



Project-Team LIS

# ***Logical Information Systems***

*Rennes*

*Activity Report*

*2014*



## 1 Team

### Head of the team

Olivier Ridoux, Professor, Université Rennes 1

### Administrative assistant

Aurélie Patier

### Université Rennes 1 personnel

Yves Bekkers, Professor, until July

Sébastien Ferré, Associate Professor, HDR

Annie Foret, Associate Professor, HDR

Benjamin Sigonneau, Research Engineer (part time: 70%), until June

### Insa personnel

Mireille Ducassé, Professor

Peggy Cellier, Associate Professor

### PhD students

Mouhamadou Ba, ARED/Insa grant, since October 2012

Soda Marème Cissé, MENRT/Rennes 1 grant, since October 2012

### Associate members

Erwan Quesseveur, Associate Professor, Université Rennes 2, Associate Member

François Le Prince, President of ALKANTE and Associated Professor, Université Rennes 2, Associate Member

## 2 Overall Objectives

### 2.1 Overview

The LIS team aims at developing *formal* methods for handling complex data sets in a *flexible* and *precise* way. “Flexible” means that the content determines the shape of the container. Very often, it is the opposite that is observed; e.g., the tree-like shape of a hierarchical file system enforces the tree-like shape of software packages. “Precise” means that any subset of the data set can be easily characterized. Again, it is the opposite that is often observed; e.g., in a hierarchical file system only sub-trees can be easily characterized. More and more information is available on the Web, and more and more information can be stored on a single machine. However, whereas the related low-level technology is developing, and performance is increasing, little is done for organizing the ever-growing amount of information. Therefore, the LIS team addresses the issues of organizing and querying information in general. The solutions are to be

both formal and practical. Operational issues such as index technologies are important, but we are convinced that their scope is too limited to solve the crucial issues.

At a formal level, *queries* and *answers* are two key notions. It is nowadays standard to consider queries as logical formulas and answers as special models of queries. Computing the preferred model of a query in some context is conceptually easy, and it warrants flexibility. However, the opposite is not that easy in general; given a subset of the data, how can we compute a query of which it is a model? Given two different subsets of the data, how can we compute a query that explains the difference? Knowing this would warrant precision. The LIS team proved that *formal concept analysis* (FCA [GW99]) is a powerful framework for analyzing  $\langle \text{query}, \text{answer} \rangle$  pairs. *Formal concepts* formalize the association between a query and its answers. Formal concepts are structured into a lattice which provides navigation links between concepts.

However, standard FCA cannot deal with queries considered as logical formulas (recall that this is the key for flexibility). Therefore a variant of FCA for logical description has been developed [7] altogether with the generic notion of *Logical information system* (LIS) that provided a reconstruction of all information system operations based on logical concept analysis. In particular, some data-mining operations are native in LIS [7, 5].

The mottoes of the LIS research are:

1. *Never impose a priori a structure on information.* E.g., do not use hierarchical structures. Imposing *a priori* a structure causes the *tyranny of the dominant decomposition*[TOHS99]. For instance, the usual class-based organisation of source code makes highly visible the connections between methods of the same class, but masks the possible connections between methods in different classes.

Instead, consider pieces of information as a bulk. Structure should emerge *a posteriori* from the contents or the point of view. As a consequence, updating the contents may change the structure: we accept it.

2. *Consider every possible rational classification, and permit changes at any time.* Here, rational means that what makes a piece of information belong or not to a class depends on the very piece of information, not on other pieces. The concept lattice induced by FCA is precisely a means to grasp all possible rational classifications.
3. *Rare events are as important as the frequent ones.* One cannot say *a priori* if a piece of information is interesting because it presents a frequent pattern, or because it presents a rare pattern.

So, rare events must not be masked by statistical artefacts. Statistics is not forbidden, but it is only a complement of a symbolic logic approach.

4. *Queries should be possible answers.*

---

[GW99] B. GANTER, R. WILLE, *Formal Concept Analysis — Mathematical Foundations*, Springer, 1999.

[TOHS99] P. TARR, H. OSSHER, W. HARRISON, S. SUTTON, “N Degrees of Separation: Multi-Dimensional Separation of Concerns”, *in: ICSE*, IEEE Computer Society, p. 107–119, 1999.

In usual information systems (say relational databases or Web browsers) there is a strict dichotomy between queries (they are intensional expressions), and answers (they are strictly extensional expressions, i.e. sets of things). We contend that a good answer must be a mix of extensional and intensional answers. E.g. the good answer to “I would like to buy a book” is seldom the whole catalog of the bookshop; it is more relevant to answer such a query with other queries, like “Is this for a child” or “Do you prefer novels or documents”.

Note that hierarchical file systems already do that. Queries (i.e., *filepaths*) yield answers that contain other queries (i.e., *sub-directories*). One of the LIS achievements is a formalization of this behaviour that does not rely on an *a priori* hierarchical structure.

Our research is intended to be *vertical* in the sense that all aspects of information systems are of interest: design, implementation, and applications.

On the implementation side, the LIS team develops systems that present the LIS abstraction at the user level [9].

On the application side, the LIS team explores the application of LIS to *Geographical information systems* (GIS). The intuition here is that the traditional layered organization of information in GIS suffers a rigid structure of thematic layers. Moreover, GIS applications usually cope with highly heterogeneous information and large amount of data; this makes them an interesting challenge for LIS. The team also works on a data-mining interpretation of bug tracking. In this case, the intuition is that pieces of information relevant to software engineering, e.g. programs, specifications or tests, can be explored very systematically by a LIS. More generally, applications to software engineering are important for the team. A recent trend of application is the assistance to group decision and negotiation. In particular LIS provide new technological support to *social choice*. For example, in committee decision making, navigating with LIS in the facts recorded in a context allows decision makers to treat all candidates in a fair way [3].

## 2.2 Key Issues

In its current state, LIS studies the following key issues:

- The LIS formalism is generic w.r.t. the logic used for describing pieces of information.  
*What are the appropriate logics for the application fields that we have chosen? (GIS and error localization) Do we need a brand new logic for every application, or is there something that different applications can share?*
- Genericity of LIS w.r.t. logic opens the door for creating ad hoc logics for describing pieces of information of an application. We already have proposed the framework of *logic functors* for helping a user build safely ad hoc logics. Logic functors are certified logic components that can be composed to form certified implementations of a logic.  
*What are the useful logic functors? How can we be sure that a toolbox of logic functors is complete for a given purpose?*

*Can the idea of certified composition be applied to another domain?* Given a domain *foo*, *foo* functors would be certified *foo* components that can be assembled to form certified implementations of *foo* systems.

*Is it possible to certify other properties than meta-logical properties?* E.g. is it possible to characterize complexity, or other non-functional properties like security?

- A family of non-commutative logics has developed over the years in the domain of computational linguistics, e.g. Lambek logic, pregroups. As for LIS, a great amount of creativity is expected for extending this family with ad hoc logics that would tackle fine-grained linguistics phenomena.

*Is it possible to build up an implementation of these logics using logic functors?*

Some LIS applications deal with objects that are sequential by nature (say, texts).

*Can these non-commutative logics primarily developed for computational linguistics help in LIS applications?*

- Hierarchical file systems have a preferred metaphor which is the tree.

*What is the proper metaphor for LIS?*

The tree is also the graphical metaphor of hierarchical file systems.

*What is the graphical metaphor for LIS?*

Knowing this is crucial for the acceptance of LIS in end-user applications.

- Geographical information systems also suffer the *tyranny of the dominant decomposition*. Here, the dominant decomposition is in rigid thematic layers that inherit from plastic sheets of ancient map design. These layers are omnipresent in the design and interface of GIS applications.

*How can LIS abstract these layers, and still display layers when needed?*

Mining geographical information is difficult because of the layers and because it must cope with complex spatial relations.

*What is the proper modeling of these relations that will permit efficient LIS operations, including data-mining?*

- Up to now, mining execution traces for bug tracking has used poor trace representations and ad hoc algorithms.

*How can the theoretical and practical framework of LIS help benefit from the wide range of information of program development environments?*

- The file system implementation of LIS can handle around 1 million elementary pieces of information, which corresponds approximately to a full homedir with 10 to 20 thousands files. This is rather small compared to relational database capabilities, but already large compared to other approaches based on formal concept analysis.

*How can it handle more? Can we reach 100 million in the next few years?*

## 3 Scientific Foundations

### 3.1 Logics for Information Systems

**Keywords:** Syntax, interpretation, semantics, subsumption.

**Glossary :**

**Syntax** Definition of the well-formed statements of a language. Statements are finite.

**Interpretation** Complete description of a world. Interpretations can be arbitrary mathematical constructs, and so can be infinite. Interpretations are models of statements, namely the worlds in which the statement is true. Statements are features of interpretations, namely the statements that are true in the world.

**Semantics** A binary relation between syntactic statements and interpretations.

**Subsumption** A relation which states that a property is more specific than another property.

Logic is the core of Logical Information Systems. However, this does not say everything because every particular usage of logic is also a point of view on logic. For instance, logic in Logic Programming is not the same as in Description Logics. This section describes the point of view on logic from information systems.

Logic is a wide domain that is concerned with formal representation and reasoning. The point of view on logic in logical information systems can be characterized by two things. Firstly, we are interested in the individual description of objects (e.g., files, pictures, program functions or methods), so that we need to represent concrete domains and data structures. This entails two levels of statements: (1) statements about objects, and (2) statements about the world (e.g., ontologies and *subsumption*). Subsumption helps to decide when an object is an answer to a query. Secondly, we need automated reasoning facilities as the subsumption must be decided between any object and a query in information retrieval. This forces us to only consider decidable logics, unless consistency or completeness are weakened.

**Properties of a Logic** A characteristic of logic is the ability to derive new statements from known statements. Such a derivation is valid w.r.t. semantics only if every model of the known statements is also a model of the new statements. This ability opens the room for *reasoning*, i.e. the production of valid statements by working at the syntactical level only. Reasoning is formalized by *inference systems* (e.g., axioms and rules). An inference system is *consistent* if it produces only valid statements; it is *complete* if it produces all valid statements. Reasoning is *decidable* if a consistent and complete inference system can be realized by an algorithm.

**Examples of Logics for Information Systems** Proposition logic is a possible logic for an information system, but it needs a lot of encoding for handling structured information. Instead, non-standard logics have been defined for some structured domains.

A large family of logics that comes into our scope is the family of Description Logics

(DL) [Bra79,CLN98], which have been widely studied, implemented, and applied in knowledge and information management. Moreover, their semantic structure is especially well-suited to be used in a LIS. The semantics of proposition logic is often exposed in terms of truth values and truth tables. To the contrary, the semantics of description logic is defined in terms of sets of objects that are close to answers to a query. DL are, therefore, of a special interest for the LIS team.

Another family of interest is *categorial grammars*. Many substructural logics come into this scope, among which non-commutative linear logic or Lambek Calculus<sup>[Lam58]</sup> that handle various concatenation principles (or ordered conjunction) in categorial grammars where logic is used both for attaching formulas to objects and for parsing seen as deduction.

At an empirical level, the categorial approach comes very close to the LIS approach. Categorial grammars correspond to LIS contents, because they both attach formulas to objects, and sentence types correspond to queries. The difference is that the answer to a LIS query is an unordered set, whereas a sentence generated by a categorial grammar is an ordered sequence. We expect a cross-fertilization of both theories in the future, especially in the LIS applications where the objects are naturally ordered.

### 3.2 Concept Analysis

**Keywords:** Objects, descriptors, context, instance, property, extension, intension, concept.

**Glossary :**

**Objects** A set of distinguished individuals.

**Descriptors** A set of distinguished properties.

**Context** A set of objects associated with descriptors.

**Instance** An object is an *instance* of a descriptor if it is associated with it in a given context.

**Property** A descriptor is a *property* of an object if it is associated with it in a given context.

**Extension** The *extension* of a collection of descriptors is the set of their common instances. Extent is a synonym.

**Intension** The *intension* of a collection of objects is the set of their common properties. Intent is a synonym.

**Concept** Given a context, and extensions and intensions taken from it, a *concept* is a pair  $(E, I)$  of an extension  $E$  and an intension  $I$  that are mutually complete; i.e.,  $I$  is the intension of the extension, and  $E$  is the extension of the intension.

---

[Bra79] R. J. BRACHMAN, “On the Epistemological Status of Semantic Nets”, *in: Associative Networks: Representation of Knowledge and Use of Knowledge by Examples*, N. V. Findler (editor), Academic Press, New York, 1979.

[CLN98] D. CALVANESE, M. LENZERINI, D. NARDI, “Description Logics for Conceptual Data Modeling”, *in: Logics for Databases and Information Systems*, J. Chomicki, G. Saake (editors), Kluwer, p. 229–263, 1998.

[Lam58] J. LAMBEK, “The Mathematics of Sentence Structure”, *American Mathematical Monthly* 65, 1958, p. 154–170.

**Formal Concept Analysis** Formal Concept Analysis (FCA) is part of the mathematical branch of applied lattice theory [Bir40,DP90]. It can be seen as a reformulation by Wille of Galois lattices [BM70] that emphasizes lattices as conceptual hierarchies [Wil82]. The mathematical foundations of FCA have been extensively studied by Ganter and Wille [GW99].

FCA mainly aims at the automatic construction of *concepts* and their classification according to a generalization ordering, given a flat representation of data. The adjective *formal* means that concepts are given a mathematical definition, which reflects the usual philosophical meaning of a “concept”. The basic notions of FCA are those of *formal context*, and *formal concept*.

A *formal context* is a binary relation between a set of objects, and a set of attributes. Through this relation attributes can be seen as properties of objects, and reciprocally, objects can be seen as instances of attributes. This is a very general settings that applies to various domains such as data analysis, information retrieval, data-mining or machine learning. In all these domains, the objects of interest are described by sets of attributes, and the objective is to relate in some way sets of objects and sets of attributes. In information retrieval a set of attributes is a query, whose answers is a set of objects. In machine learning a set of objects is a set of positive examples, whose characterization is a set of attributes.

A *formal concept* is the association of a set of objects, the *extent*, and a set of attributes, the *intent*. This comes close to the classical definition of concept in philosophy, but in FCA the relationship between extent and intent is formally defined. The extent must be the set of instances shared by all attributes of the intent; and the intent must be the set of properties shared by all objects in the extent.

The fundamental theorem of FCA says that the set of all concepts forms a complete lattice when they are ordered according to the set inclusion on extents (or intents). This is called the *concept lattice*, and it can be computed automatically from the formal context. The concept lattice is the structure that is implicit in any formal context. It contains all the information contained in the formal context; the latter can be rebuilt from the former. In data analysis, the concept lattice permits a flexible classification of data (where a concept is a class), because concepts are not organized as a strict hierarchy. In information retrieval and data-mining it is used as a search space for answers.

**Logical Concept Analysis** In Formal Concept Analysis (FCA) object properties are restricted to Boolean attributes. In many applications there is a need for richer properties, where properties are not independent. For instance, if a book has been published in 2000, it can be given the property `year = 2000`, and has then the implicit properties `year in 1990..2000` and `year in 2000..2010`. This means that properties are statements about objects that can

- 
- [Bir40] G. BIRKHOFF, *Lattice Theory*, American Mathematical Society, 1940.
  - [DP90] B. A. DAVEY, H. A. PRIESTLEY, *Introduction to Lattices and Order*, Cambridge University Press, 1990.
  - [BM70] M. BARBUT, B. MONJARDET, *Ordre et classification — Algèbre et combinatoire (2 tomes)*, Hachette, Paris, 1970.
  - [Wil82] R. WILLE, *Ordered Sets*, Reidel, 1982, ch. Restructuring lattice theory: an approach based on hierarchies of concepts, p. 445–470.
  - [GW99] B. GANTER, R. WILLE, *Formal Concept Analysis — Mathematical Foundations*, Springer, 1999.

be subject to reasoning, exactly like logical statements. Other examples of useful properties are strings and string patterns, spatial descriptions for locating objects, or patterns over the programming type of functions and methods.

FCA has been extended by other authors to handle multi-valued contexts <sup>[GW99]</sup>, but this extension takes the form of a preprocessing stage that results in a standard formal context, and forgets all logical relations between properties. Moreover it is limited in practice to valued attributes with finite domains of attributes. In 2000 we proposed a logical generalization of FCA, named Logical Concept Analysis (LCA) [7], that is the abstraction of FCA w.r.t. object descriptions and concept intents. This makes LCA an abstract component, and makes FCA the composition of LCA with a logic component. LCA makes the theory of concept analysis easily reusable in various applications.

For good composability of LCA and logics, they must agree on the specification of logics. What LCA needs from a logic is:

- a language of formulas (or statements),  $L$ , for the representation of object descriptions and concept intents,
- a procedure,  $\sqsubseteq$ , for deciding the subsumption between 2 formulas;  $\sqsubseteq$  means “is subsumed by”, “is more specific than”, “entails”,
- a procedure,  $\sqcup$ , for computing the least common subsumer of 2 formulas; it is a kind of logical disjunction,
- a formula,  $\perp$ , that is the most specific according to subsumption (logical contradiction).

This specification provides everything required to extend fundamental results of FCA to LCA (formal context, extent, intent, concept, complete lattice of concepts). For information retrieval and the expression of queries, it is useful to add, to this specification, operations such as logical conjunction, and logical tautology (the most general formula).

Any formal context defines a logic whose subsumption relation is isomorphic to the concept lattice that is derived from the formal context. An interesting result is that the *contextualized logic* (the logic defined by the logical context) is a refinement or extension of the logic used by LCA. Everything true in the logic is also true in the contextualized logic (because it is *eternal truth*); and everything true only in the contextualized logic says something that is true in the context, but not in general (because it is *instant truth*). Thus, contextualized logic forms the basis for data-mining and machine learning tasks, whose aim is to discover outstanding regularities in a given context [7, 5].

### 3.3 Logical Querying, Navigation, and Data-mining

**Keywords:** Querying, navigation, data-mining.

#### Glossary :

**Querying** The process that takes a query (e.g., a logical formula), and returns the collection of objects that satisfy the query (e.g., the extent of the query).

---

[GW99] B. GANTER, R. WILLE, *Formal Concept Analysis — Mathematical Foundations*, Springer, 1999.

**Navigation** The process of moving from place to place, where each place indicates objects they contain (i.e. *local objects*) and other places where it is possible to move (i.e. *neighbouring places*).

**Data-mining** The process of extracting outstanding regularities from data (e.g., a context) hoping to discover new and useful knowledge.

In most information systems, querying and navigation are two disconnected means for information retrieval. With querying, users formulate queries which belong to more or less complex querying languages, from simple words as in Google to highly structured languages like SQL. The system returns a set of answers to the query. This permits expressive search criteria over large amounts of data, but lacks interactivity because the dialogue is only one-way. If the answers are not satisfying, users have to imagine new queries and formulate them, which requires *a priori* knowledge of both querying language and data. With navigation, users move from place to place following links. The most common systems are folder hierarchies (e.g., file systems, bookmarks, emails), and hypertext. As opposed to querying, navigation provides interactivity by making suggestions at each step, but offers limited expressivity because navigation structures are rigid. In a hierarchy, selection criteria are presented in a fixed order. For instance, if pictures are classified first by date, then by type, one cannot easily find all landscape pictures.

The need for combining querying and navigation has already been recognized. Most proposals, however, are unsatisfying. Indeed, either querying and navigation cannot be mixed freely in a same search, or consistency of querying is not maintained. An example of the former is SFS [GJSO91], once a querying step is done, there is no more navigation. An example of the latter is HAC [GM99], some query answers may not satisfy the query. A proposal based on FCA has not these drawbacks [GMA93], and we have generalized it to work within LCA, which allows us to use logical formulas for object description and queries [7]. Logic brings expressivity in querying, and concept analysis brings the concept lattice as a navigation structure (i.e., navigation places are formal concepts). The advantages of this navigation structure is that (1) it is automatically derived from data, the logical context (see motto 1), (2) it is complete as navigation alone makes it possible to reach any object (see motto 3), and (3) it is flexible because selection criteria can be chosen in any order, thus allowing user to express their preferences (see motto 2). Querying and navigation can be freely mixed (see motto 4) in a same search because every logical formula points to a formal concept, and every formal concept is labelled by a logical formula. Put concretely, this means that a user can at each step of his search: either modify by hand the current query and reach a new place, or follow a suggested link that will modify the current query and reach a new place.

The critical operation is the computation of navigation links, which correspond to edges in

- 
- [GJSO91] D. K. GIFFORD, P. JOUVELOT, M. A. SHELDON, J. W. J. O'TOOLE, "Semantic file systems", *in: ACM Symp. Operating Systems Principles*, ACM SIGOPS, p. 16–25, 1991.
  - [GM99] B. GOPAL, U. MANBER, "Integrating Content-Based Access Mechanisms with Hierarchical File Systems", *in: third symposium on Operating Systems Design and Implementation*, USENIX Association, p. 265–278, 1999.
  - [GMA93] R. GODIN, R. MISSAOUI, A. APRIL, "Experimental Comparison of Navigation in a Galois Lattice with Conventional Information Retrieval Methods", *International Journal of Man-Machine Studies* 38, 5, 1993, p. 747–767.

the concept lattice. Indeed, the worst-case time complexity for computing the concept lattice is exponential in the number of objects, which makes it intractable in most interesting cases. We demonstrated both in theory and practice that this computation is not necessary. A key feature of LIS is that its semantics is expressed in terms of LCA, though it is not required to actually build the concept lattice. This is opposed to most (all?) previous proposals for using LCA in information retrieval.

The concept lattice upon which our navigation is based is also a rich structure for data-mining and machine learning <sup>[Kuz04]</sup>. Here again, we have combined existing techniques with logic [7, 5, 2], and applied them to the automatic classification of emails, and the prediction of the function of proteins from their sequence [5].

### 3.4 Genericity and Components

**Keywords:** Abstraction, reusability, composability, component.

#### Glossary :

**Abstraction** a mechanism and practice to reduce and factor out details so that one can focus on few concepts at a time.

**Reusability** the likelihood a segment of structured code can be used again to add new functionalities with slight or no modification. Reusable code reduces implementation time, it increases the likelihood that prior testing and use has eliminated bugs and it localizes code modifications when a change in implementation is required.

**Composability** a system design principle that deals with the inter-relationships of components. A highly composable system provides recombinant components that can be selected and assembled in various combinations to satisfy specific user requirements.

**Component** a unit of composition with contractually specified interfaces and explicit context dependencies only. A software component can be deployed independently and is subject to composition by third parties.

The application scope of Logical Information Systems is very large, and we do not expect that one design (e.g., one logic) will fit all possible applications. That is why we emphasize genericity, and we use the plural in “logical information systems”. The need for genericity is not limited to theoretical results and design, but extends to the concrete implementation of LIS.

Genericity requires programming facilities for *abstraction*, *composability*, and *reusability* of *software components*.

In LIS, abstraction is of the upper importance in the design of logical concept analysis; LCA is an abstraction of FCA. It is also at the heart of the *logic functor* framework and of its implementation; a logic functor is an abstraction of a logic (see Section 3.5). Reusability and composability are the expected outcomes of this framework. It is expected to make things easier to the designer of a LIS application. Composability is also at the heart of the very notion of formal context, and thus at the heart of concept analysis. Indeed, the flat structure

---

[Kuz04] S. O. KUZNETSOV, “Machine Learning and Formal Concept Analysis.”, *in: Int. Conf. Formal Concept Analysis*, P. W. Eklund (editor), *LNCS 2961*, Springer, p. 287–312, 2004.

of formal concepts makes it trivial to extend a context or merge two contexts, and the burden of giving a structure to the context is left to the construction of the concept lattice.

A generic implementation of LIS can be seen as a central component that is parameterized by several application-dependent components: at least a logic, and a transducer for importing data. These parameter components can be linked at compilation time (plugins). The central component as well as parameter components can themselves be the result of the composition of smaller components.

### 3.5 Logic Functors

**Keywords:** Logics, genericity, composability.

The genericity w.r.t. logic implies that for every new application a logic has to be found for describing objects in a logic context. Either a suitable logic is already known, or it must be created. Creating a logic requires designing a syntax, a semantics, algorithms for subsumption and other procedures, and proving that these algorithms are correct w.r.t. semantics. This definitely requires logic expertise and programming skills, especially for the subsumption procedure that is a theorem prover for which consistency and completeness must be proven. However, application developers and logic experts are likely to be different persons in most cases. Moreover, creating new logics from scratch for each application is unsatisfying w.r.t. reusability as these logics certainly share common parts. For instance, many applications need propositional reasoning, only changing the notion of what is a propositional variable.

We introduced high-level logic components, named *logic functors* [6], in order to make the creation of a new logic the mere composition of abstract and reusable components. All logics share a common specification that contains all useful procedures (e.g., subsumption); logic functors are functions from logics to logics, implemented as parameterized modules. Some logic functors take no parameter, and provide stand-alone but reusable logics: this is the case of concrete domains such as integers or strings. Other logic functors take one or several logics as parameters. For instance the functor  $\text{Prop}(X)$  is propositional logic abstracted over its atoms. This makes it possible to replace atoms in propositional logic by the formulas of another logic (e.g., valued attributes, terms from a taxonomy).

Logics are built by applying logic functors to sub-logics, which can themselves be defined as a composition of logic functors. For instance, the propositional logic where atoms are replaced by integer-valued attributes (and allowing for integer intervals) can be defined by the expression

$$L = \text{Prop}(\text{Set}(\text{Prod}(\text{Atom}, \text{Interval}(\text{Int})))) .$$

This results in a concrete software component  $L$  that is fully equipped with implementations of the logic specification procedures. This component can then be composed itself with LCA or a LIS system.

### 3.6 Categorical Grammars

**Keywords:** Categorical grammar, identification in the limit.

Categorial grammars are used for natural language modeling and processing; they mainly handle syntactic aspects, but Lambek variants also have a close link with semantics and Lambda-calculus. Formally, a *categorial grammar* is a structure  $G = (\Sigma, I, S)$  where:  $\Sigma$  is a finite alphabet (the words in the sentences);  $I$  is a function that maps a finite set of types to each element of  $\Sigma$  (the possible categories of each word, a lexicon);  $S$  is the *main type* associated to correct sentences. A *k-valued categorial grammar* is a categorial grammar where, for every word  $a \in \Sigma$ ,  $I(a)$  has at most  $k$  elements. A *rigid categorial grammar* is a 1-valued categorial grammar. Rigidity is a useful constraint to get learnable subclasses of grammars (and related algorithms).

Each variant of categorial grammar formalism is also determined by a derivability relation on types  $\vdash$  (which can be seen as a subcase of *linear logic* deduction in the case of Lambek grammars). Given a categorial grammar  $G = \langle \Sigma, I, S \rangle$ , a sentence  $w$  on the alphabet  $\Sigma$  belongs to the language of  $G$  whenever the words in  $w$  can be assigned by  $I$  a sequence of types that derive (according to  $\vdash$ ) the distinguished type  $S$ .

A simplified example is  $G_1 = (\Sigma_1, I_1, S)$  with  $\Sigma_1 = \{John, Mary, likes\}$   $I_1 = \{John \mapsto \{N\}, Mary \mapsto \{N\}, likes \mapsto \{N \setminus (S/N)\}\}$  the sentence “John likes Mary” belongs to the language of  $G_1$  because  $N, N \setminus (S/N), N \vdash S$  due to successive applications of the two elimination rules :  $X, X \setminus Y \vdash Y$  and  $Y, Y/X \vdash Y$ . Type constructors  $/$  and  $\setminus$  can be seen as oriented logic implications, the elimination rules are analogues of the “Modus Ponens” logic rule. An interesting issue is how the underlying rules or logics may compose (this is the design of logic functors) to deal with more fine-grained linguistic phenomenon.

Since they are lexicalized, such grammar formalisms seem well-adapted to automatic acquisition or completion perspectives. Such studies are performed in particular in Gold’s paradigm.

*Identification in the limit in the model of Gold* consists in defining an algorithm on a finite set of (possibly structured) sentences that converges to obtain a grammar in the class that generates the examples. Let  $\mathcal{G}$  be a class of grammars that we wish to learn from positive examples; let  $\mathcal{L}(G)$  denote the language associated with a grammar  $G$ ; a *learning algorithm* is a function  $\phi$  from finite sets of (structured) strings to  $\mathcal{G}$ , such that for any  $G \in \mathcal{G}$  and  $\langle e_i \rangle_{i \in \mathbb{N}}$  any enumeration of  $\mathcal{L}(G)$ , there exists a grammar  $G' \in \mathcal{G}$  such that  $\mathcal{L}(G') = \mathcal{L}(G)$  and  $n_0 \in \mathbb{N}$  such that  $\forall n > n_0 \phi(\{e_0, \dots, e_n\}) = G'$ .

Categorial grammars have close connections with logic. We therefore consider different ways of connecting computational linguistic data and LIS, such as implementing parameterized pregroups as a logic functor [10].

## 4 Application Domains

### 4.1 Geographical Information Systems

**Participants:** Soda Marème Cissé, Olivier Ridoux, Peggy Cellier, Erwan Quesseveur, François Le Prince, Sébastien Ferré.

Geographical Information Systems (GIS) is an important, fast developing domain of Information technology, and it is almost absent from INRIA projects. It is especially important for local communities (e.g. region and city councils).

Geographical information systems [LT92] handle information that are localized in space (*geolocalized*). GIS form an area which incorporates various technologies such as web, databases, or imaging. One characteristic of GIS is their organization as *layers*. This is inherited from the plastic sheets that were used until recently for drawing maps. A layer represents the road system, another the fluvial system, another the relief, etc. This is another instance of the tyranny of the dominant decomposition, and is not satisfactory: to which layer belong bridges, into which layer can we represent a multimodal network? Moreover, mining GIS is known to be difficult for the same reason; the layer structure makes inter layer relationships difficult to discover.

The first advantage of applying LIS to GIS is to allow cross-layer navigation. Another advantage is to permit a logical handling of scales. In current GIS systems, scales are treated as different layers, and it is difficult to keep the consistency between all layers that describe the same object. Another advantage that we have observed in a preliminary work is that LIS helps cleaning a data-base. This was not expected, and opens an interesting research area. Another characteristic of GIS is an intensive usage of topological relations (touches, overlaps, etc) and geographical relations (North, upstream, etc). Logic offers a rich language for expressing these relations and combining them.

## 4.2 Group Decision and Negotiation

**Participants:** Mireille Ducassé, Peggy Cellier.

Group decision and negotiation focuses on complex and self-organizing processes that constitute multiparticipant, multicriteria, ill-structured, dynamic, and often evolutionary problems. Group decision and negotiation refers to the whole process or flow of activities relevant to reaching a group decision, and not merely to the final choice - aspects of the process in group decision and negotiation include scanning, communication and information sharing, problem definition (representation) and evolution, alternative generation, and social-emotional interaction. Group decision support systems (GDSS) and negotiation support systems (GDNSS) are amongst the major approaches to address the problems.

In the current thread of research, we are showing that Logical Information Systems provide an innovative technological support for most of the above mentioned aspects of GDSS. In particular, the navigation and filtering capabilities of LIS help detect inconsistencies and missing knowledge during meetings. The updating capabilities of LIS enable participants to add objects, features and links between them on the fly. As a result the group has a more complete and relevant set of information. Furthermore, the compact views provided by LIS help participants embrace the whole required knowledge. The group can therefore build a shared understanding of the relevant information previously distributed amongst the participants. Lastly, the navigation and filtering capabilities of LIS are relevant to quickly converge on a reduced number of targets. A future trend of research will be to investigate how LIS can also support negotiation.

TODO: SPARKLIS

---

[LT92] R. LAURINI, D. THOMPSON, *Fundamentals of Spatial Information Systems*, Elsevier, Academic Press Limited, 1992.

## 5 Software

### 5.1 Camelis, Sewelis, and Sparklis

**Participants:** Sébastien Ferré.

Camelis is a stand-alone application that allows to store, retrieve and update objects through a graphical interface. Its main purpose is to experiment with the LIS paradigm. In particular, it has been very useful for refining the query-answer principle in special circumstances (e.g. when there are many answers, or when there are few answers). It is currently used as a personal storage device for handling photos, music, bibliographical references, etc, up to tens of thousands of objects. It implements as closely as possible the LIS paradigm. It is generic w.r.t. logics, and is compatible with our library of logic functors, LogFun (see Section 5.2). It is available on Linux and Windows, and comes with a user manual.

An important extension, Sewelis, has been developed to browse RDF(S) graphs, a Semantic Web standard. It uses a query language whose expressivity is similar to SPARQL, the reference query language of the Semantic Web. The LIS navigation has been proved safe (i.e., does not lead to dead-ends), and complete (i.e., can reach all conjunctive queries), so that users can perform complex searches easily and safely [4]. Sewelis also supports the guided creation and update of objects, according to the UTILIS approach [11].

Sparklis is a re-implementation of the querying capabilities of Sewelis as a Web application on top of SPARQL endpoints. Its main advantages compared to Sewelis are: (1) no setup, just load a page in a browser, (2) direct exploration of remote SPARQL endpoints, (3) scaling up to large triples that contain up to billions triples, like DBpedia, and (4) verbalization of queries in natural language for better readability.

### 5.2 LogFun

**Participants:** Sébastien Ferré.

The formal definition of a LIS is generic with respect to the logic used for object descriptions and for queries. The counterpart is that it is up to the user to design and implement a logic solver to plug in a LIS. This is too demanding on the average user, and we have developed a framework of *logic functors* that permits to build *certified* logic solvers (see Section 3.5).

LogFun is a library of *logic functors* and a *logic composer*. A user defines a logic using the logic functors, and produces a certified software implementation of the logic (i.e., parser, printer, prover) by applying the logic composer to the definition. For instance, using a functor *Interval* for reasoning on intervals (e.g.  $x \in [2, 5] \implies x \in [0, 10]$ ), and a functor *Prop* for propositional reasoning (e.g.  $a \wedge b \implies a$ ), a user can define logic *Prop(Interval)*. In this logic, a theorem like  $x \in [2, 5] \vee x \in [7, 9] \implies x \in [0, 10]$  can be proven. Note that  $[2, 5] \cup [7, 9]$  is not an interval, so that *Prop(Interval)* is an actual extension over *Interval*.

What the logic composer does when building logic *Prop(Interval)* is to compose the solver of *Interval* and the generic solver of *Prop*, and build a solver for *Prop(Interval)*. It also type-checks *Prop(Interval)* to produce its certificate using the certificates of *Interval* and *Prop*. In this example, the certificate says that *Prop(Interval)* is complete: everything that could be

deduced from the meaning of  $Prop(Interval)$  can be proved by its solver. In other circumstances, the certificate indicates that the logic defined by the user is incomplete, w.r.t. the semantics and solvers that come with the functors. In this case, the certificate also indicates what hypotheses are missing for completeness; this may help users to define a more complete variant of their logic.

Logic functors offer basic bricks and a building rule to safely design new logics. For instance, in a recent application of LIS to geographical information system, a basic reasoning capability on locations was needed. The designer of the application, not a LIS or LogFun author, could build a relevant ad hoc logic safely and rapidly.

### 5.3 Portalis

**Participants:** Yves Bekkers, Benjamin Sigonneau.

Portalis provides *on the shelf* software bricks that can be used to construct web services built on top of our logical information systems. The purpose of this sub-project is to facilitate scientific popularization and industrial transfer of tools produced by the LIS team. The logical information systems are written in OCaml, and need adapters to more standard interfaces.

The first software brick is an HTTP server that wraps Camelis core functions and provides a dedicated API to access them. We call it the *LIS server*. Each Camelis function is called using an HTTP request. The answer is then marshalled in a dedicated XML format and sent back as an HTTP response.

A second brick implements a Java layer on top of the first brick. It is meant to easily build clients. It provides a Java function for each Camelis function. The Java function serializes its parameters, calls the corresponding HTTP request, and deserializes the XML answer in the form of Java objects.

On top of that, a portal offers features that are essential to several real-world use cases: using different Camelis contexts at once, handling users and their access rights. It also exports an HTTP API to allow people who would want to write clients in a language different from Java to benefit from these functions.

### 5.4 Typed grammars

**Participants:** Denis Béchet [LINA-Nantes], Annie Foret [contact point].

A Pregroup ToolBox is under development on the gforge Inria as a collaborative work with LINA. It includes a generic pregroup parser (LINA) and grammar lexicon definitions and manipulation tools based on XML. An interface with Camelis has been developed (from Camelis to the Pregroup XML format, and the other way round). It has been used to define and experiment grammar prototypes for different natural languages.

### 5.5 SQUALL: a Semantic Query and Update High-Level Language

**Participants:** Sébastien Ferré.

SQUALL (Semantic Query and Update High-Level Language) is a controlled natural language (CNL) for querying and updating RDF graphs [4]. The main advantage of CNLs is to reconcile the high-level and natural syntax of natural languages, and the precision and lack of ambiguity of formal languages. SQUALL has a strong adequacy with RDF, and covers all constructs of SPARQL, and most constructs of SPARQL 1.1. Its syntax completely abstracts from low-level notions such as bindings and relational algebra. It features disjunction, negation, quantifiers, built-in predicates, aggregations with grouping, and n-ary relations through reification.

SQUALL is available as a Web application at <http://lisfs2008.irisa.fr/ocsigen/squall/> under two forms: one that translates SQUALL sentences to SPARQL, and another one that directly return query answers from a SPARQL endpoint.

## 5.6 PEW: Possible World Explorer

**Participants:** Sébastien Ferré, Sebastian Rudolph.

The Possible World Explorer (PEW) targets ontology designers, and aims to help them correct and complete their ontologies. It reuses the query-based faceted search principles of Sewelis for exploring the “possible worlds” (i.e., models) of an OWL ontology. Users are guided in the incremental construction of class expressions, such that only satisfiable classes are reachable. All classes made of qualified existential restrictions, nominals, intersections, unions, and atomic negations are reachable.

PEW not only supports the exploration of an ontology’s possible worlds, but also supports its completion by the addition of axioms [8]. When a class is found satisfiable, and this contradicts domain knowledge (e.g., a man that is not a person), the undesirable possible worlds can be excluded (“pew pew!”) by asserting an axiom saying that this class is unsatisfiable (e.g., every man is a person). This could be made a game, where the player would strive to exclude as many undesirable worlds as possible. The benefits are to complete the ontology with more knowledge, and therefore to improve its deduction power.

In addition to asserting negative axioms (about things that should not exist), PEW also allows for the definition of named classes (OWL equivalent class axioms), and for the creation of named individuals as instances of class expressions (OWL class assertion axioms).

## 6 New Results

### 6.1 Abstract Conceptual Navigation (ACN)

**Participants:** Sébastien Ferré.

What follows is the abstract of Ferré’s habilitation thesis [1]. It is a synthesis of various works aiming at reconciling expressivity and usability in information access. It introduces Abstract Conceptual Navigation (ACN) as a common framework to those works. The habilitation was defended on November 6th in front of the following committee: Olivier Pivert (Univ. Rennes 1, president), Norbert E. Fuchs (Univ. Zurich, referee), Marianne Huchard

(Univ. Montpellier 2, referee), Eero Hyvönen (Univ. Aalto, referee), Karell Bertet (Univ. La Rochelle, examiner), and Fabien Gandon (INRIA Valbonne, examiner).

In many domains where information access plays a central role, there is a gap between expert users who can ask complex questions through formal query languages (e.g., SQL), and lay users who either are dependent on expert users, or must restrict themselves to ask simpler questions (e.g., keyword search). Because of the formal nature of those languages, there seems to be an unescapable trade-off between expressivity and usability in information systems. The objective of this thesis is to present a number of results and perspectives that show that the expressivity of formal languages can be reconciled with the usability of widespread information systems (e.g., browsing, Faceted Search (FS)). The final aim of this work is to empower people with the capability to produce, explore, and analyze their data in a powerful way.

We have proposed a number of theories and implementations to better reconcile expressivity and usability, and applied them to a number of contexts going from file systems to the Semantic Web. In this thesis, we introduce an unifying framework inspired by Formal Concept Analysis (FCA) to factor out the main ideas of all those results: Abstract Conceptual Navigation (ACN). The principle of ACN is to guide users by letting them *navigate* in a conceptual space where places are *concepts* connected by navigation links. Concepts are characterized by a formal query, and are made of two parts: an *extension* and an *intension*. The extension is made of query results while the intension is made of the query itself and an index of query increments over results. Finally, navigation links are formally defined as query transformations. The conceptual space is not static but is induced by concrete data, and evolves with it. ACN therefore combines the *expressivity* of formal query languages with the *guidance* of conceptual navigation. The *readability* of queries is improved by verbalizing them to (or parsing them from) a Controlled Natural Language (CNL). Readability and guidance together support usability by speaking user's language, and by providing a systematic assistance.

## 6.2 Expressive and Scalable Query-Based Faceted Search over SPARQL Endpoints

**Participants:** Sébastien Ferré.

Linked data is increasingly available through SPARQL endpoints, but exploration and question answering by regular Web users largely remain an open challenge. Users have to choose between the expressivity of formal languages such as SPARQL, and the usability of tools based on navigation and visualization. In a previous work, we have proposed Query-based Faceted Search (QFS) as a way to reconcile the expressivity of formal languages and the usability of faceted search. In this work [10], we further reconciled QFS with scalability and portability by building QFS over SPARQL endpoints. We also improved expressivity and readability. Many SPARQL features are now covered: multidimensional queries, union, negation, optional, filters, aggregations, ordering. Queries are now verbalized in English (and French), so that no knowledge of SPARQL is ever necessary.

All of this is implemented in a portable Web application, Sparklis, and has been evaluated on many endpoints and questions. A demonstration [11] has been given at the Semantic Web

conference (ISWC). About 10,000 navigation steps have been performed by users from all over the world on various endpoints in the last six months. The demo video on YouTube has been viewed 360 times in 38 countries since April 2014.

### 6.3 Guided Composition of Bioinformatic Workflows with Logical Information Systems

**Participants:** Mouhamadou Ba, Sébastien Ferré, Mireille Ducassé.

In a number of domains, particularly in bioinformatics, there is a need for complex data analysis. For that issue, elementary data analysis operations called tasks are composed as workflows. The composition of tasks is however difficult due to the distributed and heterogeneous resources of bioinformatics. Heterogeneity of data and data formats in bioinformatics entail mismatches between inputs and outputs of different services, making it difficult to compose them into workflows. To reduce those mismatches, bioinformatics platforms propose ad hoc converters, called shims. When shims are written by hand, they are time-consuming to develop, and cannot anticipate all needs. When shims are automatically generated, they miss transformations, for example data composition from multiple parts, or parallel conversion of list elements.

This year, we worked on convertibility between input and output types for composition of services in bioinformatics. This work proposes to systematically detect convertibility from output types to input types. Convertibility detection relies on a rule system based on abstract types, close to XML Schema. Types allow to abstract data while precisely accounting for their composite structure [6, 5]. Based on convertibility detection, we presented an approach [7] that also addresses the generation of actual converters between input and output XML data. We show the applicability of our approach by abstracting concrete bioinformatics types (e.g., complex biosequences) for a number of bioinformatics services (e.g., blast). We illustrate how our automatically generated converters help to resolve data mismatches when composing workflows.

The objective is to use semantics to describe bioinformatic tasks and to adapt the guided approach of Sewelis, a LIS semantic web tool, to the composition of tasks. We aim at providing a tool that supports guided composition of semantic web services in bioinformatics, and that will support biologists in designing workflows for complex data analysis.

### 6.4 Group Decision Support Systems

**Participants:** Mireille Ducassé, Peggy Cellier.

Group work represents a large amount of time in professional life while many people feel that much of that time is wasted. This amount of time is still increasing because problems are becoming more complex and are meant to be solved in a distributed way. Each involved person has a local and partial view of the problem, no one embraces the whole required knowledge. In this thread of research we develop decision processes taking benefits of Logical Information Systems capabilities.

Reasoning on multiple criteria is a key issue in group decision to take into account the multidimensional nature of real-world decision-making problems. In order to reduce the induced information overload, in multicriteria decision analysis, criteria are in general aggregated, in many cases by a simple discriminant function of the form of a weighted sum. It requires to, a priori and completely, elicit preferences of decision makers. That can be quite arbitrary. In everyday life, to reduce information overload people often use a heuristic, called “*Take-the-best*”: they take criteria in a predefined order, the first criterion which discriminates the alternatives at stake is used to make the decision [MGG10]. Although useful, the heuristic can be biased. We propose the Logical Multicriteria Sort process to support multicriteria sorting within islands of agreement. It therefore does not require a complete and consistent a priori set of preferences, but rather supports groups to quickly identify the criteria for which an agreement exists. The process can be seen as a generalization of Take-the-best. It also proposes to consider one criterion at a time but once a criterion has been found discriminating it is recorded, the process is iterated and relevant criteria are logically combined. Hence, the biases of Take-the-best are reduced. The process is supported by Logical Information Systems, which give instantaneous feedbacks of each small decision and keep tracks of all of the decisions taken so far. The process is incremental, each step involves low information load. It guarantees some fairness because all considered alternatives are systematically analyzed along the selected criteria [3].

We also investigate how Logical Information Systems can help in multi-unit assignment problems, where a set of indivisible resources is to be allocated amongst a set of agents, the agents having multi-unit demands. One instance of such a problem is course allocation at universities, the agents are students and the resources are seats in courses. Dictatorship based processes are commonly used for single-unit assignments: agents are totally ordered, then each agent in turn chooses her best choice. It can, however, be highly unfair for multi-unit assignments: top agents get all their choices while bottom agents get none of them. Other widespread approaches use bids or preference ranking. There are cases, however, where agents want to be able to give more qualitative information about their wishes. At present, no multi-unit assignment system supports both quantitative and qualitative information. We propose an interactive process for multi-assignment problems where, in addition to bids and preferences, agents can give arguments to motivate their choices. Bids are used to automatically make pre-assignments, qualitative arguments and preferences help decision makers break ties in a founded way. A LIS system is used, it allows decision makers to easily handle bids, arguments and preferences in a unified interface. Agents are dealt with by sets, for example the set of those who have bidden a given number of points for a resource and who have given qualitative arguments. We say that a process is *p\_equitable* for a property *p* if all the agents satisfying *p* are treated equally. We have formally demonstrated that our process is *p\_equitable* for a clear list of properties on bids, arguments and preferences. It is also Pareto-efficient with respect to bids and Gale-Shapley-stable with respect to bids. A successful case study about a course assignment problem at a technical university has been led. It spans over two university years. The first year, after a contested assignment, the process has been designed and a simulation has been played [9]. The second year, the process has been revised according to the students

---

[MGG10] J. N. MAREWSKI, W. GAISSMAIER, G. GIGERENZER, “Good judgments do not require complex cognition”, *Cognitive Processing* 11, 2, 2010, p. 103–121.

and decision maker's feedbacks. An actual assignment has been achieved using the designed process. The decisions makers reported feeling serene about the process and confident about the resulting assignment. Furthermore, the students, even the ones who did not get all their wishes, considered that the process is equitable.

## 6.5 Segmentation of Geolocalized Trajectories using Exponential Moving Average

**Participants:** Soda Cissé, Peggy Cellier, Olivier Ridoux.

Following the track of the treatment of GIS operations using LIS approaches, we explore the mining of trajectories of mobile objects. Nowadays, large sets of data describing trajectories of mobile objects are made available by the generalization of geolocalisation sensors. Relevant information, for instance, the most used routes by children to go to school or the most extensively used streets in the morning by workers, can be extracted from this amount of available data allowing, for example, to reconsider the urban space. These trajectories have contextual attributes that are fairly easily amenable to a treatment by formal methods like FCA and LIS. However, they also have attributes which represent physical dimensions like time and space, or ambient parameters like temperature. A trajectory is represented by a set of points  $(x; y; t)$  where  $x$  and  $y$  are the geographic coordinates of a mobile object and  $t$  is a date. These data are difficult to explore and interpret in their raw form, i.e. in the form of points  $(x; y; t)$ , because they are noisy, irregularly sampled and too low level. A first step to make them usable is to resample the data, smooth it, and then to segment it into higher level segments (e.g. "stops" and "moves") that give a better grip for interpretation than the raw coordinates. We have proposed a method [8] for the segmentation of these trajectories in accelerate/decelerate segments which is based on the computation of exponential moving averages (EMA). We have conducted experiments where the exponential moving average proves to be an efficient smoothing function, and the difference between two EMA of different weights proves to discover significant accelerating-decelerating segments.

## 6.6 Extraction of Relations between Genes and Rare Diseases

**Participants:** Peggy Cellier, Nicolas Béchet [IUT Vannes], Thierry Charnois [University of Paris 13], Bruno Crémilleux [University of Caen].

Orphanet provides an international web-based knowledge portal for rare diseases including a collection of review articles. However, reviews and literature monitoring are manual. Thus, new documentation about a rare disease is a time-consuming process and automatically discovering knowledge from a large collection of texts is a crucial issue. This context represents a strong motivation to address the problem of extracting gene-rare diseases relationships from texts. We tackle this issue with a cross-fertilization of information extraction and data mining techniques (sequential pattern mining under constraints) [2]. Experiments show the interest of the method for the documentation of rare diseases.

## 6.7 Data Mining to Associate Scientific Papers with their Session Name

**Participants:** Peggy Cellier, Thierry Charnois [University of Caen], Solen Quiniou [University of Nantes].

We present a proposition based on data mining to tackle the DEFT 2014 challenge [14]. We focus on task 4 which consists of identifying the right conference session for scientific papers. The proposed approach is based on a combination of two data mining techniques. Sequence mining extracts frequent phrases in scientific papers in order to build paper and session descriptions. Then, those descriptions of papers and sessions are used to create a graph which represents shared descriptions. A graph mining technique is applied on the graph in order to extract a collection of homogenous sub-graphs corresponding to sets of papers associated to sessions.

## 6.8 Type-logical Grammar Formalisms

**Participants:** Annie Foret.

Type-logical Grammar formalisms are used in computational linguistics to model syntax using a strongly lexicalized style, properties (logical types) being attached directly to words. Type-logical grammars are also connected with logic, especially linear logic. For instance, parsing a sentence is similar to searching a proof, and semantics can be transparently mapped onto structure in the style of the Curry-Howard isomorphism.

A significant part of this work is connected to automatic acquisition (learning, grammatical inference) issues, and possibilities to help in the design of valuable grammars.

**Categorial grammar hierarchies.** The notion of  $k$ -valued categorial grammars in which every word is associated to at most  $k$  types is often used in the field of lexicalized grammars as a fruitful constraint for obtaining interesting properties like the existence of learning algorithms. This constraint is reasonable only when the classes of  $k$ -valued grammars correspond to a real hierarchy of generated languages. Such a hierarchy has been established earlier for the classical categorial grammars. The new contribution [12] studies an extension of Lambek grammars with respect to such hierarchies. Another interest of the extended calculus is to allow some parallels in grammar design (type assignments, acquisition methods) between both frameworks (pregroups and  $(L)$ ). Those extensions allow basic proper axioms (a preorder), this feature is of particular interest for a LIS/CAMELIS view of a grammar and manage related linguistic data.

**Acquiring and managing linguistic resources.** We have studied in [13] mappings between different grammar formalisms to allow a transfer of linguistic resources from one formalism to the other. This work has been done at ILCC, Edinburgh.

We have considered in parallel two families of categorial formalisms, with a view to providing some convergence both at the formal level and at the level of resources (via formally rooted transfers, and some hypotheses that can be experimentally tested). These families are on one side (1) the family of combinatory categorial grammars (CCG) [Steedman] and

on the other side (2) a family of Lambek-like grammars, pregroups (PG) [Lambek] that have been introduced as a simplification of Lambek calculus.

We have focussed on mappings that preserve the binary structures, for preserving parse structures, we have also discussed some possible alternatives in the underlying formalisms, with some experiments.

## 6.9 SQUALL: The expressiveness of SPARQL 1.1 made available as a controlled natural language

**Participants:** Sébastien Ferré.

The Semantic Web (SW) is now made of billions of triples, which are available as Linked Open Data (LOD) or as RDF stores. The SPARQL query language provides a very expressive way to search and explore this wealth of semantic data. However, user-friendly interfaces are needed to bridge the gap between end-users and SW formalisms. Navigation-based interfaces and natural language interfaces require no or little training, but they cover a small fragment of SPARQL's expressivity. We propose SQUALL, a query and update language that provides the full expressiveness of SPARQL 1.1 through a flexible controlled natural language (e.g., solution modifiers through superlatives, relational algebra through coordinations, filters through comparatives). A comprehensive and modular definition is given as a Montague grammar, and an evaluation of naturalness is done on the QALD challenge. SQUALL is conceived as a component of natural language interfaces, to be combined with lexicons, guided input, and contextual disambiguation. It is available as a Web service that translates SQUALL sentences to SPARQL, and submits them to SPARQL endpoints (e.g., DBpedia), therefore ensuring SW compliance, and leveraging the efficiency of SPARQL engines.

Compared to previous work on SQUALL, we have improved the coverage of SPARQL 1.1 features so that only transitive closure in property paths and a few minor things are missing. We have also extended the formalization and evaluation of SQUALL, which are now completely presented in detail in *Data & Knowledge Engineering* [4]. That paper was invited as an extended version of a previous paper at NLDB'13.

## 7 Contracts and Grants with Industry

### 7.1 IDFRAud: An Operational Automatic Framework for Identity Document Fraud Detection and Profiling (ANR)

**Participants:** Sébastien Ferré, Peggy Cellier.

Le projet ANR IDFRAud vise à permettre la reconnaissance automatique de documents d'identité et la détection de fraudes, en appliquant des techniques d'analyse de documents, de classification et de gestion de connaissances. Le porteur est Abdullah Almaksour de l'entreprise innovante AriadNEXT et les autres partenaires sont l'IRISA, l'IRCGN (institut de recherche criminelle de la gendarmerie nationale) et l'ENSP (école de police). Sébastien Ferré est le responsable scientifique pour le partenaire IRISA. Une première réunion a eu lieu le 24 septembre

et le projet doit débiter février 2015. Ce projet finance pour l'équipe LIS deux ans de postdoc, plus des frais de fonctionnement (ex., missions, machines).

## 7.2 Cour de Cassation

**Participants:** Annie Foret.

Annie Foret is a member of a group working with "The French Cour de Cassation". This multi-disciplinary project on "Développement d'un prototype de logiciel d'aide à la prise de décision lors de la rédaction de jugements" has received funds from ENS CHRC (Collège de Recherche Hubert Curien) (proposal submitted by François Schwarzentruher at ENS Rennes). The project also received funds from Rennes 1 University as a "Défi émergent" (proposal submitted by A. Foret).

## 7.3 Portalis: funding for maturation (FEDER Région Bretagne)

**Participants:** Yves Bekkers, Benjamin Sigonneau, Sébastien Ferré.

As part of the implementation of the pole "Regional Competitiveness and Employment (2007-2013)" in Brittany, we obtained a funding from FEDER and Région Bretagne for an engineer for a year to participate in Camelis transfer to industry. Benjamin Sigonneau has been appointed on this founding. The development started in October 2012 and finished in June 2014. The result is the embedding of Camelis into a HTTP server, i.e. a portal for LIS, hence the name Portalis (see 5.3). A key contribution of Portalis w.r.t. Camelis is the support of collaborative access, authoring, and administration of logical information systems. Users can be registered and given different access rights.

We had discussions with Mediadone, a company specializing in processing, indexing and image enhancement. Mediadone provides tools for interactive and enriched WebTV. The company is interested in using Portalis bricks to build an intelligent navigation tool based on the use of Camelis. Mediadone has already acquired a license for the commercial exploitation of Camelis.

# 8 Other Grants and Activities

## 8.1 International Collaborations

- In 2013-2014, during a "délégation CNRS", Annie Foret was a visitor of *The Institute for Language, Cognition and Computation (ILCC)* in Edinburgh. The visit was hosted by Mark Steedman's group, working on Categorical Grammars.
- Sébastien Ferré is a member of the management committee of the COST action MUMIA (IC1002 - Multilingual and multifaceted interactive information access). MUMIA aims to coordinate collaboration between the following disciplines: machine translation, information retrieval, and faceted search. The objectives of the action is to foster research and development for next generation search technologies. The domain of patent search has

been selected as a common use case, as it provides highly sophisticated and information intensive search tasks that have significant economic ramifications.

In 2014, he gave a course at the 3rd MUMIA Training School, which was located at the ICS FORTH lab, Heraklion, Crete. The title of the course was on “Query-based Faceted Search: Expressive and Guided Information Access to the Semantic Web”. The course was 2 hours, and part of it was dedicated to a hand-on session with Sparklis.

## 8.2 National Collaborations

- Annie Foret is an external collaborator of LINA (research lab. Nantes), in TALN team (Natural Language Processing), and member of “Agence Universitaire de la Francophonie” (AUF), LTT network on “Lexicologie, terminologie et traduction”. Annie Foret is member of ATALA (Association pour le Traitement automatique des Langues), and of SIF (Société Informatique de France).
- Peggy Cellier collaborates with the members of the ANR project HYBRID<sup>1</sup> on the part about data mining for natural language processing.
- Peggy Cellier and Sébastien Ferré are involved in an ANR project preproposal, ECONOM, in collaboration with LIPN in Paris, LI in Tours, and the Noopsis company. The title is “Cross-fertilization of data mining, natural language processing, and semantic web for business intelligence”.

## 9 Dissemination

### 9.1 Scientific Responsibilities

- Olivier Ridoux, Annie Foret, and Sébastien Ferré are members of the committee of the DKM scientific department (Data and Knowledge Management) at IRISA. Olivier Ridoux is the head of the department since October 2014.
- Annie Foret has been elected as a member of the scientific committee of ISTIC-Rennes1. She is a member of the committee "Développement Durable (Sustainable development)".
- Olivier Ridoux is a member of the EcoInfo CNRS service group on sustainable development and information technology.
- Sébastien Ferré and Peggy Cellier were members of the organization committee of EGC 2014 held in Rennes in January 2014. The presidents of the committee were Arnaud Martin and René Quiniou.
- Mireille Ducassé has served in the program committee of GDN 2014, international conference on Group Decision and Negotiation, Toulouse. She served in one PhD defense committees: president for Santiago Videla, Université Rennes 1, Torsten Schaub and Anne Siegel, PhD co-supervisors, July 2014. She has been a member of 2 recruitment

---

<sup>1</sup><http://hybride.loria.fr>

committees (“comité de sélection”) in computer science at university of Caen and INSA Rennes.

- Olivier Ridoux has served as a member of one PhD committees. He is a member of the "Conseil de laboratoire" at IRISA.
- Sébastien Ferré is a member of the Editorial Board of the International Conference on Formal Concept Analysis (ICFCA). He was also in 2014 a member of the program committee of: the conference CLA (Concept Lattices and Applications), and the workshop CNL (Controlled Natural Language). Finally, he served as an external reviewer for the journals *Informatica*, *INS* (Information Sciences), and *JCSS* (Journal of Computer and System Sciences).

Sébastien is a supervisor of the PhD of Mouhamadou Ba. He served as an examiner in the PhD defense committee (Montpellier 2, June 24th) of Rafat Almsiedeen on “Extraction de modèles de variabilité multi-vues : application aux applications mobiles”; and as an invited member in the PhD committee (La Rochelle, November 24th) of Clément Guérin on “Proposition d’un cadre pour l’analyse automatique, l’interprétation et la recherche interactive d’images de bande dessinée”. Finally, he is a member of the PhD committee of François Moreews, a PhD student in team Genscale at IRISA.

Sébastien is responsible of the organization of full-day seminars of the DKM department. In 2014, two seminars were organized: on June 16th about research softwares developed in the department, and on December 8th about the respective positioning of department teams w.r.t. to key data and knowledge management themes.

- Peggy Cellier was one of the six organizers of the ECML/PKDD workshop: DMNLP 2014 (Workshop on Interactions between Data Mining and Natural Language Processing): <http://dmnlp.loria.fr>. Peggy is a member of the Editorial Board of the International Conference on Formal Concept Analysis (ICFCA). She was also in 2014 a member of the program committee of: the conference ICCS (Concept Lattices and Applications), the conference ICFCA, and the PhD symposium of PKDD. She also served as an external reviewer for the conference ICSME (International Conference on Software Maintenance and Evolution). She also served as an examiner in the PhD defense committee (University of Grenoble, June 12th) of Sofiane Lagraa on “Nouveaux outils de profilage de MPSoC basés sur des techniques de fouille de données”. She served as a local member of the selection committee for an associate professor position at University of Rennes 1 (ref. MCF1071), and as an external member of the selection committee for an associate professor position at the University of François Rabelais de Tours (ref. MCF0610). Peggy is also a supervisor of the PhD of Soda Marème Cissé with Olivier Ridoux. She is a member of the “Conseil de Composante IRISA/INSA” at INSA.
- Annie Foret has been a *program committee* member of the *Formal Grammar* 2014 International Conference. She now belongs to the Standing committee of the *Formal Grammar* international conference. She is a *program committee* member of the 4th International Conference on Tools for Teaching Logic. She was member of a thesis committee for Ophélie Lacroix in Nantes (annual report before the PHD).

## 9.2 Involvement in the Scientific Community

- Mireille Ducassé has given invited talks at the Technical University Dresden in September 2014, and at the CREM laboratory of Caen in December 2014.
- Olivier Ridoux is active in the development of the "Maison pour la science en Bretagne" for the scientific training of primary and secondary teachers in Brittany.
- Peggy Cellier gave an invited joint-talk with Thierry Charnois and Damien Nouvel to the AFIA-ATALA Workshop "Journée TAL et IA" organized by Brigitte Grau and Pierre Zweigenbaum at INALCO (Paris)<sup>2</sup>: "Data Mining for NLP". She also gave an invited joint-talk with Thierry Charnois to a local workshop at University of Caen "Mini-Symposium on Data Mining" organized by Bruno Crémilleux<sup>3</sup>: "Sequential Pattern Mining under Constraints for NLP". Finally, she gave an invited talk at INSERM Rennes to the local seminar: "Fouille de données pour la découverte de relations entre entités biologiques dans les textes".
- Mouhamadou Ba and Sébastien Ferré took part in the ReNaBiGo seminar on workflows, organized by the GenOuest platform and closely related to the PhD subject of Ba.

## 9.3 Teaching

- Mouhamadou Ba teaches algorithms and Java at Licence 1 level, databases at Licence 1 level, and data mining at Master 2 level, all at Insa rennes.
- Yves Bekkers teaches programming languages: functional programming, logic programming (Prolog, lambdaProlog), artificial intelligence, relational databases, XML, new technologies for programming distributed applications on the Web. After being a specialist of logic programming, his current principal interest is teaching distributed application design using tools based on model technologies (Object programming, XML, SGBDR) which allows building applications from one need to the other using numerous bridges allowing passing automatically from one model to another. Yves Bekkers retired from active service at university in October 2014.
- At Insa, Peggy Cellier teaches algorithms and Java, database and system at licence level; symbolic data mining and formal methods for software engineering at Master level. She organized the bibliographic and internship defense for the Research in Computer Science (MRI) specialism. Since September, 2013 she is responsible of the internships of computer science students (Licence 3, Master 1 and Master 2). She has also been involved in the VAE diploma of INSA Rennes. She served as examiner for the evaluation committee of Pascal Lochert and as referee for Oscar Bastos. Finally, she served as local examiner for the recruitment of students at INSA.
- Soda Marème Cissé teaches office automation, imperative programming in Java at Licence 1 level of University Rennes 1. She also teaches databases at Licence 2 level of

---

<sup>2</sup><http://www.afia.asso.fr/tiki-index.php?page=Journée+commune+AFIA+-+ATALA+2014>

<sup>3</sup><https://www.greyc.fr/fr/node/2045/>

University Rennes 1. She also participated in green IT courses given by Olivier Ridoux as part of doctoral programs organized by MATISSE doctoral school.

- Mireille Ducassé is the director of international relations of the INSA of Rennes since december 2010. As such, she is a member of the direction of the Insa of Rennes. Since March 2014 she is also the coordinator of the international committee of Groupe INSA.

In 2014, she has set up a dual master degree programme in computer science with the Technical University of Dresden. She has also set up a dual master degree framework for six of the specialities of INSA Rennes with Politecnico Milano, as well as for all the specialities of INSA Rennes with ETS Montréal. She is responsible of a BRAFITEC agreement, renewed in 2014, for Groupe INSA with Universidade Federal da Paraiba (UFPB) and Universidade Federal de Campina Grande (UFCG), Brazil.

At Insa, she is responsible of three courses, taught in English: *Formal Methods for Software Engineering* (with the “B formal method”) and *Constraint Programming* at Master 1 level, as well as *Participatory Design* at Master 2 level.

- Sébastien Ferré teaches symbolic data mining and compilation at the master level. He also teaches formal methods for programming and software engineering at the license level. He started a course on the Semantic Web for master-level students in MIAGE. He is vice-director of the MIAGE at ISTIC, and is in charge of Master 1 internships (80 students).

He was a member of two committees for the validation of grades based on experience (VAE) for senior IT workers.

- Annie Foret teaches university courses including formal logic, object-oriented programming, XML technologies and databases.
- Olivier Ridoux teaches data-bases, algorithmics, the theory of formal languages and compilation in the engineering school ESIR. He is also in charge of the innovation training in the school, in which he also teaches sustainable development w.r.t. IT, and disruptive innovation (à la Clayton Christensen) w.r.t. scientific revolution (à la Thomas Kuhn). He also teaches logic and constraint programming at the Master level, and an introduction to the principles of IT systems at the Bachelor level. He is also in charge of the institutional communication of ESIR.

## 10 Bibliography

### Major publications by the team in recent years

- [1] D. BÉCHET, A. FORET, “A Pregroup Toolbox for Parsing and Building Grammars of Natural Languages”, *Linguistic Analysis Journal* 36, 2010.
- [2] P. CELLIER, S. FERRÉ, O. RIDOUX, M. DUCASSÉ, “A Parameterized Algorithm to Explore Formal Contexts with a Taxonomy”, *Int. J. Foundations of Computer Science (IJFCS)* 19, 2, 2008, p. 319–343.

- [3] M. DUCASSÉ, S. FERRÉ, “Fair(er) and (almost) serene committee meetings with Logical and Formal Concept Analysis”, *in: Proceedings of the International Conference on Conceptual Structures*, P. Eklund, O. Haemmerlé (editors), Springer-Verlag, July 2008. Lecture Notes in Artificial Intelligence 5113.
- [4] S. FERRÉ, A. HERMANN, “Reconciling faceted search and query languages for the Semantic Web”, *Int. J. Metadata, Semantics and Ontologies* 7, 1, 2012, p. 37–54.
- [5] S. FERRÉ, R. D. KING, “A dichotomic search algorithm for mining and learning in domain-specific logics”, *Fundamenta Informaticae – Special Issue on Advances in Mining Graphs, Trees and Sequences* 66, 1-2, 2005, p. 1–32.
- [6] S. FERRÉ, O. RIDOUX, “A Framework for Developing Embeddable Customized Logics”, *in: Int. Work. Logic-based Program Synthesis and Transformation*, A. Pettorossi (editor), LNCS 2372, Springer, p. 191–215, 2002.
- [7] S. FERRÉ, O. RIDOUX, “An Introduction to Logical Information Systems”, *Information Processing & Management* 40, 3, 2004, p. 383–419.
- [8] S. FERRÉ, S. RUDOLPH, “Advocatus Diaboli - Exploratory Enrichment of Ontologies with Negative Constraints”, *in: Int. Conf. Knowledge Engineering and Knowledge Management (EKAW)*, A. ten Teije et al. (editor), LNAI 7603, Springer, p. 42–56, 2012.
- [9] S. FERRÉ, “Camelis: a logical information system to organize and browse a collection of documents”, *Int. J. General Systems* 38, 4, 2009.
- [10] A. FORET, “A modular and parameterized presentation of pregroup calculus”, *Information and Computation Journal* 208, 5, may 2010, p. 395–604.
- [11] A. HERMANN, S. FERRÉ, M. DUCASSÉ, “An Interactive Guidance Process Supporting Consistent Updates of RDFS Graphs”, *in: Int. Conf. Knowledge Engineering and Knowledge Management (EKAW)*, A. ten Teije et al. (editor), LNAI 7603, Springer, p. 185–199, 2012.

### Doctoral dissertations and “Habilitation” theses

- [1] S. FERRÉ, *Reconciling Expressivity and Usability in Information Access - From Filesystems to the Semantic Web*, Habilitation thesis, Matisse, Univ. Rennes 1, 2014, Habilitation à Diriger des Recherches (HDR), defended on November 6th.

### Articles in referred journals and book chapters

- [2] N. BÉCHET, P. CELLIER, T. CHARNOIS, B. CRÉMILLEUX, “Fouille de motifs séquentiels pour la découverte de relations entre gènes et maladies rares”, *Revue d’Intelligence Artificielle* 28, 2-3, 2014, p. 245–270.
- [3] M. DUCASSÉ, P. CELLIER, “Fair and Fast Convergence on Islands of Agreement in Multicriteria Group Decision Making by Logical Navigation”, *Group Decision and Negotiation* 23, 4, July 2014, p. 673–694, <http://dx.doi.org/10.1007/s10726-013-9372-4>.
- [4] S. FERRÉ, “SQUALL: The expressiveness of SPARQL 1.1 made available as a controlled natural language”, *Data & Knowledge Engineering* 94, 2014, p. 163–188.

## Publications in Conferences and Workshops

- [5] M. BA, S. FERRÉ, M. DUCASSÉ, “Convertibility between input and output types to help compose services in bioinformatics”, *in: Colloque africain sur la recherche en informatique et mathématiques appliquées (CARI)*, p. 141–148, 2014.
- [6] M. BA, S. FERRÉ, M. DUCASSÉ, “Convertibilité entre types d’entrée et de sortie pour la composition de services en bio-informatique”, *in: Conf. Reconnaissance de Formes et Intelligence Artificielle (RFIA)*, 2014.
- [7] M. BA, S. FERRÉ, M. DUCASSÉ, “Generating Data Converters to Help Compose Services in Bioinformatics Workflows”, *in: Int. Conf. Database and Expert Systems Applications (DEXA)*, H. D. et al. (editor), *LNCS 8644*, Springer, p. 284–298, 2014.
- [8] S. CISSÉ, P. CELLIER, O. RIDOUX, “Segmentation of Geolocalized Trajectories using Exponential Moving Average”, *in: Colloque Africain sur la Recherche en Informatique et Mathématiques Appliquées (CARI)*, p. 149–156, 2014.
- [9] M. DUCASSÉ, P. CELLIER, “Using Biddings and Motivations in Multi-unit Assignments”, *in: Group Decision and Negotiation. A Process-Oriented View*, P. Zaraté, G. E. Kersten, J. E. Hernandez (editors), *Lecture Notes in Business Information Processing, 180*, Springer, p. 53–61, 2014, <http://dx.doi.org/10.1007/978-3-319-07179-4-6>.
- [10] S. FERRÉ, “Expressive and Scalable Query-Based Faceted Search over SPARQL Endpoints”, *in: The Semantic Web (ISWC)*, P. Mika, T. Tudorache, A. Bernstein, C. Welty, C. A. Knoblock, D. Vrandečić, P. T. Groth, N. F. Noy, K. Janowicz, C. A. Goble (editors), *LNCS 8797*, Springer, p. 438–453, 2014. Nominee for the best research paper award.
- [11] S. FERRÉ, “SPARKLIS: a SPARQL Endpoint Explorer for Expressive Question Answering”, *in: ISWC Posters & Demonstrations Track*, M. Horridge, M. Rospocher, J. van Ossenbruggen (editors), *CEUR Workshop Proceedings, 1272*, CEUR-WS.org, p. 45–48, 2014.
- [12] A. FORET, “On Associative Lambek Calculus Extended with Basic Proper Axioms”, *in: Categories and Types in Logic, Language, and Physics - Essays Dedicated to Jim Lambek on the Occasion of His 90th Birthday*, C. Casadio, B. Coecke, M. Moortgat, P. Scott (editors), *LNCS 8222*, Springer, p. 172–187, 2014, [http://dx.doi.org/10.1007/978-3-642-54789-8\\_10](http://dx.doi.org/10.1007/978-3-642-54789-8_10).
- [13] A. FORET, “On Harmonic CCG and Pregroup Grammars”, *in: Int. Conf. Logical Aspects of Computational Linguistics (LACL)*, N. Asher, S. Soloviev (editors), *LNCS 8535*, Springer, p. 83–95, 2014, [http://dx.doi.org/10.1007/978-3-662-43742-1\\_7](http://dx.doi.org/10.1007/978-3-662-43742-1_7).
- [14] S. QUINIOU, P. CELLIER, T. CHARNOIS, “Fouille de données pour associer des noms de sessions aux articles scientifiques”, *in: Défi Fouille de Textes - DEFT 2014 (Atelier TALN)*, B. Bigi (editor), Laboratoire Parole et Langage, 2014. ISBN: 978-2-9518233-6-5.