



IN PARTNERSHIP WITH:
CNRS

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

Activity Report 2013

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
**Data and Knowledge Representation
and Processing**

Table of contents

1. Members	1
2. Overall Objectives	1
2.1. Introduction	1
2.2. Highlights of the Year	2
3. Research Program	3
3.1. Data Management	3
3.2. Distributed Data Management	3
3.3. Cloud Data Management	5
3.4. Big Data	6
3.5. Uncertain Data Management	6
3.6. Big data Integration	7
3.7. Data Mining	7
3.8. Content-based Information Retrieval	8
4. Application Domains	9
5. Software and Platforms	10
5.1. WebSmatch (Web Schema Matching)	10
5.2. SON (Shared-data Overlay Network)	11
5.3. P2Prec (P2P recommendation service)	11
5.4. ProbDB (Probabilistic Database)	11
5.5. Pl@ntNet-mobile	11
5.6. Pl@ntNet-DataManager	12
5.7. SnoopIm	12
5.8. SciFloware	12
6. New Results	13
6.1. Big Data Integration	13
6.1.1. Probabilistic Data Integration	13
6.1.2. Open Data Integration	13
6.1.3. Pricing Integrated Data	14
6.2. Distributed Indexing and Searching	14
6.2.1. P2P Search and Recommendation	14
6.2.2. Spatial Queries in Wireless Data Broadcasting	15
6.3. Big Data Analysis	15
6.3.1. Big Data Analysis using Algebraic Workflows	15
6.3.2. Big Data Partitioning	15
6.4. Data Stream Mining	16
6.4.1. Mining Uncertain Data Streams	16
6.4.2. Itemset Mining over Tuple-Evolving Data Streams	16
6.5. Scalable Data Analysis	17
6.5.1. Scalable Mining of Small Visual Objects	17
6.5.2. Rare Events Identification for Large-Scale Applications	17
6.5.3. Large-scale content-based plants identification from social image data	18
7. Bilateral Contracts and Grants with Industry	18
7.1. Microsoft (2013-2017)	18
7.2. EDF R&D (2013-2014)	19
8. Partnerships and Cooperations	19
8.1. Regional Initiatives	19
8.1.1. Labex NUMEV, Montpellier	19
8.1.2. Institut de Biologie Computationnelle (IBC), Montpellier	19
8.2. National Initiatives	19

8.2.1.	ANR	19
8.2.2.	PIA	20
8.2.2.1.	Datascale (2013-2015), 250Keuros	20
8.2.2.2.	Xdata (2013-2015), 125Keuros	20
8.2.3.	Others	20
8.2.3.1.	RTRA PI@ntNet (2009-2013), 1Meuros	20
8.2.3.2.	CIFRE INA/Inria (2011-2013), 100Keuros	20
8.2.3.3.	CIFRE INA/Inria (2013-2016), 100Keuros	20
8.2.3.4.	CNRS INS2I Mastodons (2013), 30Keuros	21
8.3.	European Initiatives	21
8.4.	International Initiatives	21
8.4.1.	Inria Associate Teams	21
8.4.2.	International Benchmarks	22
8.4.3.	Inria International Partners	22
8.4.4.	Inria International Labs	22
8.4.5.	Participation In other International Programs	22
8.5.	International Research Visitors	23
8.5.1.	Visits of International Scientists	23
8.5.2.	Visits to International Teams	23
9.	Dissemination	23
9.1.	Scientific Animation	23
9.2.	Teaching - Supervision - Juries	25
9.2.1.	Teaching	25
9.2.2.	Supervision	25
9.2.3.	Juries	26
9.3.	Popularization	26
10.	Bibliography	26

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 Team: 2011 January 01, updated into Project-Team: 2012 January 01.

1. Members

Research Scientists

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

Faculty Member

Esther Pacitti [Univ. Montpellier 2, Associate Team Leader, Professor, HdR]

Engineers

Vera Bakic [Inria, Agropolis Fondation, from Jan 2013 until Sep 2013]
Emmanuel Castanier [Inria, X-Data PIA project]
Julien Champ [Inria and Inra, PI@ntNet and ARCAD projects]
Mathias Chouet [Inria, Agropolis Fondation, until Sep 2013]
Dimitri Dupuis [Inria, from Oct 2013]
Boyan Kolev [Inria, FP7 CoherentPaaS project, from Oct 2013]
Pierre Letessier [Inria, ANR OTMedia project, until May 2013]
Souheil Selmi [Inria, Agropolis Fondation, until Oct 2013]

PhD Students

Miguel Liroz-Gistau [Inria, CORDIS]
Ji Liu [Inria-MSR, from Oct 2013]
Yoann Couillec [Inria CORDIS (with the INDES team)]
Saber Salah [Inria, Hemera action]
Maximilien Servajean [Univ. Montpellier 2, NUMEV labex]
Naser Ayat [University of Amsterdam, Netherlands]
Jonas Dias [Universidade Federal de Rio de Janeiro, Brazil]

Post-Doctoral Fellows

Hervé Goëau [Inria, Agropolis Fondation, until Sep 2013]
Tristan Allard [Inria, EDF project, from Jan 2013]

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. Each discipline or organization tends to produce and manage its own data, in specific formats, with its own processes. Thus, 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: big data; 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 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

- BigdataNet: an associated team between Zenith and the distributed systems team of Profs. Amr El Abbadi and Divy Agrawal at University of California, Santa Barbara, since January 2013.
- Since October, Zenith participates in the European FP7 IP CoherentPaaS Project.

- The release of PlantNet iPhone App ¹, an image sharing and retrieval application for the identification of plants integrating several research contributions of Alexis Joly.

3. Research Program

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, big data, 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.

¹<https://itunes.apple.com/en/app/plantnet/id600547573>

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) [15]. 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 [11] 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 only 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 no 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.

3.4. Big Data

Big data has become a buzz word, with different meanings depending on your perspective, e.g. 100 terabytes is big for a transaction processing system, but small for a web search engine. It is also a moving target, as shown by two landmarks in DBMS products: the Teradata database machine in the 1980's and the Oracle Exadata database machine in 2010.

Although big data has been around for a long time, it is now more important than ever. We can see overwhelming amounts of data generated by all kinds of devices, networks and programs, e.g. sensors, mobile devices, internet, social networks, computer simulations, satellites, radiotelescopes, etc. Storage capacity has doubled every 3 years since 1980 with prices steadily going down (e.g. 1 Gigabyte for: 1M\$ in 1982, 1K\$ in 1995, 0.12\$ in 2011), making it affordable to keep more data. And massive data can produce high-value information and knowledge, which is critical for data analysis, decision support, forecasting, business intelligence, research, (data-intensive) science, etc.

The problem of big data has three main dimensions, quoted as the three big V's:

- Volume: refers to massive amounts of data, which makes it hard to store, manage, and analyze (big analytics);
- Velocity: refers to continuous data streams being produced, which makes it hard to perform online processing and analysis;
- Variety: refers to different data formats, different semantics, uncertain data, multiscale data, etc., which makes it hard to integrate and analyze.

There are also other V's like: validity (is the data correct and accurate?); veracity (are the results meaningful?); volatility (how long do you need to store this data?).

Current big data management (NoSQL) solutions have been designed for the cloud, as cloud and big data are synergistic. They 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 MapReduce and NoSQL technologies to fit our needs for scientific workflow management and scalable data analysis.

3.5. 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.

3.6. Big data 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.

3.7. 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 [14].** 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 has been 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 response time 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.

3.8. 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.
- **Distributed content-based retrieval.** Distributed content-based retrieval methods appeared recently as a promising solution to manage masses of data distributed over large 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 distributed 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. 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.

- **Personal health data analysis and privacy** The “Quantified Self” movement has gained a large popularity these past few years. Today, it is possible to acquire data on many domains related to personal data. For instance, one can collect data on her daily activities, habits or health. It is also possible to measure performances in sports. This can be done thanks to sensors, communicating devices or even connected glasses (as currently being developed by companies such as Google, for instance). Obviously, such data, once acquired, can lead to valuable knowledge for these domains. For people having a specific disease, it might be important to know if they belong to a specific category that needs particular care. For an individual, it can be interesting to find a category that corresponds to her performances in a specific sport and then adapt her training with an adequate program. Meanwhile, for privacy reasons, people will be reluctant to share their personal data and make them public. Therefore, it is important to provide them with solutions that can extract such knowledge from everybody’s data, while guaranteeing that their data won’t leave their computer and won’t be disclosed to anyone.
- **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 systems. 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. 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 and Platforms

5.1. WebSmatch (Web Schema Matching)

Participants: Emmanuel Castanier, Rémi Coletta, 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. SON (Shared-data Overlay Network)

Participants: Esther Pacitti, Didier Parigot [contact], Patrick Valduriez.

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.

5.3. P2Prec (P2P recommendation service)

Participants: Esther Pacitti [contact], Didier Parigot, Maximilien Servajean.

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 SON.

5.4. ProbDB (Probabilistic Database)

Participants: Reza Akbarinia [contact], Patrick Valduriez.

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.

5.5. Pl@ntNet-mobile

Participants: Vera Bakic, Souheil Selmi, Hervé Goëau, Alexis Joly [contact].

URL: <http://goo.gl/CpSrr3>

PI@ntNet-mobile is an image sharing and retrieval application for the identification of plants built in the continuity of the former web application PI@ntNet-Identify ² (presented in last year activity report). It is developed in the context of the PI@ntNet project that involves four French research organisations (Inria, Cirad, INRA, IRD) and the members of Tela Botanica social network. The key feature of this free app is to help identifying plant species from photographs, through a server-side visual search engine based on several results of ZENITH team on content-based information retrieval. Since its first release in March 2013 on the apple store, the application was downloaded by around 80K users in about 150 countries (between 200 and 2000 active users daily with peaks occurring during the week-ends). The collaborative training set that allows the content-based identification is continuously enriched by the users of the application and the members of Tela Botanica social network. At the time of writing, it includes about 80K images covering more than 3500 French plant species about 2/3 of the whole French flora (this is actually the widest identification tool built anytime).

5.6. PI@ntNet-DataManager

Participants: Mathias Chouet [contact], Alexis Joly.

PI@ntNet-DataManager ³ is a software dedicated to managing and sharing distributed heterogeneous botanical data. It is developed jointly by Zenith, the AMAP UMR team (CIRAD) and Telabotanica 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.

5.7. SnoopIm

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

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

SnoopIm is a content-based search engine allowing to discover and retrieve small visual patterns or objects in large collections of pictures and to derive statistics from them (frequency, visual cover, size variations, etc.). It is implemented in Javascript on top of a C++ library developed in collaboration with INA ⁴. The software is used at INA by archivists and sociologists in the context of the Transmedia Observatory project. It is also being experimented in several contexts including a logo retrieval application set up in collaboration with the French Press Agency, an experimental plant identification tool mixing textual and visual information retrieval (in the context of the PI@ntNet project) and a research project on high-throughput analysis of root architecture images.

5.8. SciFloware

Participants: Dimitri Dupuis, Didier Parigot [contact], Patrick Valduriez.

URL: <http://www-sop.inria.fr/members/Didier.Parigot/pmwiki/Scifloware>

²<http://identify.plantnet-project.org>

³<http://data.plantnet-project.org/>

⁴<http://www.ina-sup.com/>

SciFloware is an action of technology development (ADT Inria) with the goal of developing a middleware for the execution of scientific workflows in a distributed and parallel way. It capitalizes on our experience with SON and an innovative algebraic approach to the management of scientific workflows. SciFloware provides a development environment and a runtime environment for scientific workflows, interoperable with existing systems. We will validate SciFloware with workflows for analyzing biological data provided by our partners CIRAD, INRA and IRD.

6. New Results

6.1. Big Data Integration

6.1.1. Probabilistic Data Integration

Participants: Reza Akbarinia, Naser Ayat, Patrick Valduriez.

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.

An important problem that arises in big data integration is that of Entity Resolution (ER). ER is the process of identifying tuples that represent the same real-world entity. The problem of *entity resolution over probabilistic data* (which we call ERPD) arises in many distributed application domains that have to deal with probabilistic data, ranging from sensor databases to scientific data management. The ERPD problem can be formally defined as follows. Let e be an uncertain entity represented by multiple possible alternatives, i.e. tuples, each with a membership probability. Let D be an uncertain database composed of a set of tuples each associated with a membership probability. Then, given e , D , and a similarity function F , the problem is to find the entity-tuple pair (t, t_i) (where $t \in e, t_i \in D$) such that (t, t_i) has the highest cumulative probability to be the most similar in all possible worlds. This entity-tuple pair is called the *most probable match pair* of e and D , denoted by $MPMP(e, D)$.

Many real-life applications produce uncertain data distributed among a number of databases. Dealing with the ERPD problem for distributed data is quite important for such applications. A straightforward approach for answering distributed ERPD queries is to ask all distributed nodes to send their databases to a central node that deals with the problem of ER by using one of the existing centralized solutions. However, this approach is very expensive and does not scale well neither in the size of databases, nor in the number of nodes.

In [20], we proposed FD (Fully Distributed), a decentralized algorithm for dealing with the ERPD problem over distributed data, with the goal of minimizing bandwidth usage and reducing processing time. It has the following salient features. First, it uses the novel concepts of *Potential* and *essential-set* to prune data at local nodes. This leads to a significant reduction of bandwidth usage compared to the baseline approaches. Second, its execution is completely distributed and does not depend on the existence of certain nodes. We validated FD through implementation over a 75-node cluster and simulation using both synthetic and real-world data. The results show very good performance, in terms of bandwidth usage and response time.

6.1.2. Open Data Integration

Participants: Emmanuel Castanier, Patrick Valduriez.

Working with open data sources can yield high value information but raises major problems in terms of metadata extraction, data source integration and visualization. For instance, Data Publica provides more than 12 000 files of public data. However, even though data formats become richer and richer in terms of semantics and expressivity (e.g. RDF), most data producers do not use them much in practice, because they require too much upfront work, and keep using simpler tools like Excel. Unfortunately, no integration tool is able to deal in an effective way with spreadsheets. Only few initiatives (OpenII and Google Refine) deal with Excel files. However, their importers are very simple and make some strict restrictions over the input spread-sheets.

In [31], we describe a demonstration of WebSmatch, a flexible environment for Web data integration. 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. Pricing Integrated Data

Participant: Patrick Valduriez.

Data is a modern commodity, being bought and sold. Electronic data market places and independent vendors integrate data and organize their online distribution. Yet the pricing models in use either focus on the usage of computing resources, or are proprietary, opaque, most likely ad hoc, and not conducive of a healthy commodity market dynamics. In [39], we propose a generic data pricing model that is based on minimal provenance, i.e. minimal sets of tuples contributing to the result of a query. We show that the proposed model fulfills desirable properties such as contribution monotonicity, bounded-price and contribution arbitrage-freedom. We present a baseline algorithm to compute the exact price of a query based on our pricing model. We show that the problem is NP-hard. We therefore devise, present and compare several heuristics. We conduct a comprehensive experimental study to show their effectiveness and efficiency.

In most data markets, prices are prescribed and accuracy is determined by the data. Instead, we consider a model in which accuracy can be traded for discounted prices: “what you pay for is what you get”. The data market model consists of data consumers, data providers and data market owners. The data market owners are brokers between the data providers and data consumers. A data consumer proposes a price for the data that she requests. If the price is less than the price set by the data provider, then she gets an approximate value. The data market owners negotiate the pricing schemes with the data providers. They implement these schemes for the computation of the discounted approximate values. In [38], we propose a theoretical and practical pricing framework with its algorithms for the above mechanism. In this framework, the value published is randomly determined from a probability distribution. The distribution is computed such that its distance to the actual value is commensurate to the discount. The published value comes with a guarantee on the probability to be the exact value. The probability is also commensurate to the discount. We present and formalize the principles that a healthy data market should meet for such a transaction. We define two ancillary functions and describe the algorithms that compute the approximate value from the proposed price using these functions. We prove that the functions and the algorithm meet the required principles.

6.2. Distributed Indexing and Searching

6.2.1. P2P Search and Recommendation

Participants: Esther Pacitti, Maximilien Servajean.

In crossdiscipline domains, users belonging to different communities produce various scientific material that they own, share, or endorse. In that context, we are interested in querying and recommending scientific material in the form of documents. Such documents cover various topics such as models for plant phenotyping, statistics on specific kinds of plants, or biological experiments.

In [40], we investigate profile diversity, a novel idea in searching scientific documents. Combining keyword relevance with popularity in a scoring function has been the subject of different forms of social relevance. On the other hand, content diversity has been thoroughly studied in search and advertising, database queries, and recommendations.

We introduce profile diversity for scientific document search as a complement to traditional content diversity. Profile diversity combines the discipline and communities to which a user belongs. We propose an adaptation of Fagin’s threshold-based algorithms to return the most relevant and most popular documents that satisfy content and profile diversities. To validate our scoring function, DivRSci, we ran experiments that use two benchmarks: a realistic benchmark with scientists and TREC’09. We show that DivRSci presents the best compromise between all requirements we have identified. DivRSci also shows to be the best generating list of

inter-disciplinary and inter-community documents. Finally, it yields very good gains (by a factor of 6), suited for profile diversification

6.2.2. *Spatial Queries in Wireless Data Broadcasting*

Participant: Patrick Valduriez.

The main requirements for spatial query processing via mobile terminals include rapid and accurate searching and low energy consumption. Most location-based services (LBSs) are provided using an on-demand method, which is suitable for light-loaded systems where contention for wireless channels and server processing is not severe. However, as the number of users of SBSs increases, performance deteriorates rapidly since the servers' capability to process queries is limited. Furthermore, the response time of a query may significantly increase with the concentration of users' queries in a server at the same time. That is because the server has to check the locations of users and potential objects for the final result and then individually send answers to clients via a point-to-point channel. At this time, an inefficient structure of spatial index and searching algorithm may incur an extremely large access latency.

To address this problem, we propose in [27] the Hierarchical Grid Index (HGI), which provides a light-weight sequential location-based index structure for efficient SBSs. We minimize the index size through the use of hierarchical location-based identifications. And we support efficient query processing in broadcasting environments through sequential data transfer and search based on the object locations. We also propose Top-Down Search and Reduction-Counter Search algorithms for efficient searching and query processing. HGI has a simple structure through elimination of replication pointers and is therefore suitable for broadcasting environments with one-dimensional characteristics, thus enabling rapid and accurate spatial search by reducing redundant data. Our performance evaluation shows that our proposed index and algorithms are accurate and fast and support efficient spatial query processing.

6.3. Big Data Analysis

6.3.1. *Big Data Analysis using Algebraic Workflows*

Participants: Jonas Dias, Patrick Valduriez.

Analyzing big data requires the support of dataflows with many activities to extract and explore relevant information from the data. Recent approaches such as Pig Latin propose a high-level language to model such dataflows. However, the dataflow execution is typically delegated to a MapReduce implementation such as Hadoop, which does not follow an algebraic approach, thus it cannot take advantage of the optimization opportunities of PigLatin algebra.

In [35], we propose an approach for big data analysis based on algebraic workflows, which yields optimization and parallel execution of activities and supports user steering using provenance queries. We illustrate how a big data processing dataflow can be modeled using the algebra. Through an experimental evaluation using real datasets and the execution of the dataflow with Chiron, an engine that supports our algebra, we show that our approach yields performance gains of up to 19.6% using algebraic optimizations in the dataflow and up to 39% of time saved on a user steering scenario.

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

6.3.2. *Big Data Partitioning*

Participants: Reza Akbarinia, Miguel Liroz, Esther Pacitti, Patrick Valduriez.

The amount of data that is captured or generated by modern computing devices has augmented exponentially over the last years. For processing this *big data*, parallel computing has been a major solution in both industry and research. This is why, the MapReduce framework, which provides automatic distribution parallelization and fault-tolerance in a transparent way over lowcost machines, has become one of the standards in big data analysis.

For processing a big dataset over a cluster of nodes, one main step is data partitioning (or fragmentation) to divide the dataset to the nodes. In our team, we study the problem of data partitioning in two different contexts: (1) in scientific databases that are continuously growing and (2) in the MapReduce framework. In both cases, we propose automatic approaches, which are performed transparently to the users, in order to free them from the burden of complex partitioning.

In [25], we consider applications with very large databases, where data items are continuously appended. Thus, the development of efficient data partitioning is one of the main requirements to yield good performance. In particular, this problem is harder in the case of some scientific databases, such as astronomical catalogs. The complexity of the schema limits the applicability of traditional automatic approaches based on the basic partitioning techniques. The high dynamicity makes the usage of graph-based approaches impractical, as they require to consider the whole dataset in order to come up with a good partitioning scheme. In our work, we propose *DynPart* and *DynPartGroup*, two dynamic partitioning algorithms for continuously growing databases [25]. These algorithms efficiently adapt the data partitioning to the arrival of new data elements by taking into account the affinity of new data with queries and fragments. In contrast to existing static approaches, our approach offers constant execution time, no matter the size of the database, while obtaining very good partitioning efficiency. We validate our solution through experimentation over real-world data; the results show its effectiveness.

In [37] and [43], we address the problem of high data transfers in MapReduce, and propose a technique that repartitions tuples of the input datasets. Our technique optimizes the distribution of key-values over mappers, and increases the data locality in reduce tasks. It captures the relationships between input tuples and intermediate keys by monitoring the execution of a set of MapReduce jobs which are representative of the workload. Then, based on those relationships, it assigns input tuples to the appropriate chunks. With this data repartitioning and a smart scheduling of reducer tasks, our approach significantly contributes to the reduction of transferred data between mappers and reducers in job executions. We evaluate our approach through experimentation in a Hadoop deployment on top of Grid5000 using standard benchmarks. The results show high reduction in data transfer during the shuffle phase compared to Native Hadoop.

6.4. Data Stream Mining

6.4.1. Mining Uncertain Data Streams

Participants: Reza Akbarinia, Florent Masseglia.

Discovering Probabilistic Frequent Itemsets (PFI) is very challenging since algorithms designed for deterministic data are not applicable in probabilistic data. The problem is even more difficult for probabilistic data streams where massive frequent updates need to be taken into account while respecting data stream constraints. In [28], we propose FEMP (Fast and Exact Mining of Probabilistic data streams), the first solution for exact PFI mining in data streams with sliding windows. FEMP allows updating the frequentness probability of an itemset whenever a transaction is added or removed from the observation window. Using these update operations, we are able to extract PFI in sliding windows with very low response times. Furthermore, our method is exact, meaning that we are able to discover the exact probabilistic frequentness distribution function for any monitored itemset, at any time. We implemented FEMP and conducted an extensive experimental evaluation over synthetic and real-world data sets; the results illustrate its very good performance.

6.4.2. Itemset Mining over Tuple-Evolving Data Streams

Participant: Florent Masseglia.

In many data streaming applications today, tuples inside the streams may get revised over time. This type of data stream brings new issues and challenges to the data mining tasks. In [42] we present a theoretical analysis for mining frequent itemsets from sliding windows over such data. We define conditions that determine whether an infrequent itemset will become frequent when some existing tuples inside the streams have been updated. We design simple but effective structures for managing both the evolving tuples and the candidate frequent itemsets. Moreover, we provide a novel verification method that efficiently computes the counts of candidate itemsets. Experiments on real-world datasets show the efficiency and effectiveness of our proposed method.

6.5. Scalable Data Analysis

6.5.1. 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 small compared to the whole image. In a previous work, we revisited formally 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. This year, we continued working on the subject by improving our high-dimensional data hashing strategy, that works first at the visual level, and then at the geometric level. We conducted new experiments on a dedicated evaluation dataset⁵ and we did show that our the recall or our approach definitely outperforms the reference method [46].

Based on this contribution, we then address the problem of suggesting object-based visual queries in a multimedia search engine [22], [36]. State-of-the-art visual search systems are usually based on the query-by-window paradigm: a user selects any image region containing an object of interest and the system returns a ranked list of images that are likely to contain other instances of the query object. User's perception of these tools is however affected by the fact that many submitted queries actually return nothing or only junk results (complex non-rigid objects, higher-level visual concepts, etc.). In [22], we addressed the problem of suggesting only the object's queries that actually contain relevant matches in the dataset. This requires to first discover accurate object's clusters in the dataset (as an offline process); and then to select the most relevant objects according to user's intent (as an on-line process). We therefore introduce a new object's instances clustering framework based on a bipartite shared-neighbours clustering algorithm that is used to gather object's seeds discovered by our visual mining method. Shared nearest neighbours methods were not studied beforehand in the case of bipartite graphs and never used in the context of object discovery. Experiments show that this new method outperforms state-of-the-art object mining and retrieval results on the Oxford Building dataset. We finally describe two real-word object-based visual query suggestion scenarios using the proposed framework and show examples of suggested object queries. A demo was presented at ACM Multimedia 2013 [36].

This method was finally integrated within a visual-based media event detection system in the scope of a French project called the Transmedia Observatory [33]. 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.

6.5.2. Rare Events Identification for Large-Scale Applications

Participant: Florent Masseglia.

While significant work in data mining has been dedicated to the detection of single outliers in the data, less research has approached the problem of isolating a group of outliers, i.e. rare events representing micro-clusters of less – or significantly less – than 1% of the whole dataset. This research issue is critical for example in medical applications. The problem is difficult to handle as it lies at the frontier between outlier detection and clustering and distinguishes by a clear challenge to avoid missing true positives. In [41], we address this challenge and propose a novel two-stage framework, based on a backward approach, to isolate abnormal groups of events in large datasets. The key of our backward approach is to first identify the core of the dense regions and then gradually augment them based on a density-driven condition. The framework outputs a small subset of the dataset containing both rare events and outliers. We tested our framework on a biomedical application to find micro-clusters of pathological cells. The comparison against two common clustering (DBSCAN) and outlier detection (LOF) algorithms show that our approach is a very efficient alternative to the detection of rare events – generally a recall of 100% and a higher precision, positively correlated with the size of the rare event – while also providing a $\mathcal{O}(N)$ solution to the existing algorithms dominated by a $\mathcal{O}(N^2)$ complexity.

⁵<http://www-sop.inria.fr/members/Alexis.Joly/BelgaLogos/FlickrBelgaLogos.html>

6.5.3. Large-scale content-based plants identification from social image data

Participants: Hervé Goëau, Alexis Joly, Julien Champ, Saloua Litayem.

Speeding up the collection and integration of raw botanical observation data is a crucial step towards a sustainable development of agriculture and the conservation of biodiversity. Initiated in the context of a citizen sciences project in collaboration with the botanists of the AMAP UMR team and Tela Botanica social network, the overall contribution of this work [23] is an innovative collaborative workflow focused on image-based plant identification as a mean to enlist new contributors and facilitate access to botanical data. Since 2010, hundreds of thousands of geo-tagged and dated plant photographs were collected and revised by hundreds of novice, amateur and expert botanists of a specialized social network. An image-based identification tool - available as both a web and a mobile application - is synchronized with that growing data and allows any user to query or enrich the system with new observations. Extensive experiments of the visual search engine as well as system-oriented and user-oriented evaluations of the application did show that it is very helpful to determine a plant among hundreds or thousands of species [23]. As a concrete result, more than 80K people in about 150 countries did download the iPhone end point of the application [32].

From a data management and data analysis perspective, our main contribution concerns the scalability of the system. At the time of writing, the content-based search engine actually works on 120K images covering more than 5000 species (which already makes it the largest identification tool built anytime). The resulting training dataset contains several hundreds of millions feature vectors, each with several hundreds of float attributes (i.e. high-dimensional feature vectors describing the visual content). At query time, thousands of such feature vectors are extracted from the query pictures and have to be searched online in the training set to find the most similar pictures. The underlying search of approximate nearest neighbors is speed-up thanks to a data-dependent high-dimensional hashing framework based on Random Maximum Margin Hashing (RMMH), a new hash function family that we introduced in 2011. RMMH is used for both compressing the original feature vectors into compact binary hash codes and for partitioning the data into a well balanced hash table. Search is then performed through adaptive multi-probe accesses in the hash table and a top-k search refinement step on the full binary hash codes. Last improvements brought in 2013 include a multi-threaded version of the search, the use of a probabilistic asymmetric distance instead of the Hamming distance and the integration of a query optimization training stage in the compressed feature space instead of the original space. A beta version of Pl@ntNet visual search engine based on these new contributions is currently being tested and is about 8 times faster than the one used in production.

Besides scalability and efficiency, we also did work on improving the identification performances of the system [29]. We notably improved the quality of the top-K returned images by weighting each match according to its Hamming distance to the query rather than using a simple vote. We then improved the multi-cue fusion strategy by indexing separately each type of visual features rather than concatenating them in an early phase. We finally did train the optimal selection of features for each of the considered plant organ (flower, leaf, bark, fruit). Beyond the use of the visual content itself, we explored the usefulness of associated metadata and we did prove that some of them like the date can improve the identification performances (contrary to the geo-coordinates that surprisingly degraded the results). Overall, as a result of our participation to ImageCLEF plant identification benchmark [34], we obtained the second best run among 12 international groups and a total of 33 submitted runs.

7. Bilateral Contracts and Grants with Industry

7.1. Microsoft (2013-2017)

Participants: Ji Liu, Esther Pacitti, Patrick Valduriez.

This joint project is on advanced data storage and processing for cloud workflows with the Kerdata team in the context of the Joint Inria – Microsoft Research Centre. The project addresses the problem of advanced data storage and processing for supporting scientific workflows in the cloud. The goal is to design and implement a framework for the efficient processing of scientific workflows in clouds. The validation will be performed using synthetic benchmarks and real-life applications from bioinformatics: first on the Grid5000 platform in a preliminary phase, then on the Microsoft Azure cloud environment.

7.2. EDF R&D (2013-2014)

Participants: Tristand Allard, Florent Masegla, Esther Pacitti.

This project aims at developing new data mining techniques for P2P networks. The main goal is to preserve data privacy, while achieving good performance of analysis processes on the tackled data. More precisely, each participant in the P2P network has its own individual data (e.g. results of experiments for a scientific partner) and all the participants would like to acquire knowledge computed on the whole dataset (i.e. the union of all the individual data on the peers). Meanwhile, participants want a guarantee that no other participant will be able to see their data. The P2P protocol we are developing will then be able to extract knowledge from the whole set of distributed data, while avoiding centralization, and guaranteeing data privacy for all peers.

8. Partnerships and Cooperations

8.1. Regional Initiatives

8.1.1. Labex NUMEV, Montpellier

URL: <http://www.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.

8.2. National Initiatives

8.2.1. ANR

8.2.1.1. 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. PIA

8.2.2.1. Datascale (2013-2015), 250Keuros

Participants: Reza Akbarinia, Florent Masseglia, Saber Salah, Patrick Valduriez.

The Datascale project is a “projet investissements d’avenir” on big data with Bull (leader), CEA, ActiveEon SAS, Armadillo, Twenga, IPGP, Xedix and Inria (Zenith) . The goal of the project is to develop the essential technologies for big data, including efficient data management, software architecture and database architecture, and demonstrate their scalability with representative applications. In this project, the Zenith team works on data mining with Hadoop MapReduce.

8.2.2.2. Xdata (2013-2015), 125Keuros

Participants: Emmanuel Castanier, Patrick Valduriez.

The X-data project is a “projet investissements d’avenir” on big data with Data Publica (leader), Orange, La Poste, EDF, Cinequant, Hurence and Inria (Indes, Planete and Zenith) . The goal of the project is to develop a big data platform with various tools and services to integrate open data and partners’s private data for analyzing the location, density and consuming of individuals and organizations in terms of energy and services. In this project, the Zenith team heads the workpackage on data integration.

8.2.3. Others

8.2.3.1. RTRA PI@ntNet (2009-2013), 1Meuros

Participants: Alexis Joly, Hervé Goëau, Julien Champ, 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 ⁶, IMEDIA project team at Inria, and the French botanical network TelaBotanica ⁷), 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.

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

Participants: Alexis Joly, Pierre Letessier.

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

8.2.3.3. CIFRE INA/Inria (2013-2016), 100Keuros

Participants: Alexis Joly, Valentin Leveau, Patrick Valduriez.

⁶<http://amap.cirad.fr/en/>

⁷<http://www.tela-botanica.org/>

This CIFRE contract with INA allows funding a 3-years PhD (Valentin Leveau). This PhD addresses research challenges related to large-scale supervised content-based retrieval notably in distributed environments.

8.2.3.4. CNRS INS2I Mastodons (2013), 30Keuros

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

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

8.3.1.1. CoherentPaaS

Project title: A Coherent and Rich Platform as a Service with a Common Programming Model

Instrument: Integrated Project

Duration: 2013 - 2016

Total funding: 5 Meuros (Zenith: 500Keuros)

Coordinator: U. Madrid, Spain

Partner: FORTH (Greece), ICCS (Greece), INESC (Portugal) and the companies MonetDB (Netherlands), QuartetFS (France), Sparsity (Spain), Neurocom (Greece), Portugal Telecom (Portugal).

Inria contact: Patrick Valduriez

Accessing and managing large amounts of data is becoming a major obstacle to developing new cloud applications and services with correct semantics, requiring tremendous programming effort and expertise. CoherentPaaS addresses this issue in the cloud PaaS landscape by developing a PaaS that incorporates a rich and diverse set of cloud data management technologies, including no SQL data stores, such as key-value data stores and graph databases, SQL data stores, such as in-memory and column-oriented databases, hybrid systems, such as SQL engines on top on key-value data stores, and complex event processing data management systems. It uses a common query language to unify the programming models of all systems under a single paradigm and provides holistic coherence across data stores using a scalable, transactional management system. CoherentPaaS will dramatically reduce the effort required to build and the quality of the resulting cloud applications using multiple cloud data management technologies via a single query language, a uniform programming model, and ACID-based global transactional semantics. CoherentPaaS will design and build a working prototype and will validate the proposed technology with real-life use cases. In this project, Zenith is in charge of designing an SQL-like query language to query multiple databases (SQL, NoSQL) in a cloud and implementing a compiler/optimizer and query engine for that language.

8.4. International Initiatives

8.4.1. Inria Associate Teams

8.4.1.1. BIGDATANET

Title: A hybrid P2P/cloud for big data

Inria principal investigator: Patrick Valduriez

International Partner (Institution - Laboratory - Researcher):

University of California at Santa Barbara (United States) - Distributed Systems Lab. - Amr El Abbadi and Divy Agrawal

Duration: 2013 -2015

See also: <https://team.inria.fr/zenith/projects/international-projects/bigdatanet/>

The main objective of this research and scientific collaboration is to develop a hybrid architecture of a computational platform that leverages the cloud computing and the P2P computing paradigms. The resulting architecture will enable scalable data management and data analysis infrastructures that can be used to host a variety of next-generation applications that benefit from computing, storage, and networking resources that exist not only at the network core (i.e., data-centers) but also at the network edge (i.e., machines at the user level as well as machines available in CDNs – content distribution networks hosted in ISPs).

8.4.2. International Benchmarks

8.4.2.1. ImageCLEF

Title: The CLEF Cross Language Image Retrieval Track

Inria principal investigator: Alexis Joly

International Partners (Institution - Laboratory - Researcher): HES-SO (Switzerland), Yahoo! Research (Spain), IBrandenburg Technical University (Germany), diap Research Institute (Switzerland), University of Alicante (Spain), Universidad Politécnic de Valencia (Spain), UMR AMAP (France)

Duration: 2011 -2013

See also: <http://www.imageclef.org>

Since its first edition in 2003, ImageCLEF has become one of the key initiatives promoting the benchmark evaluation of algorithms for the cross-language annotation and retrieval of images in various domains, such as public and personal images, to data acquired by mobile robot platforms and botanic collections. Over the years, by providing new data collections and challenging tasks to the community of interest, the ImageCLEF lab has achieved a unique position in the multi lingual image annotation and retrieval research landscape. As an illustration of its impact, the 2013 edition attracted more than 100 registered team world-wide and 42 of them did cross the finish line by submitting runs of their system [30]. Zenith, through the implication of Alexis Joly and Hervé Goëau, is one of the co-organizer of the lab and the initiator of the plant retrieval task since 2011

8.4.3. Inria International Partners

8.4.3.1. Informal International Partners

We have regular scientific relationships with research laboratories in

- North America: Univ. of Waterloo (Tamer Özsu), Mc Gill, Montreal (Bettina Kemme).
- 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 Muñoz)

8.4.4. Inria International Labs

The Bigdatanet associated team takes part in the Inria@SiliconValley lab.

8.4.5. Participation In other International Programs

We are involved in the following international actions:

- 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.

8.5. International Research Visitors

8.5.1. Visits of International Scientists

Dennis Shasha (NYU, USA) gave a seminar on “Storing Clocked Programs Inside DNA: A Simplifying Framework for Nanocomputing” in january.

Prof. Marta Mattoso (UFRJ, Rio de Janeiro) gave a seminar in the context of IBC on “Big Data Workflows – how provenance can help” in march and “Algebraic Dataflows for Big Data Analysis” in november.

Aravind Venkatesan (NTNU, Trondheim, Norway) gave a seminar in the context of IBC on “Bringing Semantic Web Technology to the Lab Bench” in october.

Sihem Amer-Yahia (LIG) gave a seminar on “New Perspectives in Social Data Management” in november.

Themis Palpanas (Univ. Trento, Italy) gave a seminar on “Enabling Exploratory Analysis on Very Large Scientific Data” in december.

8.5.2. Visits to International Teams

Reza Akbarinia and Florent Massglaia visited UCSB (Prof. Divy Agrawal and Amr El Abbadi) in may. Esther Pacitti and Patrick Valduriez also visited UCSB and Lawrence Berkeley Laboratory, Berkeley (Dr. Arie Shoshani and Deb. Agrawal) in june.

9. Dissemination

9.1. Scientific Animation

Participation in the editorial board of scientific journals:

- VLDB Journal: P. Valduriez.
- Journal of Transactions on Large Scale Data and Knowledge Centered Systems, R. Akbarinia.
- 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 to the organization of conferences and workshops:

- Alexis Joly and Julien Champs were chair of the video program track of ACM Multimedia 2013
- Alexis Joly co-organized ImageCLEF 2013 workshop (within CLEF 2013 conference)
- Alexis Joly co-organized a special session on Image Processing and Pattern Recognition for Ecological Applications at ICIP 2013
- Alexis Joly organized a workshop in the context of the French ISIS working group on large-scale multimedia data linking and mining (Liage, structuration et fouille de grandes données multimédia)
- Alexis Joly co-organized the OT-Media Workshop on Big Data for Media Analysis

Participation in conference program committees :

- ACM SIGMOD Conf. 2014: R. Akbarinia, P. Valduriez (area chair)
- Int. Conf. on Extending DataBase Technologies (EDBT), 2014: R. Akbarinia, A. Joly, F. Massegli
- IEEE Int. Conf. on Data Engineering (ICDE) 2014: P. Valduriez
- Int. Conf. on Middleware 2012: E. Pacitti
- IDEAS 2013: P. Valduriez
- Brazilian e-Science Workshop 2013: E. Pacitti
- Journées Bases de Données Avancées (BDA) 2013: A. Joly
- Int. Workshop on Open Data, 2013: P. Valduriez
- ACM Multimedia, 2013: A. Joly
- ACM Int. Conf. on Multimedia Retrieval (ICMR), 2013: A. Joly
- CLEF 2013, 2014: A. Joly
- Int. Conf. on Multimedia Modeling 2013: Alexis Joly
- ACM Symposium On Applied Computing (ACM SAC, Data Stream track), 2013: F. Massegli
- Conférence Internationale Francophone sur l'Extraction et la Gestion de Connaissance (EGC), 2013: F. Massegli
- IEEE Int. Conf. on Data Mining (ICDM), 2013: F. Massegli
- Int. Pacific-Asia Conf. on Knowledge Discovery and Data Mining (PAKDD), 2013: F. Massegli
- Pacific-Rim Conference on Multimedia (PCM), 2013: A. Joly
- International Conference on Similarity Search and Applications (SISAP), 2013: A. Joly
- Ubiquitous Data Mining (UDM), held by the International Joint Conference on Artificial Intelligence (IJCAI) 2013 : F. Massegli
- Extraction et Gestion des Connaissances : F. Massegli

Journal reviewing:

- VLDB Journal: A. Joly
- IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI): A. Joly
- IEEE transactions on multimedia: A. Joly
- Journal of Mathematical Imaging and Vision (JMIV): A. Joly
- Transactions on Knowledge and Data Engineering: A. Joly
- EURASIP Journal on Image and Video Processing: A. Joly
- Computer Vision and Image Understanding Journal: A. Joly
- ACM Trans. on Database Systems: E. Pacitti
- Knowledge and Information Systems (KAIS): F. Massegli

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.
- Patrick Valduriez gave a seminar on “Parallel Techniques for Big Data” at IBC, Montpellier, and a tutorial on the same topic at BDA 2013. He also gave a seminar in the Morgenstern series in Sophia-Antipolis on “Cloud and Bigdata: a perfect marriage?”.

Other activities (international):

- Alexis Joly was one of the 6 panelists of ACM Multimedia 2013 panel on cross-media analysis and mining.
- Reza Akbarinia and Florent Masseglia gave a talk on BigdataNet project whose objective is to develop a hybrid P2P-Cloud platform, at the Berkeley Inria Stanford (BIS2013) workshop at Stanford in May.
- Esther Pacitti and Patrick Valduriez gave talks on search and recommendation, and scientific workflow management, respectively, at LBL, Berkeley, in June.

9.2. Teaching - Supervision - Juries

9.2.1. Teaching

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

Reza Akbarinia:

Master Research: Large scale data management, 7.5h, level M2, Faculty of Science, UM2

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

Licence: Computing Tools, 54h, Level L3, Faculty of Science, UM2

Florent Masseglia:

Master Research: Large scale data management, 3h, 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

Industry internship committee, 50h, level M2, Polytech' Montpellier

Master Research: Large scale data management, 4,5h, level M2, Faculty of Science, UM2

Didier Parigot:

Master Research: Large scale data management, 9h, level M2, Faculty of Science, UM2

Patrick Valduriez:

Master Research: Large scale data management, 12h, level M2, Faculty of Science, UM2

Professional: Distributed Information Systems, 50h, level M2, Capgemini Institut

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

Alexis Joly:

Master Research: Large scale data management, 6h, level M2, Faculty of Science, UM2

9.2.2. Supervision

- 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 : Ji Liu, Scientific Workflows in Multisite Cloud, started oct. 2013, UM2, Advisors: Esther Pacitti and Patrick Valduriez
- PhD in progress : Saber Salah, Optimizing a Cloud for Data Mining Primitives, started nov. 2012, UM2, Advisor: Florent Masseglia, 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
- PhD in progress : Valentin Leveau, Supervised content-based information retrieval in big multimedia data, started April 2013, University of Montpellier, Advisor: Patrick Valduriez, co-advisor: Alexis Joly and Olivier Buisson

9.2.3. Juries

Members of the team participated to the following Ph.D. committees:

- R. Akbarinia: Miguel Liroz (UM2);
- F. Masseglia: Cassio Melo (Ecole Centrale Paris, committee chair), Yasin Amanullah (Polytech Nantes, Reviewer);
- E. Pacitti: Mohamed Reda Bouadjenek (UVSQ), Miguel Liroz (UM2), Housseem Chihoub (Univ. Rennes 1) ;
- P. Valduriez: Jonas Dias (UFRJ, Rio de Janeiro), Gylfi por Gudmundsson (Univ. Rennes 1, Reviewer), Miguel Liroz (UM2), Asterios Katsifodimos (Univ. Paris Sud), Viet-Trung Tran (ENS Rennes).

9.3. Popularization

Zenith participated to the following events in France:

- “La fête de la science” (Montpellier), with the animation of a stand at Genopolys (a science village). We received 7 groups from Oct. 10 to Oct. 11, and then a wider audience on Oct. 12. With this stand, we have participated in the effort of relaunching LIRMM’s activities for this major event in science popularization.
- “Futur en Seine 2013” (Paris, <http://www.futur-en-seine.fr/fens2013/>) with the co-animation of the stand for the presentation of Interstices’ activities.
- annual seminar of the “mecsci” workgroup of Inria, Paris (CNAM) on January 30 and 31 [47].

F. Masseglia is a member of the editorial committee, and contributor, of interstices (<https://interstices.info>), and a member of the scientific committee for the edition of “Datagramme”, the game from Inria on science popularization. He has given a 3 days training to elementary school teachers. This event, called “Graines de sciences” (Ecole de Physique des Houches, <http://www.fondation-lamap.org/fr/graines-de-sciences>) was funded by “La main à la pâte” (<http://www.fondation-lamap.org/>), a French foundation dedicated to science popularization based on manual activities. The goal, here, was to give a 3h training to groups of a dozen of teachers, by mainly using games and activities they would be able to transfer in their classrooms.

A. Joly presented the PI@ntNet project and the resulting iPhone application in several events including the "Salon de l’agriculture 2013", "Les assises de la biodiversité", "Journée thématique du CNRS autour des usages mobiles de l’information scientifique dans l’Enseignement Supérieur et la Recherche", "Rencontres Agropolis", "journée des nouveaux arrivants de l’Inria", "journées INRA/Inria", etc. Many news media articles, blogs, TV news and twitts also reported on PI@ntNet (typing "PI@ntNet" on Google is the best way to get most of them). Finally a movie was realized in collaboration with Interstices (⁸).

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>

⁸<http://videotheque.inria.fr/videotheque/media/browse/2>

- [2] R. AKBARINIA, F. MASSEGLIA. *Fast and Exact Mining of Probabilistic Data Streams*, in "PKDD'2013: European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases", Prague, Czech Republic, Lecture Notes in Computer Science, Springer, 2013, pp. 493-508 [DOI : 10.1007/978-3-642-40988-2_32], <http://hal.inria.fr/lirmm-00838618>
- [3] R. AKBARINIA, E. PACITTI, P. VALDURIEZ. *Best Position Algorithms for Efficient Top-k Query Processing*, in "Information Systems", 2011, vol. 36, n^o 6, pp. 973-989, <http://hal.inria.fr/lirmm-00607882/en>
- [4] R. AKBARINIA, P. VALDURIEZ, G. VERGER. *Efficient Evaluation of SUM Queries Over Probabilistic Data*, in "IEEE Transactions on Knowledge and Data Engineering", 2013, vol. 25, n^o 4, pp. 764-775, <http://hal.inria.fr/lirmm-00652293>
- [5] N. AYAT, R. AKBARINIA, H. AFSARMANESH, P. VALDURIEZ. *Entity Resolution for Distributed Probabilistic Data*, in "Distributed and Parallel Databases", 2013, vol. 31, n^o 4, pp. 509-542, <http://hal.inria.fr/lirmm-00879631>
- [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, pp. 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, pp. 2351-2365 [DOI : 10.1109/TPDS.2012.24], <http://hal.inria.fr/lirmm-00748992>
- [8] A. JOLY, H. GOEAO, P. BONNET, V. BAKIC, J. BARBE, S. SELMI, I. YAHIAOUI, J. CARRÉ, E. MOUYSSET, J.-F. MOLINO, N. BOUJEMAA, D. BARTHÉLÉMY. *Interactive plant identification based on social image data*, in "Ecological Informatics", 2013 [DOI : 10.1016/J.ECOINF.2013.07.006], <http://www.sciencedirect.com/science/article/pii/S157495411300071X>
- [9] P. LETESSIER, O. BUISSON, A. JOLY. *Scalable Mining of Small Visual Objects*, in "Proceedings of the 20th ACM International Conference on Multimedia", New York, NY, USA, MM '12, ACM, 2012, pp. 599-608, <http://doi.acm.org/10.1145/2393347.2393431>
- [10] 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, pp. 1328-1339, <http://hal.inria.fr/hal-00640431/en>
- [11] E. PACITTI, R. AKBARINIA, M. EL DICK. , *P2P Techniques for Decentralized Applications*, Morgan & Claypool Publishers, 2012, 104 p. , <http://hal.inria.fr/lirmm-00748635>
- [12] 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, pp. 273-281, <http://hal.inria.fr/inria-00473481/en>
- [13] 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, pp. 649-674, <http://hal.archives-ouvertes.fr/hal-00374999/fr/>

- [14] C. ZHANG, F. MASSEGLIA, Y. LECHEVALLIER. *ABS: The Anti Bouncing Model for Usage Data Streams*, in "IEEE Int. Conf. on Data Mining (ICDM)", 2010, pp. 1169-1174
- [15] T. M. ÖZSU, P. VALDURIEZ. , *Principles of Distributed Database Systems, third edition*, Springer, 2011, 845 p. , <http://hal.inria.fr/hal-00640392/en>

Publications of the year

Doctoral Dissertations and Habilitation Theses

- [16] J. DIAS. , *Exécution interactive pour expériences computationnelles à grande échelle*, Universidade Federal de Rio de Janeiro, December 2013, <http://hal.inria.fr/tel-00939266>
- [17] P. LETESSIER. , *Découverte et exploitation d'objets visuels fréquents dans des collections multimédias*, Telecom ParisTech, March 2013, <http://hal.inria.fr/tel-00912992>
- [18] M. LIROZ-GISTAU. , *Partitionnement dans les Systèmes de Gestion de Données Parallèles*, Université Montpellier II - Sciences et Techniques du Languedoc, December 2013, <http://hal.inria.fr/tel-00920615>

Articles in International Peer-Reviewed Journals

- [19] R. AKBARINIA, P. VALDURIEZ, G. VERGER. *Efficient Evaluation of SUM Queries Over Probabilistic Data*, in "IEEE Transactions on Knowledge and Data Engineering", 2013, vol. 25, n^o 4, pp. 764-775, <http://hal.inria.fr/lirmm-00652293>
- [20] N. AYAT, R. AKBARINIA, H. AFSARMANESH, P. VALDURIEZ. *Entity Resolution for Distributed Probabilistic Data*, in "Distributed and Parallel Databases", 2013, vol. 31, n^o 4, pp. 509-542, <http://hal.inria.fr/lirmm-00879631>
- [21] W. K. DEDZOE, P. LAMARRE, R. AKBARINIA, P. VALDURIEZ. *As-Soon-As-Possible Top-k Query Processing in P2P Systems*, in "Transactions on Large-Scale Data- and Knowledge-Centered Systems (TLDKS)", 2013, 28 p. , <http://hal.inria.fr/lirmm-00821929>
- [22] A. HAMZAOUI, P. LETESSIER, A. JOLY, O. BUISSON, N. BOUJEMAA. *Object-based visual query suggestion*, in "Journal of Multimedia Tools and Applications", January 2013, pp. 1-26 [DOI : 10.1007/s11042-012-1340-5], <http://hal.inria.fr/hal-00806635>
- [23] A. JOLY, H. GOËAU, B. PIERRE, B. VERA, B. JULIEN, S. SOUHEIL, Y. ITHÉRI, C. JENNIFER, M. ELISE, J.-F. MOLINO, B. NOZHA, D. BARTHÉLÉMY. *Interactive plant identification based on social image data*, in "Ecological Informatics", August 2013, <http://hal.inria.fr/hal-00908872>
- [24] M. JORGE AUGUSTO, A. EDUARDO CUNHA DE, S. GERSON, T. YVES LE, P. VALDURIEZ. *Stress Testing of Transactional Database Systems*, in "Journal of Information and Data Management", 2013, vol. 4, n^o 3, pp. 279-294, <http://hal.inria.fr/lirmm-00905200>
- [25] M. LIROZ-GISTAU, R. AKBARINIA, E. PACITTI, F. PORTO, P. VALDURIEZ. *Dynamic Workload-Based Partitioning Algorithms for Continuously Growing Databases*, in "Transactions on Large-Scale Data- and Knowledge-Centered Systems", 2014, 105 p. , <http://hal.inria.fr/lirmm-00906966>

- [26] E. OGASAWARA, D. JONAS, V. SILVA, C. FERNANDO, O. DANIEL DE, F. PORTO, M. MATTOSO, P. VALDURIEZ. *Chiron: A Parallel Engine for Algebraic Scientific Workflows*, in "Journal of Concurrency and Computation: Practice and Experience", 2013, vol. 25, n^o 16, pp. 2327-2341, <http://hal.inria.fr/lirmm-00806557>
- [27] K. PARK, P. VALDURIEZ. *A Hierarchical Grid Index (HGI), spatial queries in wireless data broadcasting*, in "Distributed and Parallel Databases", February 2013, vol. 31, n^o 3, pp. 413-446 [DOI : 10.1007/s10619-013-7121-Y], <http://hal.inria.fr/lirmm-00797095>

International Conferences with Proceedings

- [28] R. AKBARINIA, F. MASSEGLIA. *Fast and Exact Mining of Probabilistic Data Streams*, in "PKDD'2013: European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases", Prague, Czech Republic, Lecture Notes in Computer Science, Springer, 2013, pp. 493-508 [DOI : 10.1007/978-3-642-40988-2_32], <http://hal.inria.fr/lirmm-00838618>
- [29] V. BAKIC, S. MOUINE, S. OUERTANI-LITAYEM, A. VERROUST-BLONDET, I. YAHIAOUI, H. GOËAU, A. JOLY. *Inria's participation at ImageCLEF 2013 Plant Identification Task*, in "CLEF (Online Working Notes/Labs/Workshop) 2013", Valencia, Spain, September 2013, <http://hal.inria.fr/hal-00874330>
- [30] C. BARBARA, M. HENNING, T. BART, V. MAURICIO, R. PAREDES, D. ZELLHOFER, H. GOËAU, A. JOLY, B. PIERRE, J. MARTINEZ GOMEZ, I. GARCIA VAREA, M. CAZORLA. *Imageclef 2013: the vision, the data and the open challenges*, in "CLEF 2013", Valencia, Spain, September 2013, <http://hal.inria.fr/hal-00908881>
- [31] E. CASTANIER, R. COLETTA, P. VALDURIEZ, C. FRISCH. *WebSmatch: a tool for Open Data*, in "WOD'2013: International Workshop on Open Data", Paris, France, 2013, <http://hal.inria.fr/lirmm-00858602>
- [32] H. GOËAU, P. BONNET, A. JOLY, V. BAKIC, J. BARBE, S. SELMI, J. CARRÉ, D. BARTHÉLÉMY, N. BOUJEMAA, J.-F. MOLINO, G. DUCHÉ, A. PERRONET. *PI@ntNet Mobile App*, in "ACM Multimedia", Barcelone, Spain, ACM, October 2013, pp. 423-424 [DOI : 10.1145/2502081.2502251], <http://hal.inria.fr/hal-00908910>
- [33] N. HERVÉ, M.-L. VIAUD, J. THIEVRE, A. SAULNIER, P. LETESSIER, J. CHAMP, O. BUISSON, A. JOLY. *OTmedia: The French Transmedia News Observatory*, in "ACM Multimedia", Barcelone, Spain, October 2013 [DOI : 10.1145/2502081.2502260], <http://hal.inria.fr/hal-00908925>
- [34] A. JOLY, H. GOËAU, P. BONNET, V. BAKIC, J.-F. MOLINO, D. BARTHÉLÉMY, N. BOUJEMAA. *The Imageclef Plant Identification Task 2013*, in "International workshop on Multimedia analysis for ecological data", Barcelone, Spain, October 2013, <http://hal.inria.fr/hal-00908934>
- [35] D. JONAS, E. OGASAWARA, O. DANIEL DE, F. PORTO, P. VALDURIEZ, M. MATTOSO. *Algebraic Dataflows for Big Data Analysis*, in "BigData'2013: International Conference on Big Data", Santa Clara, United States, IEEE, October 2013, 6 p. , <http://hal.inria.fr/lirmm-00857221>
- [36] P. LETESSIER, N. HERVÉ, C. JULIEN, A. JOLY, O. BUISSON, A. HAMZAOUI. *Small objects query suggestion in a large web-image collection*, in "ACM Multimedia", Barcelone, Spain, ACM, October 2013 [DOI : 10.1145/2502081.2502248], <http://hal.inria.fr/hal-00908891>

- [37] M. LIROZ-GISTAU, R. AKBARINIA, D. AGRAWAL, E. PACITTI, P. VALDURIEZ. *Data Partitioning for Minimizing Transferred Data in MapReduce*, in "Globe'2013: 6th International Conference on Data Management in Cloud, Grid and P2P Systems", Prague, Czech Republic, A. HAMEURLAIN, W. RAHAYU, D. TANIAR (editors), LNCS, Springer, August 2013, pp. 1-12 [DOI : 10.1007/978-3-642-40053-7_1], <http://hal.inria.fr/lirmm-00879527>
- [38] T. RUIMING, S. DONGXU, S. BRESSAN, P. VALDURIEZ. *What you Pay for is What you Get*, in "DEXA'2013: 24th International Conference on Database and Expert Systems Applications", Prague, Czech Republic, H. DECKER, L. LHOTSKA, S. LINK (editors), Springer, 2013, pp. 395-409, <http://hal.inria.fr/lirmm-00831864>
- [39] T. RUIMING, W. HUAYU, B. ZHIFENG, B. STEPHANE, P. VALDURIEZ. *The Price is Right: Models and Algorithms for Pricing Data*, in "DEXA'2013: 24th International Conference on Database and Expert Systems Applications", Czech Republic, H. DECKER, L. LHOTSKA, S. LINK (editors), Springer, 2013, pp. 380-394, <http://hal.inria.fr/lirmm-00831859>
- [40] M. SERVAJEAN, E. PACITTI, S. AMER-YAHIA, P. NEVEU. *Profile Diversity in Search and Recommendation*, in "SRS 2013: 4th International Workshop on Social Recommender Systems (in conjunction WWW 2013)", Rio de Janeiro, Brazil, I. GUY, L. CHEN, M. X. ZHOU (editors), IW3C2, May 2013, pp. 973-980, Paper Session: User Models - WWW'13 Companion (dl.acm.org/citation.cfm?id=2488094), <http://hal.inria.fr/lirmm-00806676>
- [41] E. SZEKELY, P. PONCELET, F. MASSEGLIA, M. TEISSEIRE, R. CEZAR. *A Density-Based Backward Approach to Isolate Rare Events in Large-Scale Applications*, in "DS'13: Discovery Science", Singapore, Singapore, J. FÜRNKRANZ, E. HÜLLERMEIER, T. HIGUCHI (editors), Lecture Notes in Computer Science, Springer, October 2013, pp. 249-264 [DOI : 10.1007/978-3-642-40897-7_17], <http://hal.inria.fr/lirmm-00907893>
- [42] C. ZHANG, Y. HAO, M. MAZURAN, C. ZANIOLO, H. MOUSAVI, F. MASSEGLIA. *Mining frequent itemsets over tuple-evolving data streams*, in "SAC'13: Symposium on Applied Computing", Coimbra, Portugal, March 2013, pp. 267-274, <http://hal.inria.fr/lirmm-00830923>

Conferences without Proceedings

- [43] M. LIROZ-GISTAU, R. AKBARINIA, D. AGRAWAL, E. PACITTI, P. VALDURIEZ. *MR-Part : Minimizing Data Transfers Between Mappers and Reducers in MapReduce*, in "BDA'2013 : 29èmes journées Bases de Données Avancées", Nantes, France, October 2013, <http://hal.inria.fr/lirmm-00879531>
- [44] M. SERVAJEAN, E. PACITTI, S. AMER-YAHIA, P. NEVEU. *Profile Diversity for Phenotyping Data Search and Recommendation*, in "BDA 2013 - 29e Journées Bases de Données Avancées", Nantes, France, October 2013, 20 p. , Session: Applications innovantes, <http://hal.inria.fr/lirmm-00879575>

Scientific Books (or Scientific Book chapters)

- [45] M. JAWAD, P. SERRANO-ALVARADO, P. VALDURIEZ. *Supporting Data Privacy in P2P Systems*, in "Security and Privacy Preserving in Social Networks", R. CHBEIR, B. A. BOUNA (editors), Springer, August 2013, 50 p. , <http://hal.inria.fr/hal-00807625>

Research Reports

- [46] P. LETESSIER, O. BUISSON, A. JOLY. , *Scalable Mining of Small Visual Objects (with new experiments)*, December 2013, <http://hal.inria.fr/hal-00912560>
- [47] A. ROUSSEAU, A. DARNAUD, B. GOGLIN, C. ACHARIAN, C. LEININGER, C. GODIN, C. HOLIK, C. KIRCHNER, D. RIVES, E. DARQUIE, E. KERRIEN, F. NEYRET, F. MASSEGLIA, F. DUFOUR, G. BERRY, G. DOWEK, H. ROBAK, H. XYPAS, I. ILLINA, I. GNAEDIG, J. JONGWANE, J. EHREL, L. VIENNOT, L. GUION, L. CALDERAN, L. KOVACIC, M. COLLIN, M.-A. ENARD, M.-H. COMTE, M. QUINSON, M. OLIVI, M. GIRAUD, M. DORÉMUS, M. OGOUCHI, M. DROIN, N. LACAUX, N. ROUGIER, N. ROUSSEL, P. GUITTON, P. PETERLONGO, R.-M. CORNUS, S. VANDERMEERSCH, S. MAHEO, S. LEFEBVRE, S. BOLDO, T. VIÉVILLE, V. POIREL, A. CHABREUIL, A. FISCHER, C. FARGE, C. VADEL, I. ASTIC, J.-P. DUMONT, L. FÉJOZ, P. RAMBERT, P. PARADINAS, S. DE QUATREBARBES, S. LAURENT. , *Médiation Scientifique : une facette de nos métiers de la recherche*, March 2013, 34 p. , <http://hal.inria.fr/hal-00804915>
- [48] M. SERVAJEAN, E. PACITTI, S. AMER-YAHIA, A. EL ABBADI, P. NEVEU. , *Profile Diversity for P2P Search and Recommendation*, June 2013, <http://hal.inria.fr/lirmm-00829308>
- [49] M. SERVAJEAN, E. PACITTI, S. AMER-YAHIA, P. NEVEU. , *Profile Diversity in Search and Recommendation*, February 2013, <http://hal.inria.fr/lirmm-00794814>
- [50] E. SZEKELY, P. PONCELET, F. MASSEGLIA, M. TEISSEIRE, R. CEZAR. , *Isolating rare events in large-scale applications using a backward approach*, January 2013, <http://hal.inria.fr/lirmm-00798074>

Patents and standards

- [51] S. BRINGAY, R. CEZAR, D. IENCO, A. MAS, F. MASSEGLIA, P. PONCELET, P. PUDLO, M. TEISSEIRE, J.-P. VENDRELL, E. SZEKELY. , *Process for identifying rare events*, 2013, <http://hal.inria.fr/lirmm-00913008>

Other Publications

- [52] T. ALLARD, B. NGUYEN, P. PUCHERAL. , *Comment préserver l'anonymat ?*, November 2013, Pour la Science, nr 433, <http://hal.inria.fr/hal-00937554>