology transfer issues and the potential of benchmarking campaigns as a tool to foster it (conduction of a survey of about hundred people from both academy and industry, organization of a think-tank with about 20 stakeholders, writing of a recommendation report for the EU commission).

8.4. International Initiatives

8.4.1. Inria International Partners

We have regular scientific relationships with research laboratories in

³<http://amap.cirad.fr/en/>

⁴<http://www.tela-botanica.org/>

- North America: Univ. of Waterloo (Tamer Özsu), Univ. of California, Santa Barbara (Divy Agrawal, Amr El Abbadi).
- Asia: National Univ. of Singapore (Beng Chin Ooi, Stéphane Bressan), Wonkwang University, Korea (Kwangjin Park)
- Europe: Univ. of Amsterdam (Naser Ayat, Hamideh Afsarmanesh), Univ. of Madrid (Ricardo Jiménez-Periz), UPC Barcelona (Josep Lluís Larriba Pey, Victor Munoz)

8.4.2. Participation In International Programs

We are involved in the following international actions:

- CNPq-Inria project DatLuge (Data & Task Management in Large Scale, 2009-2012) with UFRJ (Marta Mattoso, Vanessa Braganholo, Alexandre Lima), LNCC, Rio de Janeiro (Fabio Porto), and UFPR, Curitiba (Eduardo Almeida) to work on large scale scientific workflows;
- FAPERJ-Inria project SwfP2Pcloud (Data-centric workflow management in hybrid P2P clouds, 2011-2013) with UFRJ (Marta Mattoso, Vanessa Braganholo, Alexandre Lima) and LNCC, Rio de Janeiro (Fabio Porto) to work on large scale scientific workflows in hybrid P2P clouds;
- CNPq-Inria project Hoscar (HPC and data management, 2012-2015) with LNCC (Fabio Porto), UFC, UFRGS (Philippe Navaux), UFRJ (Alvaro Coutinho, Marta Mattoso) to work on data management in high performance computing environments;
- EGIDE Osmoze project SECC (SERVICES for CURRICULA Comparison, 2011-2012), with Riga Technical University (Janis Grundspenkis, Marité Kirikova) to work on automatic analysis and mapping of conceptual trees and maps acquired from digital documents.

8.5. International Research Visitors

8.5.1. Visits of International Scientists

Prof. Jens Dittrich (Univ. Saarland, Germany) gave a seminar at LIRMM on data management with MapReduce.

Prof. Marta Mattoso (UFRJ, Rio de Janeiro) gave a seminar at LIRMM in the context of IBC on data provenance in scientific workflows.

8.5.2. Visits to International Teams

Esther Pacitti and Patrick Valduriez were invited researchers at the National University of Singapore in July.

9. Dissemination

9.1. Scientific Animation

Participation in the editorial board of scientific journals:

- VLDB Journal: P. Valduriez.
- Proceedings of the VLDB Endowment (PVLDB): E. Pacitti.
- Distributed and Parallel Databases, Kluwer Academic Publishers: P. Valduriez.
- Internet and Databases: Web Information Systems, Kluwer Academic Publishers: P. Valduriez.
- Journal of Information and Data Management, Brazilian Computer Society Special Interest Group on Databases: P. Valduriez.
- Book series “Data Centric Systems and Applications” (Springer-Verlag): P. Valduriez.
- Ingénierie des Systèmes d’Information, Hermès : P. Valduriez.

Participation in conference program committees :

- Int. Conf. on VLDB 2012: E. Pacitti
- Int. Conf. on Extending DataBase Technologies (EDBT): Z. Bellahsene, E. Pacitti, 2012, 2013
- Int. Conf. on Information and Knowledge Management (CIKM) 2012: E. Pacitti
- IEEE Int. Conf. on Data Engineering (ICDE) 2012: Z. Bellahsene
- Int. Conf. on Middleware 2012: E. Pacitti
- Journées Bases de Données Avancées (BDA), 2012: R. Akbarinia, Z. Bellahsene
- Int. Workshop on Open Data, 2012: P. Valduriez
- ACM Int. Conf. on Multimedia Retrieval 2012: Alexis Joly
- CLEF-labs 2012, 2013: Alexis Joly
- Int. Conf. on Multimedia Modeling 2013: Alexis Joly
- Int. Conf. on Cooperative Information systems (CoopIS): Z. Bellahsene, 2012
- Int. Conf. on Advanced Information Systems Engineering (CAiSE) : Z. Bellahsene, 2012, 2013
- European Semantic Web Conference (ESWC) 2012: Z. Bellahsene
- ACM Symposium On Applied Computing (ACM SAC, data stream track), 2012: F. Masegla
- Int. Conf. on Advances in Databases, Knowledge, and Data Applications (DBKDA), 2012: F. Masegla
- Conférence Internationale Francophone sur l'Extraction et la Gestion de Connaissance (EGC), 2012: F. Masegla
- IEEE Int. Conf. on Data Mining (ICDM), 2012: F. Masegla
- Int. Conf. on Artificial Intelligence in Medicine (LEMEDS@AIME), 2011: F. Masegla
- Int. Workshop on Multimedia Data Mining (MDM@KDD), 2011: F. Masegla
- Int. Pacific-Asia Conf. on Knowledge Discovery and Data Mining (PAKDD), 2012: F. Masegla
- IEEE Int. Conf. on Advanced Information Networking and Applications (AINA) 2012: H. Bouziane
- Rencontres Francophones sur les Aspects Algorithmiques de Telecommunications (AlgoTel), 2012: H. Bouziane
- IEEE Int. Parallel and Distributed Processing Symposium (IPDPS) 2013: H. Bouziane

Other activities (national):

- The members of Zenith have always been strongly involved in organizing the French database research community, in the context of the I3 GDR and the BDA conference.
- Esther Pacitti organized a workshop on bigdata and life science in the context of the Mastodons project in Montpellier (<http://www.lirmm.fr/~pacitti/Mastodons.html>).
- Patrick Valduriez gave a talk on "Open Data Integration with WebSmatch" at SophiaConf2012, organized by Telecom Valley in Sophia-Antipolis. He also chaired the ModSiCS2020 workshop on Data, Models and Theories for Complex Systems: new challenges and opportunities, organized by UM2 in march.

Other activities (international):

- Since 2011, Alexis Joly organizes a system-oriented evaluation initiative dedicated to plant's identification in the context of the CLEF evaluation forum <http://clef2012.org/>. This initiative is the first one dealing with the evaluation of content-based technologies in the context of ecological data. Finally, he chaired a panel on "Scientific Multimedia Data" at the 10th international workshop on Content-Based Multimedia Indexing (CBMI2012). He also chaired a session on "Search in Large-Scale Multimedia Data" at the International Workshop on Search Computing (sponsored by the EU commission).
- Patrick Valduriez chaired a panel on "Big Data: Scientific and Societal Challenges" at the Berkeley Inria Stanford (BIS2012) workshop in June.

9.2. Teaching - Supervision - Juries

9.2.1. Teaching

Licence : Enseignant, titre du cours, nombre d'heures en équivalent TD, niveau (L1, L2, L3), université, pays

Master : Enseignant, titre du cours, nombre d'heures en équivalent TD, niveau (M1, M2), université, pays

Doctorat : Enseignant, titre du cours, nombre d'heures en équivalent TD, université, pays

Most permanent members of Zenith teach at the Licence and Master degree levels at UM2.

Reza Akbarinia:

Master Research: Data sharing in P2P, 9h, level M2, Faculty of Science, UM2

Master : Distributed Data Management, 9h, level M1, Faculty of Science, UM2

Licence Pro: JavaScript/jQuery, 10h, level L3, IUT, UM2

Licence: Computer Tools, 36h, level L2, Faculty of Science, UM2

Zohra Bellahsène:

IUT 1: Relational Databases , 100h, level L1, IUT, UM2

IUT 2: Object-relational Databases, 60h, Level L2, IUT, UM2

Master Computer Science : Data Integration, 15h, level M2, Faculty of Science, UM2

Hinde Bouziane:

Licence: Networks, 36h, level L3, Faculty of Science, UM2

Master DECOL: Networks and communication, 90h, level M1, Faculty of Science, UM2

Master DECOL: Distributed Data Management, 6h, level M1, Faculty of Science, UM2

Master IPS: Introduction to operating systems and networks, 36h, level M1, Faculty of Science, UM2

Master Research DECOL: Data sharing in P2P, 15h, level M2, Faculty of Science, UM2

Esther Pacitti:

IG3: Database design, physical organization, 54h, level L3, Polytech' Montpellier, UM2

IG4: Networks, 42h, level M1, Polytech' Montpellier, UM2

IG4: Object-relational databases, 32h, level M1, Polytech' Montpellier, UM2

IG5: Distributed systems, virtualization, 27h, level M2, Polytech' Montpellier, UM2

Master Research: Data sharing in P2P, 4,5h, level M2, Faculty of Science, UM2

Didier Parigot:

Master Research: Data sharing in P2P, 9h, level M2, Faculty of Science, UM2

Patrick Valduriez:

Master Research: Data sharing in P2P, 13h, level M2, Faculty of Science, UM2

Professional: Distributed Information Systems, 55h, level M2, Capgemini Institut

Professional: XML, 40h, level M2, Orsys Formation

9.2.2. Supervision

- PhD: Vincenzo Gulisano, StreamCloud: An Elastic Parallel-Distributed Stream Processing Engine, 20 dec. 2012, Univ. Madrid, Advisors: Ricardo Jiménez Peris, Patrick Valduriez
- PhD: Ayoub Ait-lahcen, Developing component-based applications with a data-centric approach and within a service-oriented P2P architecture: specification, analysis and middleware, 15 dec. 2012, Univ. Mohammed V Rabat and Univ. Nice, Advisors: Driss Aboutajdine and Didier Parigot
- PhD: Duy Hoa Ngo, A generic approach to ontology matching, 14 dec. 2012, UM2, Advisor: Zohra Bellahsène, co-advisor: Rémi Coletta
- PhD: Imen Mami, Une approche déclarative pour la modélisation et la résolution du problème de la sélection de vues à matérialiser, 15 nov. 2012, UM2, Advisor: Zohra Bellahsène, co-advisor: Rémi Coletta
- PhD: Toufik Sarni, Real-time support for software transactional memory in multicore systems, 16 oct. 2012, Univ. Nantes, Advisor: Patrick Valduriez, co-advisor: Audrey Queudet
- PhD in progress: Yoann Couillec, Langages de programmation et données ouvertes, started oct. 2012, Univ. Nice Sophia-Antipolis, Advisors: Manuel Serrano and Patrick Valduriez
- PhD in progress: Jonas Dias, Interactive Workflows for Data-Centric Experiments, started nov. 2011, Universidade Federal de Rio de Janeiro, Brazil, Advisors: Marta Mattoso and Patrick Valduriez
- PhD in progress : Pierre Letessier, Frequent Visual Objects Discovery in Multimedia Collections, started nov. 2009, Telecom ParisTech, Advisor: Nozha Boujemaa, co-advisors: Olivier Buisson and Alexis Joly
- PhD in progress : Miguel Liroz, Massive Data Management for Scientific Applications, started oct. 2010, UM2, Advisors: Esther Pacitti and Patrick Valduriez, co-advisor: Reza Akbarinia
- PhD in progress : Saber Salah, Optimizing a Cloud for Data Mining Primitives, started nov. 2012, UM2, Advisors: Florent Massegli, co-advisor: Reza Akbarinia
- PhD in progress : Maximilien Servajean, Decentralized and Personalized Recommendation Protocols for Content Sharing: application to phenotyping, started oct. 2011, UM2, Advisor: Esther Pacitti, co-advisor: Pascal Neveu
- PhD in progress : Naser Ayat, Uncertain Data Integration, started sept. 2010, University of Amsterdam, Netherlands, Advisors: Hamideh Afsarmanesh and Patrick Valduriez, co-advisor: Reza Akbarinia

9.2.3. Juries

Members of the team participated to the following Ph.D. committees: R. Akbarinia: Asma Souihli (TELECOM ParisTech); Z. Bellahsène: Imen Mami (UM2), Duy Hoa Ngo (UM2); F. Massegli: Pierre-Nicolas Mougél (Univ. Lyon); E. Pacitti: Fady Draïdi (UM2), Hien Truong (Univ. Nancy); D. Parigot: Ayoub Ait Lahcen (Univ. Rabat, Univ. Nice); P. Valduriez: Fady Draïdi (UM2), Toufiq Sarni (Univ. Nantes), Vincenzo Gulisano (Univ. Madrid).

E. Pacitti participated to the PES (Prime d'Excellence Scientifique) committee of the CNU (Conseil National des Universités).

9.3. Popularization

ZENITH has animated Inria's booth during 3 days at WWW 2012 conference in Lyon through demos of Pl@ntNet tools <http://www.inria.fr/actualite/actualites-inria/inria-a-www2012>.

Several movies were realized in collaboration with INA in order to disseminate the results of the French ANR project OTMedia <http://www.inria.fr/actualite/actualites-inria/inria-a-www2012>. One of the movie was used as a visual support for an oral presentation of Alexis Joly at ACM Multimedia 2012 in the context of our participation to the Multimedia Grand Challenge <http://www.acmmm12.org/>.

10. Bibliography

Major publications by the team in recent years

- [1] P. PONCELET, F. MASSEGLIA, M. TEISSEIRE (editors). *Data Mining Patterns: New Methods and Applications*, Premier Reference Source, Idea Group, 2007, ISBN 978-1599041629, <http://hal.inria.fr/lirmm-00365419/en>.
- [2] R. AKBARINIA, E. PACITTI, P. VALDURIEZ. *Best Position Algorithms for Efficient Top-k Query Processing*, in "Information Systems", 2011, vol. 36, n^o 6, p. 973-989, <http://hal.inria.fr/lirmm-00607882/en>.
- [3] R. AKBARINIA, P. VALDURIEZ, G. VERGER. *Efficient Evaluation of SUM Queries Over Probabilistic Data*, in "IEEE Transactions on Knowledge and Data Engineering", 2013, To appear, <http://hal.inria.fr/lirmm-00652293/en>.
- [4] Z. BELLAHSENE, A. BONIFATI, E. RAHM. *Schema Matching and Mapping*, Springer, March 2011, 320, <http://hal.inria.fr/lirmm-00581346/en>.
- [5] A. BENOIT, H. L. BOUZIANE, Y. ROBERT. *Optimizing the reliability of streaming applications under throughput constraints*, in "International Journal of Parallel Programming", 2011, vol. 39, n^o 5, p. 584-614 [DOI : 10.1007/s10766-011-0165-6], <http://hal.inria.fr/inria-00574555/en>.
- [6] M. EL DICK, E. PACITTI, R. AKBARINIA, B. KEMME. *Building a Peer-to-Peer Content Distribution Network with High Performance, Scalability and Robustness*, in "Information Systems", 2011, vol. 36, n^o 2, p. 222-247, <http://hal.inria.fr/lirmm-00607898/en>.
- [7] V. GULISANO, R. JIMENEZ-PERIS, M. PATINO-MARTÍNEZ, C. SORIENTE, P. VALDURIEZ. *StreamCloud: An Elastic and Scalable Data Streaming System*, in "IEEE Transactions on Parallel and Distributed Systems", December 2012, vol. 23, n^o 12, p. 2351-2365 [DOI : 10.1109/TPDS.2012.24], <http://hal.inria.fr/lirmm-00748992>.
- [8] E. OGASAWARA, D. DE OLIVEIRA, P. VALDURIEZ, D. DIAS, F. PORTO, M. MATTOSO. *An Algebraic Approach for Data-Centric Scientific Workflows*, in "Proceedings of VLDB", 2011, vol. 4, n^o 11, p. 1328-1339, <http://hal.inria.fr/hal-00640431/en>.
- [9] E. PACITTI, R. AKBARINIA, M. EL DICK. *P2P Techniques for Decentralized Applications*, Morgan & Claypool Publishers, 2012, 104, <http://hal.inria.fr/lirmm-00748635>.
- [10] E. PACITTI, P. VALDURIEZ, M. MATTOSO. *Grid Data Management: Open Problems and New Issues*, in "Journal of Grid Computing", 2007, vol. 5, n^o 3, p. 273-281, <http://hal.inria.fr/inria-00473481/en>.
- [11] J.-A. QUIANÉ-RUIZ, P. LAMARRE, P. VALDURIEZ. *A Self-Adaptable Query Allocation Framework for Distributed Information Systems*, in "The VLDB Journal", 2009, vol. 18, n^o 3, p. 649-674, <http://hal.archives-ouvertes.fr/hal-00374999/fr/>.
- [12] C. ZHANG, F. MASSEGLIA, Y. LECHEVALLIER. *ABS: The Anti Bouncing Model for Usage Data Streams*, in "IEEE Int. Conf. on Data Mining (ICDM)", 2010, p. 1169-1174.

- [13] T. M. ÖZSU, P. VALDURIEZ. *Principles of Distributed Database Systems, third edition*, Springer, 2011, 845, <http://hal.inria.fr/hal-00640392/en>.

Publications of the year

Doctoral Dissertations and Habilitation Theses

- [14] A. AIT LAHCEN. *Développement d'Applications à Base de Composants avec une Approche Centrée sur les Données et dans une Architecture Service et Pair-à-Pair : Spécification, Analyse et Intergiciel*, Université de Nice Sophia-Antipolis and Université MoHammed V - Agdal-Rabat, December 2012, <http://hal.inria.fr/tel-00766329>.
- [15] F. DRAIDI. *Recommandation Pair-à-Pair pour Communautés en Ligne à Grande Echelle*, Université Montpellier II - Sciences et Techniques du Languedoc, March 2012, <http://hal.inria.fr/tel-00766963>.
- [16] V. GULISANO. *StreamCloud: un moteur de traitement de streams parallèle et distribué*, Universidad Politécnica de Madrid, December 2012, <http://hal.inria.fr/tel-00768281>.
- [17] I. MAMI. *Une approche déclarative pour la modélisation et la résolution du problème de la sélection de vues à matérialiser*, Université Montpellier II - Sciences et Techniques du Languedoc, November 2012, <http://hal.inria.fr/tel-00760992>.
- [18] D. H. NGO. *Amélioration de l'alignement d'ontologies par les techniques d'apprentissage automatique, d'appariement de graphes et de recherche d'information*, Université Montpellier II - Sciences et Techniques du Languedoc, December 2012, <http://hal.inria.fr/tel-00767318>.
- [19] T. SARNI. *Vers une mémoire transactionnelle temps réel*, Université de Nantes, October 2012, <http://hal.inria.fr/tel-00750637>.

Articles in International Peer-Reviewed Journals

- [20] R. AKBARINIA, P. VALDURIEZ, G. VERGER. *Efficient Evaluation of SUM Queries Over Probabilistic Data*, in "IEEE Transactions on Knowledge and Data Engineering", 2013, To appear, <http://hal.inria.fr/lirmm-00652293>.
- [21] V. GULISANO, R. JIMENEZ-PERIS, M. PATINO-MARTÍNEZ, C. SORIENTE, P. VALDURIEZ. *StreamCloud: An Elastic and Scalable Data Streaming System*, in "IEEE Transactions on Parallel and Distributed Systems", December 2012, vol. 23, n^o 12, p. 2351-2365 [DOI : 10.1109/TPDS.2012.24], <http://hal.inria.fr/lirmm-00748992>.
- [22] V. GULISANO, R. JIMENEZ-PERIS, M. PATIÑO-MARTINEZ, C. SORIENTE, P. VALDURIEZ. *A Big Data Platform for Large Scale Event Processing*, in "ERCIM News", April 2012, vol. 2012, n^o 89, 2, <http://hal.inria.fr/lirmm-00748582>.
- [23] I. MAMI, Z. BELLAHSENE. *A survey of view selection methods*, in "Sigmod Record", March 2012, vol. 41, n^o 1, p. 20-29 [DOI : 10.1145/2206869.2206874], <http://hal.inria.fr/lirmm-00720157>.
- [24] E. PACITTI, P. VALDURIEZ. *Zenith: Scientific Data Management on a Large Scale*, in "ERCIM News", April 2012, vol. 2012, n^o 89, 2, <http://hal.inria.fr/lirmm-00748563>.

- [25] J.-A. QUIANÉ-RUIZ, P. LAMARRE, P. VALDURIEZ. *Satisfaction-based query replication - An automatic and self-adaptable approach for replicating queries in the presence of autonomous participants*, in "Distributed and Parallel Databases", February 2012, vol. 30, n^o 1, p. 1-26 [DOI : 10.1007/s10619-011-7086-7], <http://hal.inria.fr/lirmm-00748553>.
- [26] A. REBAI, A. JOLY, N. BOUJEMAA. *BLasso for object categorization and retrieval: Towards interpretable visual models*, in "Pattern Recognition", June 2012, vol. 45, n^o 6, p. 2377-2389 [DOI : 10.1016/J.PATCOG.2011.11.022], <http://hal.inria.fr/hal-00739706>.

International Conferences with Proceedings

- [27] A. AIT LAHCEN, S. MOULINE. *Defining and Analyzing P2P Applications with a Data-Dependency Formalism*, in "PDCAT'12: Parallel and Distributed Computing, Applications and Technologies", Beijing, China, IEEE, December 2012, 8, <http://hal.inria.fr/lirmm-00757286>.
- [28] A. AIT LAHCEN, D. PARIGOT. *A Lightweight Middleware for developing P2P Applications with Component and Service-Based Principles*, in "CSE'12: International Computational Science and Engineering", Pathos, Cyprus, December 2012, 8, <http://hal.inria.fr/lirmm-00757105>.
- [29] R. AKBARINIA, F. MASSEGLIA. *FMU: Fast Mining of Probabilistic Frequent Itemsets in Uncertain Data Streams*, in "BDA2012: Journées Bases de Données Avancées", France, 2012, 18, <http://hal.inria.fr/lirmm-00748605>.
- [30] N. AYAT, H. AFSARMANESH, R. AKBARINIA, P. VALDURIEZ. *An Uncertain Data Integration System*, in "ODBASE'2012: 11th International Conference on Ontologies, DataBases, and Applications of Semantics", Roma, Italy, 2012, 18, <http://hal.inria.fr/lirmm-00748621>.
- [31] N. AYAT, R. AKBARINIA, H. AFSARMANESH, P. VALDURIEZ. *Entity Resolution for Uncertain Data*, in "BDA'2012: Bases de Données Avancées", France, 2012, 20, <http://hal.inria.fr/lirmm-00748625>.
- [32] V. BAKIC, I. YAHIAOUI, S. MOUINE, S. LITAYEM OUERTANI, W. OUERTANI, A. VERROUST-BLONDET, H. GOËAU, A. JOLY. *Inria IMEDIA2's Participation at ImageCLEF 2012 Plant Identification Task*, in "CLEF (Online Working Notes/Labs/Workshop) 2012", Rome, Italy, September 2012, <http://hal.inria.fr/hal-00744901>.
- [33] E. CASTANIER, R. COLETTA, P. VALDURIEZ, C. FRISCH. *Public Data Integration with WebSmatch*, in "BDA'2012: Bases de Données Avancées", France, 2012, 4, <http://hal.inria.fr/lirmm-00750910>.
- [34] E. CASTANIER, R. COLETTA, P. VALDURIEZ, C. FRISCH, D. H. NGO, Z. BELLAHSENE. *Public Data Integration with WebSmatch*, in "WOD'2012: International Workshop on Open Data", Nantes, France, 2012, 6, <http://hal.inria.fr/lirmm-00750927>.
- [35] E. CASTANIER, R. COLETTA, P. VALDURIEZ, C. FRISCH, D. H. NGO, Z. BELLAHSENE. *Public Data Integration with WebSmatch*, in "Proceedings of the First International Workshop on Open Data", 2012, vol. abs/1205.2555, 8, <http://hal.inria.fr/lirmm-00750889>.
- [36] F. CHIRIGATI, V. SILVA, E. OGASAWARA, J. DIAS, F. PORTO, P. VALDURIEZ, M. MATTOSO. *Evaluating Parameter Sweep Workflows in High Performance Computing*, in "SWEET'12: 1st International Workshop

- on Scalable Workflow Enactment Engines and Technologies", United States, ACM, May 2012, 10, <http://hal.inria.fr/lirmm-00749968>.
- [37] H. GOËAU, P. BONNET, B. JULIEN, V. BAKIC, A. JOLY, J.-F. MOLINO. *Multi-Organ Plant Identification*, in "ACM International Workshop on Multimedia Analysis for Ecological Data", Nara, Japan, October 2012, <http://hal.inria.fr/hal-00739724>.
- [38] A. HAMZAOUI, A. JOLY, N. BOUJEMAA. *Plant Leaves Morphological Categorization with Shared Nearest Neighbours Clustering*, in "ACM International Workshop on Multimedia Analysis for Ecological Data", Nara, Japan, October 2012, <http://hal.inria.fr/hal-00739715>.
- [39] A. JOLY, J. CHAMP, P. LETESSIER, N. HERVÉ, O. BUISSON, M.-L. VIAUD. *Visual-Based Transmedia Events Detection*, in "ACM Multimédia 2012", Nara, Japan, ACM, November 2012, p. 1351-1352 [DOI : 10.1145/2393347.2396480], <http://hal.inria.fr/hal-00755696>.
- [40] P. LETESSIER, A. JOLY, O. BUISSON. *Scalable Mining of Small Visual Objects*, in "ACM Multimédia 2012", Nara, Japan, October 2012, <http://hal.inria.fr/hal-00739735>.
- [41] M. LIROZ-GISTAU, R. AKBARINIA, E. PACITTI, F. PORTO, P. VALDURIEZ. *Dynamic Workload-Based Partitioning for Large-Scale Databases*, in "DEXA'2012: 23rd International Conference on Database and Expert Systems Applications", Vienna, Austria, LNCS, 2012, p. 183-190 [DOI : 10.1007/978-3-642-32597-7_16], <http://hal.inria.fr/lirmm-00748549>.
- [42] M. LIROZ-GISTAU, R. AKBARINIA, E. PACITTI, F. PORTO, P. VALDURIEZ. *DynPart: Dynamic Partitioning for Large-Scale Databases*, in "BDA'2012: 28e journées Bases de Données Avancées", Clermont-Ferrand, France, 2012, <http://hal.inria.fr/lirmm-00748585>.
- [43] S. LITAYEM OUERTANI, A. JOLY, N. BOUJEMAA. *Hash-Based Support Vector Machines Approximation for Large Scale Prediction*, in "BMVC - British Machine Vision Conference - 2012", Surrey, United Kingdom, British Machine Vision Conference (BMVC), 2012, <http://hal.inria.fr/hal-00733912>.
- [44] I. MAMI, Z. BELLAHSENE, R. COLETTA. *View Selection under Multiple Resource Constraints in a Distributed Context*, in "DEXA'2012: Database and Expert Systems Applications", Vienne, Austria, S. W. LIDDLE, K.-D. SCHEWE, A. M. TJOA, X. ZHOU (editors), Lecture Notes in Computer Science, Springer, September 2012, p. 281-296 [DOI : 10.1007/978-3-642-32597-7_25], <http://hal.inria.fr/lirmm-00736722>.
- [45] D. H. NGO, Z. BELLAHSENE. *YAM++ : A multi-strategy based approach for Ontology matching task*, in "Knowledge Engineering and Knowledge Management", Galway City, Ireland, M. D'AQUIN, A. NIKOLOV (editors), 2012, 5, <http://hal.inria.fr/lirmm-00720639>.
- [46] D. PARIGOT, A. AIT LAHCEN, S. MOULINE. *Toward Data-Centric View on Service-Based Component Systems: Formalism, Analysis and Execution*, in "PDP'11: Euromicro International Conference on Parallel, Distributed and Network-Based Computing", Garching, Germany, February 2012, <http://hal.inria.fr/lirmm-00648265>.
- [47] R. TRAD, A. JOLY, N. BOUJEMAA. *Distributed KNN-graph approximation via hashing*, in "ICMR - International Conference on Multimedia Retrieval - 2012", Hong-Kong, Hong Kong, ACM, June 2012 [DOI : 10.1145/2324796.2324847], <http://hal.inria.fr/hal-00739713>.

- [48] C. ZHANG, F. MASSEGLIA, X. ZHANG. *Discovering Highly Informative Feature Set Over High Dimensions*, in "ICTAI'2012: 24th International Conference on Tools with Artificial Intelligence", Greece, IEEE, November 2012, 1, <http://hal.inria.fr/lirmm-00753807>.
- [49] C. ZHANG, F. MASSEGLIA, X. ZHANG. *Modeling and Clustering Users with Evolving Profiles in Usage Streams*, in "TIME'2012: 19th International Symposium on Temporal Representation and Reasoning", United Kingdom, September 2012, p. 133-140, <http://hal.inria.fr/lirmm-00753791>.

National Conferences with Proceeding

- [50] F. DRAIDI, E. PACITTI, D. PARIGOT, G. VERGER, P. VALDURIEZ, R. COLETTA, R. AKBARINIA, E. CASTANIER. *P2PShare: a Social-based P2P Data Sharing System*, in "BDA'2012: Bases de Données Avancées", Clermont-Ferrand, France, October 2012, 5, <http://hal.inria.fr/lirmm-00757169>.
- [51] I. MAMI, Z. BELLAHSENE, R. COLETTA. *A Constraint Satisfaction based Approach to View Selection in a Distributed Context*, in "BDA'2012: Bases de Données Avancées", Clermont-ferrand, France, September 2012, 10, <http://hal.inria.fr/lirmm-00736723>.
- [52] D. H. NGO, Z. BELLAHSENE. *YAM++ : (not) Yet Another Matcher for Ontology Matching Task*, in "Bases de Données Avancées", France, N. ANCIAUX (editor), 2012, 5, <http://hal.inria.fr/lirmm-00720648>.
- [53] R. TRAD, A. JOLY, N. BOUJEMAA. *Distributed approximate KNN Graph construction for high dimensional Data*, in "BDA - 28e journées Bases de Données Avancées - 2012", Clermont-Ferrand, France, October 2012, <http://hal.inria.fr/ha-00756624>.

Conferences without Proceedings

- [54] H. GOËAU, P. BONNET, A. JOLY, I. YAHIAOUI, D. BARTHÉLÉMY, N. BOUJEMAA, J.-F. MOLINO. *The IMAGECLEF 2012 Plant identification Task*, in "CLEF 2012", Rome, Italy, September 2012, <http://hal.inria.fr/ha-00739740>.
- [55] D. H. NGO, Z. BELLAHSENE. *YAM++ - Results for OAEI 2012*, in "International Semantic Web Conference", United States, 2012, <http://hal.inria.fr/lirmm-00758720>.

Scientific Books (or Scientific Book chapters)

- [56] E. PACITTI, R. AKBARINIA, M. EL DICK. *P2P Techniques for Decentralized Applications*, Morgan & Claypool Publishers, 2012, 104, <http://hal.inria.fr/lirmm-00748635>.

Other Publications

- [57] P. VALDURIEZ, N. ARNAUD, F. BRIAND, O. GASCUEL, O. GIMENEZ, C. GODIN, H. JOURDE, P. KOSUTH, P. NEVEU, E. PACITTI, A. PARMEGGIANI. *Final Report of the ModSysC2020 Working Group - Data, Models and Theories for Complex Systems: new challenges and opportunities*, January 2012, <http://hal.inria.fr/lirmm-00749078>.