



RESEARCH CENTER  
Lille - Nord Europe

FIELD

Activity Report 2013

# Section New Results

Edition: 2014-03-19



## ALGORITHMICS, PROGRAMMING, SOFTWARE AND ARCHITECTURE

- 1. ATEAMS Project-Team ..... 4
- 2. DREAMPAL Team ..... 6

## APPLIED MATHEMATICS, COMPUTATION AND SIMULATION

- 3. DOLPHIN Project-Team ..... 11
- 4. MODAL Project-Team ..... 16
- 5. NON-A Project-Team ..... 20
- 6. SequeL Project-Team ..... 26
- 7. SIMPAF Project-Team ..... 35

## DIGITAL HEALTH, BIOLOGY AND EARTH

- 8. BONSAI Project-Team ..... 37
- 9. SHACRA Project-Team ..... 38

## NETWORKS, SYSTEMS AND SERVICES, DISTRIBUTED COMPUTING

- 10. ADAM Project-Team ..... 43
- 11. FUN Project-Team ..... 44
- 12. RMOD Project-Team ..... 51

## PERCEPTION, COGNITION AND INTERACTION

- 13. LINKS Team ..... 54
- 14. MAGNET Team ..... 56
- 15. MINT Project-Team ..... 58

## ATEAMS Project-Team

# 5. New Results

## 5.1. Empirical analyses of source code

Rascal was used to perform empirical investigations of existing source code bases. First of all, Davy Landman performed an analysis of project management source code to investigate if domain knowledge is present in source code and, if so, how easy it is to extract that knowledge [26]. An earlier experiment in static analysis of PHP code was finalized by Mark Hills. The result is a deep study of feature usage in a large number of well-known PHP projects [25]. Vadim Zaytsev conducted an experiment to recognize micro-patterns in grammars and meta-models [32]. Finally, Jeroen van den Bos performed a deep empirical study to find out as to how far a domain-specific language facilitates evolution [34]. The results showed that the Derric DSL did indeed cover most evolution scenarios, but there is still room for improving the language. In all cases Rascal proved to be instrumental in performing the experiments.

## 5.2. Better parsing and disambiguation

Ali Afrozeh worked on a new implementation of GLL parsing, called Iguana. Unlike traditional parser generators, Iguana adopted the interpretive approach that is also used in the Ensō parser. This experiment is still ongoing, but the new parser is expected to be integrated into Rascal beginning of 2014. Additionally, a longstanding problem of disambiguation using operator precedence was solved [23]. Traditional approaches are either not safe (i.e. they make the language smaller), or they do not support complex precedence rules as found in, for instance, OCaml.

## 5.3. Extensible Programming

Modular and extensible implementation of languages could have major impact on how DSLs will be implemented. Anastasia Izmaylova continued here work on improving the extensibility of Rascal's module system, by providing open recursive function combinators.

Extensible programming is traditionally plagued by what has become known as “the expression problem”, which captures the fact that most programming languages either support extension of data variants, or extension of operations, but not both. Object Algebras are simple solution to this problem. In [30] we have extended this model to support feature-oriented programming. These results are currently being integrated into the Ensō system.

## 5.4. DSLs for Games

In collaboration with the Hogeschool van Amsterdam, Riemer van Rozen developed a workbench for MicroMachinations, a DSL for game economies [28]. Completely built using Rascal, this DSL environment features syntax highlighting, static analysis, interactive simulation, and SPIN-based model-checking of process models describing the economy of a game. The project shows the versatility of Rascal as a language workbench for the development of DSLs.

## 5.5. DSLs for Questionnaires

In the context of computational auditing we have intensified our research on DSLs for questionnaires. It was proposed by Tijs van der Storm as the benchmark task for the Language Workbench Challenge 2013 (LWC'13), which has resulted in a thorough overview and qualitative comparison of language workbenches [24]. As a side-effect, there are now two publicly available Rascal implementations of the questionnaire DSL (QL-R-Kemi and Demoqles). A first step has been made to collect all implementations to create a “chrestomathy” for further study and dissemination of language workbench concepts and DSL implementation patterns. Other results include a formal semantics of the dynamics of questionnaires [21], and an initial prototype of a questionnaire model for modeling the Dutch Tax Income filing application by Pablo Inostroza Valdera.

## 5.6. Live Programming

Live programming aims to bring the dynamic execution of programs closer to the programmer, ideally almost obliterating the gap between editing and executing the program. We are working on applying such principles in the context of DSLs. This has led to two results: a live programming environment for a DSL for questionnaires [36], and Trinity, a data-driven IDE for Derric [35]. Riemer van Rozen has worked on applying similar techniques to MicroMachinations, so that game economies can be adapted at runtime.

## 5.7. Visualization and interaction

Atze van der Ploeg worked on designing new algorithms and abstractions in the domain of visualization and abstraction. His first result is a fast algorithm for drawing non-layered, tidy trees [20]. DeForm is a library for the declarative specification of resolution-independent 2D graphics [27]. In [31] he proposed a reformulation of the traditional functional reactive programming (FRP) framework, which is both simple and efficient to implement.

## 5.8. Guarded Coroutines

Anastasia Izmaylova and Paul Klint have built an initial version of a compiler for Rascal. The performance improvements with respect to the interpreter are impressive. Moreover, the design of compiler is based on a new construct for implementing languages with complex backtracking and pattern matching semantics: guarded coroutines. This construct will be instrumental in extending the Rascal language with new kinds of control-flow and concurrency.

## 5.9. Data structures for meta programming

The efficiency of many meta programs is dependent on the internal data structures used to represent collections, trees, relations etc. Michael Steindorfer has worked on comparing the performance of various persistent collection libraries (e.g., those used in Rascal, Clojure, and Scala). This has led to a redesign of the PDB collection library that underlies the data structures of Rascal. Furthermore, he developed the Orpheus tool, an object redundancy profiler to assess the effects of maximal sharing.

## **DREAMPAL Team**

# **6. New Results**

## **6.1. Language-Independent Symbolic Execution, Program Equivalence, and Program Verification**

A significant part of our research project consists in applying formal techniques for symbolically executing and formally verifying HiHope programs, as well as for formally proving the equivalence of HiHope programs with the corresponding HoMade assembly and machine-code programs obtained by compilation of HiHope.

- Symbolic execution will detect bugs (e.g., stack underflow) in HiHope programs. Additionally, symbolic execution is the natural execution manner of HiHope programs as soon as they contain (typically, underspecified) hardware IPs;
- program verification will guarantee the absence of bugs (with respect to specified properties, e.g., no stack underflow, no invocation of unavailable IPs, ...);
- program equivalence will guarantee that such above-mentioned bugs are also absent from the HoMade assembly and machine-code programs obtained by compilation of HiHope source code.

Since these languages (especially HiHope) are not completely defined yet, we decided to work (together with our colleagues from Univ. Iasi, Romania) on language-independent symbolic execution, program-equivalence, and program-verification techniques. In this way, when all the languages in our project become stable, we will be readily able to instantiate the above generic techniques on (the K formal definitions of) the languages in question. We note that all the techniques described below are also independent of K: they are applicable to other language-definition frameworks that use similar rewriting-based formal operational semantics.

### **6.1.1. Symbolic Execution**

In [9] we propose a language-independent symbolic execution framework for languages endowed with a formal operational semantics based on term rewriting. Starting from a given definition of a language, a new language definition is automatically generated, which has the same syntax as the original one but whose semantics extends data domains with symbolic values and adapts semantical rules to deal with these values. Then, the symbolic execution of concrete programs is, by definition, the execution of programs with the new symbolic semantics, on symbolic input data. We prove that the symbolic execution thus defined has the properties naturally expected from it. A prototype implementation of our approach was developed in the K framework. We demonstrate the genericity of our tool by instantiating it on several languages, and show how it can be used for the symbolic execution and model checking of several programs.

### **6.1.2. Program Equivalence**

In [12] we propose a logic and a deductive system for stating and automatically proving the equivalence of programs in deterministic languages having a rewriting-based operational semantics. The deductive system is circular in nature and is proved sound and weakly complete; together, these results say that, when it terminates, our system correctly solves the program-equivalence problem as we state it. We show that our approach is suitable for proving the equivalence of both terminating and non-terminating programs, and also the equivalence of both concrete and symbolic programs. The latter are programs in which some statements or expressions are symbolic variables. By proving the equivalence between symbolic programs, one proves in one shot the equivalence of (possibly, infinitely) many concrete programs obtained by replacing the variables by concrete statements or expressions. We also report on a prototype implementation of the proposed deductive system in the K framework.

### 6.1.3. Program Verification

In [14] we present an automatic and language-independent program verification approach based on symbolic execution. The specification formalism we consider is Reachability Logic, a language-independent logic that constitutes an alternative to Hoare logics. Reachability Logic has a sound and relatively complete deduction system, which offers a lot of freedom (but no guidelines) for constructing proofs. Hence, we propose symbolic execution as a strategy for proof construction. We show that, under reasonable conditions on the semantics of programming languages, our symbolic-execution based Reachability-Logic formula verification is sound. We present a prototype implementation of the resulting language-independent verifier as an extension of a generic symbolic execution engine that we are developing in the K framework. The verifier is illustrated on programs written in languages also formally defined in K.

## 6.2. Master-Slave Control Structure for MP-SoC Architectures

Our Synchronous Communication Asynchronous Computation (SCAC) model is a data-parallel execution model dedicated to the Massively Parallel System-on-Chip. This model proposes a novel control structure, referred to as master-slave control [11]. Its concept departs from the centralized configuration. However, instead of a uni-processor master controlling a set of parallel processing elements (PE), the master cooperates with a grid of parallel slave controllers which supervises the activities of cluster of PEs.

The control structure in SCAC model is presented by two hierarchical control levels:

- The Master Control Unit (MCU), which controls the order execution in the whole system. It is a simple processor, which fetches and decodes program instruction and broadcasts execution orders to Slave Control Unit. It controls the end execution to establish synchronous communication.
- The Slave Control Unit (SCU), which controls: local node and PEs activities, parallel instructions execution and synchronous communication. It is a crucial component in the master-slave control structure. The SCUs grid allows independent parallel execution.

The hardware architecture is composed of a single MCU and multiple Slave controllers (SCUs) combined with local processing element (PE) (or a cluster of 16 PEs), known collectively as Nodes. The MCU and SCU array are connected through single level hierarchical bus and the SCUs are connected together through X-net interconnection network [2]. This network is clocked synchronously with the SCUs and respectively with the PEs. SCU controllers in the grid care for the instruction execution activities that involve a large degree of parallelism and the communication activities that need to coordinate all the PEs in the grid. The structure of master-slave control should be distinguished from other hierarchical or clustered approaches proposed for parallel computing. Such proposals are usually motivated by memory latency considerations and the desire to build a scalable system. The use of two control levels is therefore visible to the user in its effect on the communication between various processors. With master-slave control structure, the PEs in massively parallel system can execute independently and then can communicate synchronously. Such a construction has the advantage of allowing the designer to optimize distinct processors for their intended tasks and to implement simple interconnection network without additionally buffers and complex routing algorithms.

The aim of these last works is to design a master-slaves control structure for SCAC architecture to allow autonomous processing with simple and regular communication. This control structure based on IP blocks which offers good flexibility and scalability was implemented in synthesizable VHDL code. It is simulated and synthesized for Xilinx Virtex q6 (XC6VLX240T) board. The difficulty of designing a master-slave structure is a compromise between an optimal execution time and high flexibility, while reducing power consumption and silicon area.

## 6.3. Toward a Massively Parallel and Reflective Execution Model

FPGAs are undoubtedly suited to the definition of what could be called a DSHA (Domain Specific Hardware Architecture). Similarity with the DSSA (Domain Specific Software Architecture) an assembly of functional components performs basic transformations on data, while a software / hardware infrastructure ensures the

ordering of these transformations. The HoMade processor is designed with this in mind: it can be seen as an IP integrator offering a mechanism for interprocess communication IPs via a battery and a scheduler of IPs via dedicated instructions for flow control. In this control we find two particular instructions for flow control designed for a massively parallel execution model for SPMD, and a new instruction can make HoMade reflexive. With this instruction, you can at runtime change the behavior of a virtual component by dynamically associating it to a particular HoMade instruction sentence and in particular IP triggering instructions. Some components can successively after applying this instruction, trigger a hardware IP, a software function which itself can trigger a flow of execution of hardware IPs. This intercession<sup>2</sup> feature, parts of HoMade core, is valid for one processor or for all HoMade slave components in a massively parallel architecture. We demonstrated on a FPGA board which computes the Fibonacci sequence with three different methods, but always through a single call to a unique Virtual Component.

#### **6.4. Power Estimation at System-Level for MPSoC Based Platforms**

Shifting the design entry point up to the system level is the most important countermeasure adopted to manage the increasing complexity of Multiprocessor System on Chip (MPSoC). The reason is that decisions taken at this level, early in the design cycle, have the greatest impact on the final design in terms of power and energy efficiency. However, taking decisions at this level is very difficult, since the design space is extremely wide and it has so far been mostly a manual activity. Efficient system-level power estimation tools are therefore necessary to enable proper Design Space Exploration (DSE) based on power/energy and timing. We propose a tool based on efficient hybrid system level power estimation methodology for MPSoC. In this methodology, a combination of Functional Level Power Analysis (FLPA) and system level simulation technique are used to compute the power of the whole system. Basically, the FLPA concept is proposed for processor architecture in order to obtain parameterized arithmetic power models depending on the consumption of the main functional blocks. In this work, FLPA is extended to set up generic power models for the different parts of the platform. In addition, a simulation framework is developed at the transactional level to evaluate accurately the activities used in the related power models. The combination of the above two parts leads to a hybrid power estimation, that gives a better trade-off between accuracy and speed. The proposed methodology has several benefits: It considers the power consumption of the embedded system in its entirety; and Leads to accurate estimates without a costly and complex material. The proposed methodology is also scalable for exploring complex embedded architectures. Based on the proposed methodology, our Power Estimation Tool at System-Level (PETS) is developed. The usefulness and effectiveness of our PETS tool is validated through a typical mono-processor and multiprocessor embedded system designed around the TI OMAP (3530 and 5912) and the Xilinx Virtex II Pro FPGA boards. This methodology is demonstrated and evaluated by using a variety of basic programs to complete media benchmarks. Estimated power values are compared to real board measurements for both simple and multiprocessor architectures. Our obtained power estimation results provide less than 3% of error for mono-processor, 3.8% for homogeneous multiprocessor system and 4.3% for heterogeneous multiprocessor system and 70x faster compared to the state-of-the-art power estimation tools. These results have been presented in the PhD of Santhosh Kumar Rethinagiri [2] and published in [4].

#### **6.5. Dynamically reconfigurable CPU/FPGA architecture for the testing and simulation of avionic systems**

Real-time computing systems are increasingly used in aerospace and avionic industries. In the face of power wall and real-time requirements, hardware designers are directed towards reconfigurable computing with the usage of heterogeneous CPU/FPGA systems. However, there is a lack of real-time environments able to deal with the execution of applications on such heterogeneous systems dedicated to avionic Testing and Simulation (T&S). This year, we addressed the problem of soft real-time environments for CPU/FPGA systems and we proposed first a high-performance hardware architecture used to implement intimately coupled hardware and software avionic models. Second, we developed an efficient real-time software environment for the model's

---

<sup>2</sup>Wikipedia definition: intercession is the ability of a program to modify its own execution state or alter its own interpretation or meaning.



execution, the multi-core CPU monitoring and the runtime task re-allocation to avoid the timing constraint violation. Experimental results underpin the industrial relevance of the presented approach for avionic T&S systems with real-time support. These results are presented in the PhD of George Afonso [1] and in different publications [7] [10] [8].

## **6.6. A custom reconfiguration controller for partial and dynamic reconfiguration in HoMade based systems**

In all Xilinx devices supporting dynamic reconfiguration, such a functionality is realized using a hardware reconfiguration port called ICAP, that moves bitstreams from the reconfiguration memory to the programmable logic. ICAP is initialized by a Xilinx HW controller driven exclusively by a Microblaze processor and thus connected to a PLB or AXI bus.

This makes the partial and dynamic reconfiguration a very tedious task, as it implies using several Xilinx tools (XPS, ISE, PlanAhead,..etc). PDR becomes also resources and time consuming due to the fact that it uses very large interfaces and a static Xilinx architecture (in addition to the system that we want to design) including specific processors, buses, controllers,..etc.

Our contribution is the design of a custom ICAP controller, driven only by a HoMade processor, without any additional processors, buses or controllers. This ensures that our HoMade reconfigurable systems consumes fewer resources on the FPGA and does not require other tools than the standard ISE and PlanAhead tools in order to be designed.

## **6.7. Hardware control for partially and dynamically reconfigurable systems: from modelling to implementation**

This work proposes a control design methodology for FPGA-based reconfigurable systems aiming at increasing control design productivity and guaranteeing implementation efficiency. This methodology is based on a semi-distributed control model [5] composed of a set of modular distributed controllers executing each observation, decision-making and reconfiguration tasks for a reconfigurable region of the system, and a coordinator between the distributed controllers decisions in order to respect global systems constraints and objectives. This semi-distributed decision-making is based on the mode-automata formalism. The proposed combination between modularity, control splitting and formalism-based design allows to enhance the flexibility, reusability and scalability of the control design. Another point that can be added to this combination, to enhance design productivity, is design automation. For this, the proposed methodology is based on Model-Driven Engineering approach [5] allowing to automate code generation from high-level models. This approach makes use of the UML MARTE (Modeling and Analysis of Real-Time and Embedded Systems) standard profile, allowing to make low-level technical details transparent to designers and to automate the VHDL code generation for hardware implementation of the modeled control systems in order to guarantee their performance. The generated control systems were validated using simulation. Synthesis results showed an acceptable time and resource overhead for systems having different numbers of controllers. A control system composed of four controllers and a coordinator was also validated through physical implementation in an FPGA system for an image processing application.

## **6.8. A model-based approach for dynamically reconfigurable systems design: from MARTE to RecoMARTE**

This work is done in the context of the ANR FAMOUS project. It proposes a co-design methodology of dynamically reconfigurable systems based on FPGA. Our methodology is based on the Engineering Model Driven approach (MDE) and the models specification is done in the UML MARTE profile. It aims at ensuring flexibility, reusability and automation to facilitate the work of the designer and improve his productivity. The first contribution related is identifying parts of dynamically reconfigurable FPGA that can be modeled at the high abstraction levels. So, we defined a design flow based on the MDE to ensure the automation of code generation. According to this flow, several models are created mainly through MARTE profile concepts.

However, the modeling concepts of dynamic reconfiguration on FPGAs required extensions in MARTE. Thus, we identified the missing concepts to be integrated in a new profile that extends MARTE called RECOMARTE. The second contribution allows the chain automation and experimental validation. To integrate our design flow and to automate code generation, a processing chain was used. The final model resulting from MARTE proposed design flow is given as input to this chain.

We thereby move from MARTE to RECOMARTE models via an intermediate description according to the IP-XACT standard to finally generate files describing the complete system in the Xilinx XPS environment. This automation will accelerate the design phase and avoid errors due to the direct manipulation of these details. Finally, an example of application of image processing has been developed to demonstrate and validate our methodology.

## DOLPHIN Project-Team

## 6. New Results

### 6.1. Bi-level multi-objective optimization for pricing problems in long-haul transportation

Participants: M. Diaby, L. Brotcorne and E-G. Talbi

This work is concerned with the problem of pricing for a long-haul full load goods transportation. More precisely, we are interested in the situation where each vehicle delivers single request at a time. In this environment, we study the problem of pricing and valorization of unutilized capacity between two carriers. The first carrier B, cannot serve all the transportation requests and he thus needs to use outsourcing : second carrier A or his competitors. Carrier A, has to define the prices for carrier B transportation requests. Once carrier A has given its prices for the operations, it is B's decision to turn to A or to another carrier. This sequential and non-cooperative decision-making process can be adequately represented as a bilevel program : carrier B (the follower) wants to minimize transportation cost while A (the leader) seeks to maximize the revenue. Carrier A explicitly incorporates the reaction of carrier B in his optimization process.

Two types of models have been proposed : the bilevel mono-objective model and the bilevel biobjective model. More precisely, two objectives are simultaneously considered for the leader problem : the maximization of revenue and balancing the free load length (limiting the free load distances). We propose exact methods to solve moderate size instance of the problem and the heuristics to solve large-scale instances in reasonable time.

### 6.2. Approximating multi-objective scheduling problems

Participant: El-ghazali Talbi

External participants: Said Dabia, Tom Van Woensel, Tom De Kok (Eindhoven Technical University)

In this contribution, we propose a generic approach to deal with multi-objective scheduling problems (MOSPs). The aim is to determine the set of Pareto solutions that represent the interactions between the different objectives. Due to the complexity of MOSPs, an efficient approximation based on dynamic programming is developed. The approximation has a provable worst case performance guarantee. Eventhough the approximate Pareto set consists of fewer solutions, it represents a good coverage of the true set of Pareto solutions. We consider generic cost parameters that depend on the state of the system. Numerical results are presented for the time-dependent multi-objective knapsack problem, showing the value of the approximation in the special case when the state of the system is expressed in terms of time [23].

### 6.3. Force-Based Cooperative Search Directions in Evolutionary Multi-objective Optimization

Participants: Bilel Derbel, Dimo Brockhoff, Arnaud Liefvooghe

In order to approximate the set of Pareto optimal solutions, several evolutionary multi-objective optimization (EMO) algorithms transfer the multi-objective problem into several independent single-objective ones by means of scalarizing functions. The choice of the scalarizing functions' underlying search directions, however, is typically problem-dependent and therefore difficult if no information about the problem characteristics are known before the search process. In [46], we present new ideas of how these search directions can be computed *adaptively* during the search process in a *cooperative* manner. Based on the idea of Newton's law of universal gravitation, solutions attract and repel each other *in the objective space*. Several force-based EMO algorithms are proposed and compared experimentally on general bi-objective  $\rho$ MNK landscapes with different objective correlations. It turns out that the new approach is easy to implement, fast, and competitive with respect to a  $(\mu + \lambda)$ -SMS-EMOA variant, in particular if the objectives show strong positive or negative correlations.

## 6.4. DYNAMO (DYNAMIC programming using Metaheuristic for Optimization Problems)

Participants: Sophie Jacquin, Laetitia Jourdan, El-Ghazali Talbi

DYNAMOP (DYNAMIC programming using Metaheuristic for Optimization Problems) is a new dynamic programming based on genetic algorithm to solve a hydro-scheduling problem. The representation which is based on a path in the graph of states of dynamic programming is adapted to dynamic structure of the problem and it allows to hybridize easily evolutionary algorithms with dynamic programming. DYNAMOP has been tested on two case studies of hydro-scheduling problem with different price scenarios. Experiments indicate that the proposed approach performs considerably better than classical genetic algorithms and dynamic programming.

## 6.5. MOCA-I: Multi-Objective Classification Algorithm for Imbalanced Data

Participants: Julie Jacques, Clarisse Dhaenens, Laetitia Jourdan

Dealing with Imbalanced data is a real challenge as predicting the minority class may be very difficult but has a great interest for medical applications for example. Therefore, we propose MOCA-I, a new multi-objective local search algorithm that is conceived to deal with class imbalance, double meaning of missing data, volumetry and need of highly interpretable results all together [50]. MOCAI is a Pittsburgh multi-objective partial classification rule mining algorithm, using dominance-based multi-objective local search (DMLS). In comparison to state-of-the-art classification algorithms, MOCA-I obtains the best results on the 10 data sets of literature and is statistically better on the real data sets [50].

## 6.6. Neutrality Analysis is Graph coloring problem

Participants: Aymeric Blot, Clarisse Dhaenens, Laetitia Jourdan, Marie-Eleonore Marmion

Solving neutral problems is challenging as many optimization methods have difficulty to obtain good solutions. Hence, studying the neutrality in order to provide insights on the structure of the problem to be solved may be an answer. This has been done for the graph coloring problem (GCP) for which the neutrality of some hard instances has been quantified. This neutrality property has to be detected as it impacts the search process. Indeed, local optima may belong to plateaus that represent a barrier for local search methods. Then, we also aim at pointing out the interest of exploiting neutrality during the search. Therefore, a generic local search dedicated to neutral problems, NILS, is performed on several hard instances [78].

## 6.7. Neutrality in Multi-objective Local Search

Participants: Aymeric Blot, Clarisse Dhaenens, Laetitia Jourdan

External Participants: Hernan Aguirre, Kiyoshi Tanaka - Shinshu University, Japan

In multi-objective combinatorial optimization, the dominance-based local search algorithms are faced to sets of non-comparable solutions. In the absence of preferences, these solutions are equally good from the Pareto dominance perspective and can be considered neutral in term of quality, similar to the solutions who shares the same fitness value in mono-objective optimization. We propose two ideas to use the neutrality to improve the current local search algorithms. First, we analyze the distribution of neighbors for both small fully enumerable instances and hard large instances, to understand the distribution of neutral neighbors according to the rank of the solutions. Then, we compare the results of the proposed algorithms with the standard ones according to different indicators.

## 6.8. Biclustering for GWA data

Participants: Khedidja Seridi, Laetitia Jourdan, El-Ghazali Talbi

We have examined the possibilities of applying biclustering approaches to detect association between SNP markers and phenotype traits. Therefore, we have proposed a multiobjective model for biclustering problems in GWA context. Furthermore, we have proposed an adapted heuristic and meta-heuristic to solve it. The good performances of our algorithms are assessed using synthetic data sets.

## 6.9. Fitness Landscape Analysis for Multiobjective Optimization

Participant: Arnaud Liefoghe

External participants: Hernan Aguirre, Kiyoshi Tanaka (Shinshu Univ., Japan), Sébastien Verel (Univ. Littoral Côte d'Opale, France)

In [57], we investigate the correlation between the characteristics extracted from the problem instance and the performance of a simple evolutionary multiobjective optimization algorithm. First, a number of features are identified and measured on a large set of enumerable multiobjective NK-landscapes with objective correlation. A correlation analysis is conducted between those attributes, including low-level features extracted from the problem input data as well as high-level features extracted from the Pareto set, the Pareto graph and the fitness landscape. Second, we experimentally analyze the (estimated) running time of the global SEMO algorithm to identify a  $(1 + \epsilon)$ -approximation of the Pareto set. By putting this performance measure in relation with problem instance features, we are able to explain the difficulties encountered by the algorithm with respect to the main instance characteristics.

In [38], we study the effects of population size on selection and performance scalability of two dominance-based algorithms applied to many-objective optimization. Our aim is to understand the relationship between the size of the Pareto optimal set, a characteristic of the many-objective problem at hand, the population size and the ability of the algorithm to retain Pareto optimal solutions in its population and find new ones. This work clarifies important issues of the dynamics of evolutionary algorithms on many-objective landscapes, particularly related to survival selection. It shows that optimal solutions are dropped from the population in favor of suboptimal solutions that appear non-dominated when survival selection is applied. It also shows that this selection lapse, the dropping of optimal solution, affects the discovery of new optimal solutions and is correlated to population size and the distribution of solutions that survival selection renders. Selection makes less mistakes with larger populations and when the distribution of solutions is better controlled. The results of this study will be helpful to properly set population size and have a clearer idea about the performance expectation of the algorithm.

## 6.10. On Set-based Local Search for Multiobjective Combinatorial Optimization

Participant from DOLPHIN: Arnaud Liefoghe

External participants: Matthieu Basseur, Adrien Goëffon (Univ. Angers, France), Sébastien Verel (Univ. Littoral Côte d'Opale, France)

In [42], we formalize a multiobjective local search paradigm by combining set-based multiobjective optimization and neighborhood-based search principles. Approximating the Pareto set of a multiobjective optimization problem has been recently defined as a set problem, in which the search space is made of all feasible solution-sets. We here introduce a general set-based local search algorithm, explicitly based on a set-domain search space, evaluation function, and neighborhood relation. Different classes of set-domain neighborhood structures are proposed, each one leading to a different set-based local search variant. The corresponding methodology generalizes and unifies a large number of existing approaches for multiobjective optimization. Preliminary experiments on multiobjective NK-landscapes with objective correlation validates the ability of the set-based local search principles. Moreover, our investigations shed the light to further research on the efficient exploration of large-size set-domain neighborhood structures.

## 6.11. Feature selection in high dimensional regression problems for genomics

Participants: Julie Hamon, Clarisse Dhaenens (External collaborator : Julien Jacques)

In the context of genomic selection in animal breeding, an important objective consists in looking for explicative markers for a phe-notype under study. In order to deal with a high number of markers, we propose to use combinatorial optimization to perform variable selection. Results show that our approach outperforms some classical and widely used methods on simulated and “closed to real” datasets [76]. Familial relationships have also been used in this specific context and allow to improve results.

## 6.12. Indicator-Based Multiobjective Optimization

Participant: Dimo Brockhoff

External Participants: Johannes Bader (formerly at ETH Zurich, Switzerland), Youssef Hamadi (Microsoft Research, Cambridge, UK), Souhila Kaci (Université Montpellier 2, France), Lothar Thiele (ETH Zurich, Switzerland), Heike Trautmann (University of Munster, Germany) Tobias Wagner (TU Dortmund, Germany), and Eckart Zitzler (PH Bern, Switzerland)

Indicator-based (evolutionary) multiobjective optimization algorithms have been first introduced in 2004 and typically use a quality indicator, assigning a solution set a real value, as a direct, internal performance criterion. Given that the indicator and the number  $\mu$  of desired points is fixed, the optimization goal, also denoted by the term *optimal  $\mu$ -distribution*, is then defined as the solution set(s) of size  $\mu$  which maximizes the indicator value.

In 2013, we continued to investigate, theoretically and numerically, the optimal  $\mu$ -distributions for the R2 indicator, an often recommended indicator based on scalarization functions [73]. We also proposed a new multiobjective optimizer with an R2-indicator-based selection [70]. With respect to the even more common *hypervolume indicator*, we combined the idea of the weighted hypervolume indicator with the idea of interactive algorithms and proposed a new algorithm that adapts the weighted hypervolume’s weight function according to the user’s preferences during the search. Last, we summarized our knowledge on the weighted hypervolume indicator and proposed a general framework of how to employ it within the hypervolume-based W-HypE algorithm [18].

## 6.13. A Hybrid Metaheuristic for Multiobjective Unconstrained Binary Quadratic Programming

Participant : Arnaud Liefoghe

External participants : Jin-Kao Hao (Univ. Angers, France), Sébastien Verel (Univ. Littoral Côte d’Opale, France)

The conventional Unconstrained Binary Quadratic Programming (UBQP) problem is known to be a unified modeling and solution framework for many combinatorial optimization problems. In [29], we extend the single-objective UBQP to the multiobjective case (mUBQP) where multiple objectives are to be optimized simultaneously. We propose a hybrid metaheuristic which combines an elitist evolutionary multiobjective optimization algorithm and a state-of-the-art single-objective tabu search procedure by using an achievement scalarizing function. Finally, we define a formal model to generate mUBQP instances and validate the performance of the proposed approach in obtaining competitive results on large-size mUBQP instances with two and three objectives.

## 6.14. Multi-core GPU-based parallel optimization

We have mainly investigated the design and implementation on multi-core GPU-based platforms of metaheuristics and tree-based exact optimization methods focusing on Branch and bound (B&B) algorithms (Ph.D thesis of I. Chakroun).

- **GPU-based Metaheuristics**

Participants: N. Melab, T-V. Luong, K. Boufaras and N. Melab

We came out with a pioneering work on single-solution methods. The hierarchy of parallel models has been rethought on GPU dealing with CPU-GPU data transfer optimization, thread control and automatic mapping of candidate solutions to threads. The implementation of the proposed approaches is provided through ParadisEO-GPU in [62] (nominated for Best Paper Award). High speedups have been achieved for some problems.

- **Multi-core GPU-based B&B algorithms**

For exact optimization, we have revisited the design and implementation of highly irregular B&B algorithms on GPU dealing with hierarchical device memory optimization, on GPU combined with multi-core [45] dealing with CPU-GPU data transfer optimization and work partitioning, and on GPU-enhanced computational grids. Accelerations up to  $\times 217$  are achieved on Tesla Nvidia C2050 on large Flow-Shop problems.

## 6.15. Energy-aware scheduling for clouds

Participants: Y. Kessaci, N. Melab, E-G. Talbi

High-performance computing (HPC) is moving from in-house to cloud-based HPC. One of the major issues of this later is the scheduling of HPC applications taking into account the energy criterion in addition to performance. In [54], we have addressed that issue (Ph.D thesis of Y. Kessaci). We have proposed several metaheuristics for cloud managers and experimented on OpenNebula using different (VMs) arrival scenarii and different hardware infrastructures. The results show that our approaches outperform the scheduler provided in OpenNebula by a significant margin in terms of energy consumption and number of scheduled VMs.

## 6.16. Heterogenous Multi-CPU Multi-GPU Parallel Branch-and-Bound Tree Search

Participants: Trong-Tuan Vu, Bilel Derbel, Nouredine Melab

In this work [71], we push forward the design of parallel and distributed optimization algorithms running on heterogenous systems consisting of multiple CPUs coupled with multiple GPUs. We consider parallel Branch-and-Bound (B&B), viewed as a generic algorithm searching in a dynamic tree representing a set of candidate solutions built dynamically at runtime. Given that several distributed CPUs and GPUs coming from possibly different clusters connected through a network can be used to parallelize the tree search, we give new insights into how to fully benefit from such a heterogeneous environment. More precisely, we describe a two-level generic and fully distributed parallel approach taking into account PU characteristics. In the first level, we use data streaming in order to allow parallelism between hosts and devices. The evaluation of tree nodes is done inside a GPU while the CPU-host is performing the pruning, selection and decomposition operations in parallel. In the second level, our approach incorporates an adaptive dynamic load balancing scheme based on distributed work stealing, in order to flow workloads efficiently from overloaded PUs to idle ones at runtime. We deployed our approach over a distributed system of up to 20 GPUs and 128 CPUs coming from three clusters. Different scales and configurations of PUs were experimented with the B&B algorithm and the well-known FlowShop combinatorial optimization problem as a case study. Firstly, on one single GPU, we improve on the running time of previous B&B GPUs implementation by at least a factor of two. More importantly, independently of CPUs or GPUs scale or power, our approach provides a substantial speed-up which is *nearly optimal* compared to the ideal performance one could expect in theory.

## MODAL Project-Team

### 6. New Results

#### 6.1. Resampling procedures

**Participant:** Alain Celisse.

The new deep understanding of cross-validation procedures in density estimation has been tackled with new results in terms of risk estimation and model selection [7]. This is the first step towards a fully data-driven and optimal choice of cross-validation strategy.

#### 6.2. Kernel change-point

**Participants:** Alain Celisse, Guillemette Marot, Morgane Pierre-Jean.

On the basis of theoretical arguments, an empirical analysis has been carried out to assess the influence of the choice of the kernel in the kernel change-point strategy described in [2]. This assessment has been done in the biological context of copy number variation and allele B fraction. Several talks have been given in seminars (SSB seminar in Paris,...) and workshops (JSFDS, SMPGD,...)

#### 6.3. Gaussian process in RKHS

**Participants:** Alain Celisse, Jérémie Kellner.

Since numerous papers make a Gaussian assumption for observations in the reproducing kernel Hilbert space (RKHS), it is important to be able to assess the validity of this crucial assumption. As long as it has been validated, the Gaussian framework can be further used to infer statistical properties of the population at hand (mean, variance,...).

A statistical test has been designed to address such questions at the RKHS level. It is fully computationally efficient and provides really good power in numerous settings. Theoretical properties for the test statistic have been derived as well.

#### 6.4. Model for conditionally correlated categorical data

**Participants:** Christophe Biernacki, Matthieu Marbac-Lourdelle, Vincent Vandewalle.

It is a model-based clustering where categorical data are grouped into conditionally independent blocks. The corresponding block distribution is a parsimonious multinomial distribution where the few free parameters correspond to the most likely modality crossings, while the remaining probability mass is uniformly spread over the other modality crossings. The exact computation of the integrated complete-data likelihood allows to perform the model selection, by a Gibbs sampler, reducing the computing time consuming by parameter estimation and avoiding BIC criterion biases pointed out by our experiments.

This model was presented in a conference [13] with scientific committee and in a seminar [17]. An article will be soon submitted. Furthermore, a R package is currently under development.

#### 6.5. Mixture model for mixed kind of data

**Participants:** Christophe Biernacki, Matthieu Marbac-Lourdelle, Vincent Vandewalle.

A mixture model of Gaussian copula allows to cluster mixed kind of data. Each component is composed by classical margins while the conditional dependencies between the variables is modeled by a Gaussian copula. The parameter estimation is performed by a Gibbs sampler. This model was presented in a conference [14]. Some technical points will be developed before providing an article.



## 6.6. Mixture of Gaussians with Missing Data

**Participants:** Christophe Biernacki, Vincent Vandewalle.

The generative models allow to handle missing data. This can be easily performed by using the EM algorithm, which has a closed form M-step in the Gaussian setting. This can for instance be useful for distance estimation with missing data. It has been proposed to improve the distance estimation by fitting a mixture of Gaussian distributions instead of a considering only one Gaussian component [21]. This is a joined work with Emil Eirola and Amaury Lendrasse .

A parallel work is in progress on the mixture degeneracy when considering mixture of Gaussians with missing data. It have been experimentally noticed that the degeneracy in this case is particularly slow. This behaviour is different from the usual setting of degeneracy with mixture of Gaussians which is usually rather fast. A first attempt of the theoretical characterization of this behaviour around a degenerated solution has been presented at a conference [16].

## 6.7. Transfert learning in model-based clustering

**Participant:** Christophe Biernacki.

In many situations one needs to cluster several datasets, possibly arising from different populations, instead of a single one, into partitions with identical meaning and described by similar features. Such situations involve commonly two kinds of standard clustering processes. The samples are clustered traditionally either as if all units arose from the same distribution, or on the contrary as if the samples came from distinct and unrelated populations. But a third situation should be considered: As the datasets share statistical units of same nature and as they are described by features of same meaning, there may exist some link between the samples. We propose a linear stochastic link between the samples, what can be justified from some simple but realistic assumptions, both in the Gaussian and in the  $t$  mixture model-based clustering context [26]. This is a joint work with Alexandre Lourme.

## 6.8. Gaussian Models Scale Invariant and Stable by Projection

**Participant:** Christophe Biernacki.

Gaussian mixture model-based clustering is now a standard tool to determine an hypothetical underlying structure into continuous data. However many usual parsimonious models, despite their appealing geometrical interpretation, suffer from major drawbacks as scale dependence or unsustainability of the constraints by projection. In this work we present a new family of parsimonious Gaussian models based on a variance-correlation decomposition of the covariance matrices. These new models are stable by projection into the canonical planes and, so, faithfully representable in low dimension. They are also stable by modification of the measurement units of the data and such a modification does not change the model selection based on likelihood criteria. We highlight all these stability properties by a specific geometrical representation of each model. A detailed GEM algorithm is also provided for every model inference. Then, on biological and geological data, we compare our stable models to standard geometrical ones.

This joint work with Alexandre Lourme is now published in [6].

## 6.9. Clustering and variable selection in regression

**Participants:** Christophe Biernacki, Loïc Yengo, Julien Jacques.

A new framework is proposed to address the issue of simultaneous linear regression and clustering of predictors where regression coefficients are assumed to be drawn from a Gaussian mixture distribution. Prediction is thus performed using the conditional distribution of the regression coefficients given the data, while clusters are easily derived from posterior distribution in groups given the data. This work is now published in [28]

## 6.10. An AIC-like criterion for semi-supervised classification

**Participants:** Christophe Biernacki, Vincent Vandewalle.

In semi-supervised classification, generative models take naturally into account unlabeled data and parameter estimation can be easily performed through the EM algorithm. However, traditional model selection criteria either does not take into consideration the predictive purpose (AIC or BIC criteria) or involve a high computational cost because of the EM mechanism (cross validation criteria). Alternatively, we propose the penalized model selection criterion AICcond which aims to estimate the predictive power of a generative model by approximating its predictive deviance. AICcond has similar computational cost to AIC, owns good consistency theoretical properties and highlights encouraging behaviour for variable and model selection in comparison to other standard criteria.

This joint work with Gilles Celeux and Gérard Govaert is now published in[16].

## 6.11. Consistency of a nonparametric conditional mode estimator for random fields

**Participant:** Sophie Dabo-Niang.

Sophie Dabo-Niang settled the consistency of a nonparametric conditional mode estimator for random fields, Statistical Methods and Applications [9].

## 6.12. Spatial linear models

Spatial linear models only capture global linear relationships between locations. However, in many circumstances the spatial dependency is not linear. It is, for example, the classical case where one deals with the spatial pattern of extreme events such as in the economic analysis of poverty, in the environmental science,... This leads naturally to consider nonparametric modeling.

## 6.13. Auto-associative models

Serge Iovleff gave a complete treatment of the Auto-Associative models in the semi-linear case and wrote a software for estimating these models (hal-00734070, version 1).

## 6.14. BlockCluster

Serge Iovleff has submitted a paper on the BlockCluster package in collaboration with Parmeet Bathia.

## 6.15. Rmixmod

Serge Iovleff has contributed to a paper submitted to JSS (hal-00919486, version 1) in collaboration with R. Lebre, F. Langrognet, C. Biernacki, G. Celeux, and G. Govaert.

## 6.16. Clustering for functional data

**Participants:** Julien Jacques, Cristian Preda.

In Jacques & Preda 2014 (CSDA), we propose a model-based clustering algorithm for multivariate functional data, based on multivariate functional principal components analysis. A review on clustering for functional data has also be published in Jacques & Preda 2014 (ADAC). Variable selection in high-dimensional regression  
Participants: Julie Hamon, Julien Jacques, Clarisse Dhaenens. In the context of genomic analysis, dealing with high-throughput genotyping data, we develop a genetic algorithm which looks for the best subset of variables (of given size) to predict some quantitative feature.

## 6.17. Wavelet based clustering using mixed effects functional models

**Participant:** Guillemette Marot.

The paper related to the wavelet based clustering procedure presented in the activity report from MODAL team in 2012 was published in Biometrics [22].

### **6.18. Differential meta-analysis of RNA-seq data from multiple studies**

**Participant:** Guillemette Marot.

An adaptation of meta-analysis methods initially proposed for microarray studies has been proposed for RNA-seq data. The R package metaRNASeq is available on the R Forge and the preprint of the paper is available on Arxiv [48].

### **6.19. Toxoplasma transcription factor TgAP2XI-5 regulates the expression of genes involved in parasite virulence and host invasion**

**Participant:** Guillemette Marot.

The use of peak detection methods implemented in the Bioconductor package Ringo has enabled to better understand part of the gene regulation process in *T. Gondii* parasite. The new findings in Biology have been published in *Walker (2013)*.

## NON-A Project-Team

### 6. New Results

#### 6.1. Homogeneity theory and analysis of nonlinear systems

Homogeneity is a kind of symmetry, if it is presented in a system model, then it may simplify analysis of stability and performance properties of the system. The new results obtained in 2013 are as follows:

- The notion of geometric homogeneity has been extended for differential inclusions in [44]. This kind of homogeneity provides the most advanced coordinate-free framework for analysis and synthesis of nonlinear discontinuous systems. Theorem of L. Rosier on a homogeneous Lyapunov function existence and an equivalent notion of global asymptotic stability for differential inclusions have been presented. Robustness properties (ISS) of sliding mode systems applying the homogeneity concept have been considered in [46].
- Retraction obstruction for time-varying stabilization on compact manifolds has been revisited in [13].
- Several conditions have been proposed to check different robustness properties (ISS, iISS, IOSS and OSS) for generic nonlinear systems applying the weighted homogeneity concept (global or local) in [14], [45]. The advantages of this result are that, under some mild conditions, the system robustness can be established as a function of the degree of homogeneity.
- A new algorithm for the analysis of strange attractors has been presented in [51]. An application of that results for observability-singularity manifolds in the context of chaos based cryptography has been given in [52].
- Exciting multi-DOF systems by feedback resonance has been considered in [20].
- Some conditions on existence of oscillations in hybrid systems have been established in [23], [57]. An application to a humanoid robot locomotion has been considered.
- Considering two chaotic Rossler systems, the paper [83] presents a study on the forced synchronization of two systems, bidirectionally coupled by transmitting unidirectional signals which explicitly depend on a single state variable (from the emitter) and only affect directly the dynamics corresponding to the transmitted state variable (of the receiver).
- The paper [33] is concerned with the construction of local observers for nonlinear systems without inputs, satisfying an observability rank condition. The aim of this study is, first, to define a homogeneous approximation that keeps the observability property unchanged. This approximation is further used in the synthesis of local observer which is proven to be locally convergent for Lyapunov-stable systems.
- The paper [74] addresses the problem of exact average-consensus reaching in a prescribed time. The communication topology is assumed to be defined by a weighted undirected graph and the agents are represented by integrators. A nonlinear control protocol, which ensures a finite-time convergence, is proposed. With the designed protocol, any prescribed convergence time can be guaranteed regardless of the initial conditions.
- The Implicit Lyapunov Function (ILF) method for finite-time stability analysis has been introduced in [75]. The control algorithm for finite-time stabilization of a chain of integrators has been developed. The scheme of control parameters selection has been presented by LMIs. The robustness of the finite-time control algorithm with respect to system uncertainties and disturbances has been studied. The new high order sliding mode control has been derived as a particular case of the developed finite-time control algorithm. The settling time estimate has been obtained using ILF method. The algorithm of practical implementation of the ILF control scheme has been discussed.

## 6.2. Model-free control

The model free control techniques form a new and quickly developing area of control theory. It has been established by the team members and nowadays these tools find many practical applications and attract a lot of attention due to their clear advantages for designers: they provide a control law independently in the model knowledge. The achievements obtained in 2013 are as follows:

- A new development of the model-free control theory with application to active magnetic bearing control have been presented in [53].
- "Model-free control" and the corresponding "intelligent" PID controllers (iPIDs), which already had many successful concrete applications, have been presented in [27] for the first time in a unified manner, where the new advances have been taken into account.
- In [62], it is shown that the "intelligent" controllers, which are associated to the recently introduced model-free control synthesis, may be easily implemented on cheap and small programmable devices.
- An application of the model-free control for regulation of the water level under several constraints has been reported in [40].

## 6.3. Algebraic technique for estimation, differentiation and its applications

Elementary techniques from operational calculus, differential algebra, and non-commutative algebra lead to a new algebraic approach for estimation and detection. It is investigated in various areas of applied sciences and engineering. The following lists only some applications:

- Design of a stabilizing feedback based on acceleration measurements and an algebraic state estimation method has been proposed in [54].
- An extension of the algebraic differentiation method to fractional derivatives calculation in continuous and discrete time has been studied in [88] and [89] respectively. Applications to identification and parameter estimation of fractional linear systems have been considered in [67], [68].
- Smoothing noisy data with spline functions is well known in approximation theory. Smoothing splines have been already used to deal with the problem of numerical differentiation. In [43], we extend this method to estimate the fractional derivatives of a smooth signal from its discrete noisy data. We begin with finding a smoothing spline by solving the Tikhonov regularization problem. Then, we propose a fractional order differentiator by calculating the fractional derivative of the obtained smoothing spline.
- In [81], we apply an algebraic method to estimate the amplitudes, phases and frequencies of a biased and noisy sum of complex exponential sinusoidal signals. The obtained estimates are integrals of the noisy measured signal: these integrals act as time-varying filters.

## 6.4. Observability and observer design for nonlinear systems

Observability analysis and observer design are important issues in the field of control theory. Some recent results are listed below:

- An epistemology of observation theory and its application in the design of software sensor in power electronics have been presented in [42].
- New results on observability and detectability of singular linear systems with unknown inputs have been developed in [12].
- The paper [47] supplies a new algorithm to compute the internal dynamics (or inversion dynamics) of affine MIMO control nonlinear systems.
- The design of observers for nonlinear systems with unknown, time-varying, bounded delays, on both state and input, still constitutes an open problem. In [28], we show how to solve it for a class of nonlinear systems by combining the high gain observer approach with a Lyapunov-Krasovskii functional. Sufficient conditions have been provided to prove the practical stability of the observer.

- An influence of restricted isometry property to the observability under sparse measurements has been analyzed in [65].
- The paper [38], [79] concerns the design of a nonlinear observer through a transformation of a nonlinear system into an observer form that supports a high gain observer. Sufficient geometrical condition has been deduced to guarantee the existence of change of coordinates allowing the transformation of a nonlinear system into the proposed normal form. In [80], the Partial Observability Normal Forms (PONF) of nonlinear dynamical systems have been investigated. Necessary and sufficient conditions for the existence of a diffeomorphism bringing the original nonlinear system into a PONF have been established.

## 6.5. Sliding mode control and estimation

Sliding mode algorithms are very popular for finite-time estimation and regulation. The recent results obtained by the group are as follows:

- Some constructive approximations and an alternative theoretic characterization of some classes of sliding mode control processes has been presented in [11].
- In [64] we investigate observer design under sparse measurement, i.e. under Nyquist-Shannon frequency. An analysis demonstrates that it is impossible to use only a high order sliding mode observer in the case of sparse measurement. Then it has been shown that a high order sliding mode observer coupled with an impulsive observer is a pertinent solution at least for some particular class of systems.
- Anomaly detection has been an active open problem in the networks community for several years. In [35], we aim at detecting such abnormal signals by control theory techniques. Several classes of sliding mode observers have been proposed for a fluid flow model of the TCP/internet protocol network.
- A sliding mode control has been developed for robust stabilization of fractional-order input-delay linear systems in the presence of uncertainties and external disturbances in [78]. First, a fractional-order state predictor has been used to compensate the delay in the input control. Second, a robust sliding mode control has been proposed in order to stabilize the system and to thwart the effect of model uncertainties and external disturbances. The sliding mode controller has been designed by considering a sliding surface defined by fractional order integral.

## 6.6. Non-linear, Sampled and Time-delay systems

Nonlinearities, sampling, quantization and time-delays cause serious obstructions for control and observer design in many fields of techniques and engineering (e.g. networked and internet systems, distributed systems etc.). The proposed by the team algebraic approach suits well for estimation and regulation in such a type of systems. The recent results are listed below:

- The work [59] aims at decreasing the number of sampling instants in state feedback control for perturbed linear time invariant systems. The approach is based on linear matrix inequalities obtained thanks to Lyapunov-Razumikhin stability conditions and convexification arguments that guarantee the exponential stability for a chosen decay-rate.
- A novel self-triggered control, which aims at decreasing the number of sampling instants for the state feedback control of perturbed linear time invariant systems, has been proposed in [60]. The approach is based on convex embeddings that allow for designing a state-dependent sampling function guaranteeing the system's exponential stability for a desired decay-rate and norm-bounded perturbations. One of the main contributions of the paper [60] is an LMI based algorithm that optimizes the choice of the Lyapunov function so as to enlarge the lower-bound of the sampling function while taking into account both the perturbations and the decay-rate.
- In [63], we consider the issue of stabilizing a class of linear systems using irregular sampled output measurements.

- The paper [73] is dedicated to the stability analysis of nonlinear sampled-data systems, which are affine in the input. Assuming that a stabilizing continuous-time controller exists and it is implemented digitally, we intend to provide sufficient asymptotic/exponential stability conditions for the sampled-data system. This allows to find an estimate of the upper bound on the asynchronous sampling periods. The stability analysis problem is formulated both globally and locally. The main idea of the paper is to address the stability problem in the framework of dissipativity theory. Furthermore, the result is particularized for the class of polynomial input-affine sampled-data systems, where stability may be tested numerically using sum of squares decomposition and semidefinite programming.
- The problem of output control design for linear system with unknown and time-varying input delay, bounded exogenous disturbances and bounded deterministic measurement noises has been considered in [77]. The prediction technique has been combined with Luenberger-like observer design in order to provide the stabilizing output feedback. The scheme of parameters tuning for reduction of measurement noises effect and exogenous disturbances effects has been developed using the Attractive Ellipsoids Method.
- Using the theory of non-commutative rings, the paper [39] studies the delay identification of nonlinear time-delay systems with unknown inputs. A sufficient condition has been given in order to deduce an output delay equation from the studied system. Then necessary and sufficient conditions have been proposed to judge whether the deduced output delay equation can be used to identify the delay, which is involved in this equation.

## 6.7. Interval control and estimation

In many cases due to parametric and/or signal uncertainties presented in a plant model it is not possible to design a conventional observer, which provides a point-wise estimate of state in a finite time or asymptotically. In this case it is still frequently possible to design interval observers, which generate an estimate on the interval of the admissible values of the state at the current instant of time. The recent new results in this field are listed below:

- The work [49] is devoted to interval observer design for Linear Parameter-Varying (LPV) systems under assumption that the vector of scheduling parameters is not available for measurements. Stability conditions are expressed in terms of matrix inequalities, which can be solved using standard numerical solvers. Robustness and estimation accuracy with respect to model uncertainty is analyzed. Two solutions are proposed for nonnegative systems and for a generic case. The efficiency of the proposed approach is demonstrated through computer simulations.
- Development of interval observers for time invariant [55] and time-varying [21] discrete-time systems has been presented by the members of the team.
- Interval estimation for uncertain systems with time-varying delays has been considered in [22], [56]. A reduced-order interval observer has been designed, stability and robustness conditions have been obtained.
- The paper [24] is devoted to design of interval observers for Linear Time Varying (LTV) systems and a class of nonlinear time-varying systems in the output canonical form. An interval observer design is feasible if it is possible to calculate the observer gains making the estimation error dynamics cooperative and stable. It has been shown in [24] that under some mild conditions the cooperativity of an LTV system can be ensured by a static linear transformation of coordinates.
- The problem of output stabilization of a class of nonlinear systems subject to parametric and signal uncertainties has been studied in [25]. First, an interval observer has been designed estimating the set of admissible values for the state. Next, it has been proposed to design a control algorithm for the interval observer providing convergence of interval variables to zero, that implies a similar convergence of the state for the original nonlinear system. An application of the proposed technique shows that a robust stabilization can be performed for linear time-varying and LPV systems without assumption that the vector of scheduling parameters is available for measurements.

- The paper [26] deals with the problem of joint state and parameter estimation based on a set adaptive observer design. The problem is formulated and solved for an LPV system. The resolution methodology avoids the exponential complexity obstruction usually encountered in the set-membership parameter estimation.
- The output stabilization problem for a linear system with an unknown bounded time-varying input delay has been considered in [34], [76]. The interval observation technique has been applied in order to obtain guaranteed interval estimate of the system state. The procedure of the interval observer synthesis uses lower and upper estimates of the unknown delay and requires to solve a special Sylvester's equation. The interval predictor has been introduced in order to design a linear stabilizing feedback. The control design procedure is based on LMIs.
- The paper [37] describes a robust set-membership-based Fault Detection and Isolation (FDI) technique for a particular class of nonlinear systems, the so-called flat systems. The proposed strategy consists in checking if the expected input value belongs to an estimated feasible set computed using the system model and the derivatives of the measured output vector. The output derivatives are computed using a numerical differentiator. The set-membership estimator design for the input vector takes into account the measurement noise thereby making the consistency test robust.
- The objective of the work [82] is to develop some design methods of interval observers for a class of nonlinear continuous-time systems. It has been assumed that the estimated system can be represented as a superposition of the nominal subsystem (belonged to the class of uniformly observable systems) and a Lipschitz nonlinear perturbation vanishing at the origin. Then it has been shown that there exists an interval observer for the system that estimates the set of admissible values for the state consistent with the output measurements.

## 6.8. Networked robots

The mobile robots constitute an important area of practical development for the team:

- The paper [71] presents a path planning algorithm for autonomous navigation of non-holonomic mobile robots in complex environment. The irregular contour of obstacles is represented by segments. The goal of the robot is to move towards a known target while avoiding obstacles. The velocity constraints, kinematic robot model and non-holonomic constraint are considered in the problem. The optimal path planning problem is formulated as a constrained receding horizon planning problem and the trajectory is obtained by solving an optimal control problem with constraints. Local minima are avoided by choosing intermediate objectives based on the real time environment.
- The paper [69] presents a cooperative path planning approach for the navigation of non-holonomic mobile robots in environment with obstacles. Shared information can be obtained by sharing the local information between robots, thus the trajectories can be more optimized. Visibility graph approach is used to generate a series of intermediate objectives which guarantee the robots to reach the final objective without local minima. Then the reach of intermediate objectives is ensured by the optimal path planning algorithm. The velocity constraints, kinematic constraints and non-holonomic constraints of the mobile robot are considered in the problem.
- The paper [70] presents the real-time identification of different types of non-holonomic mobile robot systems. Since the robot type is a priori unknown, the robot systems are formulated as a switched singular nonlinear system, and the problem becomes the real-time identification of the switching signal, and then the existence of the input-output functions and the distinguishability of the systems are studied.
- An intelligent PID controller (*i*-PID controller) has been applied to control the non-holonomic mobile robot with measurement disturbance in [72]. Because of the particularity of the non-holonomic systems, this paper proposes to use a switching parameter  $\alpha$  in the *i*-PID controller.



## 6.9. Applications

As it was mentioned, Non-A is a kind of "method-driven" project, which deals with different aspects of finite-time estimation and control. Thus different applications are possible, ones touched this year are as follows (skipping the networked robots considered in the previous section):

- A sensorless speed control for a DC series motor has been presented in [41] based on sliding-mode control and estimation algorithms.
- The paper [48] presents a feasibility study, which aims to demonstrate the applicability of the CNC automation philosophy for the process of AFM probe-based nano machining conducted on commercial AFM instruments.
- An oscillatory failure case detection for aircrafts using non-homogeneous sliding-mode differentiator in noisy environment has been considered in [50].
- Sensorless fault tolerant control for induction motors has been developed in [18].
- The problem of an actuator fault detection in aircraft systems has been considered in [19]. A particular attention has been paid to the oscillatory failure case study.
- In [58], we consider a vehicle equipped with active front steer and rear torque vectoring. While the former adds an incremental steer angle to the driver's input, the latter imposes a torque by means of the rear axle. The active front steer control is actuated through the front tires, while the rear torque vectoring can be actuated through the rear tires. A nonlinear controller using the super-twisting algorithm has been designed in order to track in a finite time the lateral and yaw angular velocity references.
- Systematic and multifactor risk models have been revisited via algebraic methods, which were already successfully developed in signal processing and in automatic control, in [61].
- In [84], we address the problem of approximating scattered data points by C1-smooth polynomial spline curves and surfaces using L1-norm optimization. The use of this norm helped us to preserve the shape of the data even near to abrupt changes.
- As capacitor voltages are necessary for the three-cell DC-DC chopper control, the estimation of such voltages by an observer is attractive solution in terms of cost. However, due to the hybrid behaviour of this structure, the capacitor voltages may be partially or even not observable for a given switching configuration. In other words, the observability matrix associated to the capacitor voltages never has a full rank. In order to make the observer conceivable, the paper [29] proposes a new design by establishing sufficient conditions under which the capacitor voltages can be reconstructed within appropriate specific switching sequence and not necessarily instantly.
- The problem of converters coordination of a fuel cell system involving a hydrogen fuel cell with supercapacitors for applications with high instantaneous dynamic power has been addressed in [32]. The problem is solved by using a non-linear controller based on passivity.
- The paper [66] is devoted to development of control algorithms for nonlinear parametrically uncertain systems. Original system dynamics is approximated by a set of local NARX models combined by a special mixing rule. Algorithm for local models' parameters estimation and structure adjustment has been developed. The developed technique has been applied to the problem of regulation of spark ignition engines.
- The paper [36] is dedicated to the problem of pneumatic cylinder control without pressure measurement. Based on the theory of homogeneous, finite time stable, ordinary differential equations, a state feedback nonlinear controller has been proposed. The closed loop system stability has been proven and an attraction domain of the controller has been given.

## SequeL Project-Team

## 6. New Results

### 6.1. Decision-making Under Uncertainty

#### 6.1.1. Reinforcement Learning

##### *Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model [2]*

We consider the problem of learning the optimal action-value function in discounted-reward Markov decision processes (MDPs). We prove new PAC bounds on the sample-complexity of two well-known model-based reinforcement learning (RL) algorithms in the presence of a generative model of the MDP: value iteration and policy iteration. The first result indicates that for an MDP with  $N$  state-action pairs and the discount factor  $\gamma \in [0, 1)$  only  $O(N \log(N/\delta) / [(1 - \gamma)^3 \epsilon^2])$  state-transition samples are required to find an  $\epsilon$ -optimal estimation of the action-value function with the probability (w.p.)  $1 - \delta$ . Further, we prove that, for small values of  $\epsilon$ , an order of  $O(N \log(N/\delta) / [(1 - \gamma)^3 \epsilon^2])$  samples is required to find an  $\epsilon$ -optimal policy w.p.  $1 - \delta$ . We also prove a matching lower bound of  $\Omega(N \log(N/\delta) / [(1 - \gamma)^3 \epsilon^2])$  on the sample complexity of estimating the optimal action-value function. To the best of our knowledge, this is the first minimax result on the sample complexity of RL: The upper bound matches the lower bound in terms of  $N$ ,  $\epsilon$ ,  $\delta$  and  $1/(1 - \gamma)$  up to a constant factor. Also, both our lower bound and upper bound improve on the state-of-the-art in terms of their dependence on  $1/(1 - \gamma)$ .

##### *Regret Bounds for Reinforcement Learning with Policy Advice [13]*

In some reinforcement learning problems an agent may be provided with a set of input policies, perhaps learned from prior experience or provided by advisors. We present a reinforcement learning with policy advice (RLPA) algorithm which leverages this input set and learns to use the best policy in the set for the reinforcement learning task at hand. We prove that RLPA has a sub-linear regret of  $O(\sqrt{T})$  relative to the best input policy, and that both this regret and its computational complexity are independent of the size of the state and action space. Our empirical simulations support our theoretical analysis. This suggests RLPA may offer significant advantages in large domains where some prior good policies are provided.

##### *Optimistic planning for belief-augmented Markov decision processes [11]*

This paper presents the Bayesian Optimistic Planning (BOP) algorithm, a novel model-based Bayesian reinforcement learning approach. BOP extends the planning approach of the Optimistic Planning for Markov Decision Processes (OP-MDP) algorithm [10], [9] to contexts where the transition model of the MDP is initially unknown and progressively learned through interactions within the environment. The knowledge about the unknown MDP is represented with a probability distribution over all possible transition models using Dirichlet distributions, and the BOP algorithm plans in the belief-augmented state space constructed by concatenating the original state vector with the current posterior distribution over transition models. We show that BOP becomes Bayesian optimal when the budget parameter increases to infinity. Preliminary empirical validations show promising performance.

##### *Aggregating optimistic planning trees for solving markov decision processes [16]*

This paper addresses the problem of online planning in Markov decision processes using a generative model and under a budget constraint. We propose a new algorithm, ASOP, which is based on the construction of a forest of single successor state planning trees, where each tree corresponds to a random realization of the stochastic environment. The trees are explored using a "safe" optimistic planning strategy which combines the optimistic principle (in order to explore the most promising part of the search space first) and a safety principle (which guarantees a certain amount of uniform exploration). In the decision-making step of the algorithm, the individual trees are aggregated and an immediate action is recommended. We provide a finite-sample analysis and discuss the trade-off between the principles of optimism and safety. We report numerical results on a benchmark problem showing that ASOP performs as well as state-of-the-art optimistic planning algorithms.

### ***Optimal Regret Bounds for Selecting the State Representation in Reinforcement Learning [20]***

We consider an agent interacting with an environment in a single stream of actions, observations, and rewards, with no reset. This process is not assumed to be a Markov Decision Process (MDP). Rather, the agent has several representations (mapping histories of past interactions to a discrete state space) of the environment with unknown dynamics, only some of which result in an MDP. The goal is to minimize the average regret criterion against an agent who knows an MDP representation giving the highest optimal reward, and acts optimally in it. Recent regret bounds for this setting are of order  $O(T^{2/3})$  with an additive term constant yet exponential in some characteristics of the optimal MDP. We propose an algorithm whose regret after  $T$  time steps is  $O(\sqrt{T})$ , with all constants reasonably small. This is optimal in  $T$  since  $O(\sqrt{T})$  is the optimal regret in the setting of learning in a (single discrete) MDP.

### ***Competing with an Infinite Set of Models in Reinforcement Learning [21]***

We consider a reinforcement learning setting where the learner also has to deal with the problem of finding a suitable state-representation function from a given set of models. This has to be done while interacting with the environment in an online fashion (no resets), and the goal is to have small regret with respect to any Markov model in the set. For this setting, recently the BLB algorithm has been proposed, which achieves regret of order  $T^{2/3}$ , provided that the given set of models is finite. Our first contribution is to extend this result to a countably infinite set of models. Moreover, the BLB regret bound suffers from an additive term that can be exponential in the diameter of the MDP involved, since the diameter has to be guessed. The algorithm we propose avoids guessing the diameter, thus improving the regret bound.

### ***A review of optimistic planning in Markov decision processes [30]***

We review a class of online planning algorithms for deterministic and stochastic optimal control problems, modeled as Markov decision processes. At each discrete time step, these algorithms maximize the predicted value of planning policies from the current state, and apply the first action of the best policy found. An overall receding-horizon algorithm results, which can also be seen as a type of model-predictive control. The space of planning policies is explored optimistically, focusing on areas with largest upper bounds on the value - or upper confidence bounds, in the stochastic case. The resulting optimistic planning framework integrates several types of optimism previously used in planning, optimization, and reinforcement learning, in order to obtain several intuitive algorithms with good performance guarantees. We describe in detail three recent such algorithms, outline the theoretical guarantees on their performance, and illustrate their behavior in a numerical example.

## **6.1.2. Multi-arm Bandit Theory**

### ***Automatic motor task selection via a bandit algorithm for a brain-controlled button [4]***

**Objective.** Brain-computer interfaces (BCIs) based on sensorimotor rhythms use a variety of motor tasks, such as imagining moving the right or left hand, the feet or the tongue. Finding the tasks that yield best performance, specifically to each user, is a time-consuming preliminary phase to a BCI experiment. This study presents a new adaptive procedure to automatically select (online) the most promising motor task for an asynchronous brain-controlled button. **Approach.** We develop for this purpose an adaptive algorithm UCB-classif based on the stochastic bandit theory and design an EEG experiment to test our method. We compare (offline) the adaptive algorithm to a naïve selection strategy which uses uniformly distributed samples from each task. We also run the adaptive algorithm online to fully validate the approach. **Main results.** By not wasting time on inefficient tasks, and focusing on the most promising ones, this algorithm results in a faster task selection and a more efficient use of the BCI training session. More precisely, the offline analysis reveals that the use of this algorithm can reduce the time needed to select the most appropriate task by almost half without loss in precision, or alternatively, allow us to investigate twice the number of tasks within a similar time span. Online tests confirm that the method leads to an optimal task selection. **Significance.** This study is the first one to optimize the task selection phase by an adaptive procedure. By increasing the number of tasks that can be tested in a given time span, the proposed method could contribute to reducing 'BCI illiteracy'.

***Kullback-Leibler Upper Confidence Bounds for Optimal Sequential Allocation [3]***

We consider optimal sequential allocation in the context of the so-called stochastic multi-armed bandit model. We describe a generic index policy, in the sense of Gittins (1979), based on upper confidence bounds of the arm payoffs computed using the Kullback-Leibler divergence. We consider two classes of distributions for which instances of this general idea are analyzed: The kl-UCB algorithm is designed for one-parameter exponential families and the empirical KL-UCB algorithm for bounded and finitely supported distributions. Our main contribution is a unified finite-time analysis of the regret of these algorithms that asymptotically matches the lower bounds of Lai and Robbins (1985) and Burnetas and Katehakis (1996), respectively. We also investigate the behavior of these algorithms when used with general bounded rewards, showing in particular that they provide significant improvements over the state-of-the-art.

***Sequential Transfer in Multi-armed Bandit with Finite Set of Models [14]***

Learning from prior tasks and transferring that experience to improve future performance is critical for building lifelong learning agents. Although results in supervised and reinforcement learning show that transfer may significantly improve the learning performance, most of the literature on transfer is focused on batch learning tasks. In this paper we study the problem of *sequential transfer in online learning*, notably in the multi-armed bandit framework, where the objective is to minimize the total regret over a sequence of tasks by transferring knowledge from prior tasks. Under the assumption that the tasks are drawn from a stationary distribution over a finite set of models, we define a novel bandit algorithm based on a method-of-moments approach for the estimation of the possible tasks and derive regret bounds for it. We introduce a novel bandit algorithm based on a method-of-moments approach for estimating the possible tasks and derive regret bounds for it. Finally, we report preliminary empirical results confirming the theoretical findings.

***Optimizing P300-speller sequences by RIP-ping groups apart [25]***

So far P300-speller design has put very little emphasis on the design of optimized flash patterns, a surprising fact given the importance of the sequence of flashes on the selection outcome. Previous work in this domain has consisted in studying consecutive flashes, to prevent the same letter or its neighbors from flashing consecutively. To this effect, the flashing letters form more random groups than the original row-column sequences for the P300 paradigm, but the groups remain fixed across repetitions. This has several important consequences, among which a lack of discrepancy between the scores of the different letters. The new approach proposed in this paper accumulates evidence for individual elements, and optimizes the sequences by relaxing the constraint that letters should belong to fixed groups across repetitions. The method is inspired by the theory of Restricted Isometry Property matrices in Compressed Sensing, and it can be applied to any display grid size, and for any target flash frequency. This leads to P300 sequences which are shown here to perform significantly better than the state of the art, in simulations and online tests.

***Stochastic Simultaneous Optimistic Optimization [26]***

We study the problem of global maximization of a function  $f$  given a finite number of evaluations perturbed by noise. We consider a very weak assumption on the function, namely that it is locally smooth (in some precise sense) with respect to some semi-metric, around one of its global maxima. Compared to previous works on bandits in general spaces (Kleinberg et al., 2008; Bubeck et al., 2011a) our algorithm does not require the knowledge of this semi-metric. Our algorithm, StoSOO, follows an optimistic strategy to iteratively construct upper confidence bounds over the hierarchical partitions of the function domain to decide which point to sample next. A finite-time analysis of StoSOO shows that it performs almost as well as the best specifically-tuned algorithms even though the local smoothness of the function is not known.

***Toward optimal stratification for stratified monte-carlo integration [9]***

We consider the problem of adaptive stratified sampling for Monte Carlo integration of a noisy function, given a finite budget  $n$  of noisy evaluations to the function. We tackle in this paper the problem of adapting to the function at the same time the number of samples into each stratum and the partition itself. More precisely, it is interesting to refine the partition of the domain in area where the noise to the function, or where the variations

of the function, are very heterogeneous. On the other hand, having a (too) refined stratification is not optimal. Indeed, the more refined the stratification, the more difficult it is to adjust the allocation of the samples to the stratification, i.e. sample more points where the noise or variations of the function are larger. We provide in this paper an algorithm that selects online, among a large class of partitions, the partition that provides the optimal trade-off, and allocates the samples almost optimally on this partition

### ***Thompson sampling for one-dimensional exponential family bandits [18]***

Thompson Sampling has been demonstrated in many complex bandit models, however the theoretical guarantees available for the parametric multi-armed bandit are still limited to the Bernoulli case. Here we extend them by proving asymptotic optimality of the algorithm using the Jeffreys prior for 1-dimensional exponential family bandits. Our proof builds on previous work, but also makes extensive use of closed forms for Kullback-Leibler divergence and Fisher information (and thus Jeffreys prior) available in an exponential family. This allow us to give a finite time exponential concentration inequality for posterior distributions on exponential families that may be of interest in its own right. Moreover our analysis covers some distributions for which no optimistic algorithm has yet been proposed, including heavy-tailed exponential families.

### ***Finite-Time Analysis of Kernelised Contextual Bandits [27]***

We tackle the problem of online reward maximisation over a large finite set of actions described by their contexts. We focus on the case when the number of actions is too big to sample all of them even once. However we assume that we have access to the similarities between actions' contexts and that the expected reward is an arbitrary linear function of the contexts' images in the related reproducing kernel Hilbert space (RKHS). We propose KernelUCB, a kernelised UCB algorithm, and give a cumulative regret bound through a frequentist analysis. For contextual bandits, the related algorithm GP-UCB turns out to be a special case of our algorithm, and our finite-time analysis improves the regret bound of GP-UCB for the agnostic case, both in the terms of the kernel-dependent quantity and the RKHS norm of the reward function. Moreover, for the linear kernel, our regret bound matches the lower bound for contextual linear bandits.

### ***From Bandits to Monte-Carlo Tree Search: The Optimistic Principle Applied to Optimization and Planning [33]***

This work covers several aspects of the optimism in the face of uncertainty principle applied to large scale optimization problems under finite numerical budget. The initial motivation for the research reported here originated from the empirical success of the so-called Monte-Carlo Tree Search method popularized in computer-go and further extended to many other games as well as optimization and planning problems. Our objective is to contribute to the development of theoretical foundations of the field by characterizing the complexity of the underlying optimization problems and designing efficient algorithms with performance guarantees. The main idea presented here is that it is possible to decompose a complex decision making problem (such as an optimization problem in a large search space) into a sequence of elementary decisions, where each decision of the sequence is solved using a (stochastic) multi-armed bandit (simple mathematical model for decision making in stochastic environments). This so-called hierarchical bandit approach (where the reward observed by a bandit in the hierarchy is itself the return of another bandit at a deeper level) possesses the nice feature of starting the exploration by a quasi-uniform sampling of the space and then focusing progressively on the most promising area, at different scales, according to the evaluations observed so far, and eventually performing a local search around the global optima of the function. The performance of the method is assessed in terms of the optimality of the returned solution as a function of the number of function evaluations. Our main contribution to the field of function optimization is a class of hierarchical optimistic algorithms designed for general search spaces (such as metric spaces, trees, graphs, Euclidean spaces, ...) with different algorithmic instantiations depending on whether the evaluations are noisy or noiseless and whether some measure of the "smoothness" of the function is known or unknown. The performance of the algorithms depend on the local behavior of the function around its global optima expressed in terms of the quantity of near-optimal states measured with some metric. If this local smoothness of the function is known then one can design very efficient optimization algorithms (with convergence rate independent of the space dimension),

and when it is not known, we can build adaptive techniques that can, in some cases, perform almost as well as when it is known.

## 6.2. Statistical analysis of time series

### 6.2.1. Change Point Analysis

#### *Nonparametric multiple change point estimation in highly dependent time series [17]*

Given a heterogeneous time-series sample, it is required to find the points in time (called change points) where the probability distribution generating the data has changed. The data is assumed to have been generated by arbitrary, unknown, stationary ergodic distributions. No modeling, independence or mixing are made. A novel, computationally efficient, nonparametric method is proposed, and is shown to be asymptotically consistent in this general framework; the theoretical results are complemented with experimental evaluations.

### 6.2.2. Clustering Time Series, Online and Offline

#### *A Binary-Classification-Based Metric between Time-Series Distributions and Its Use in Statistical and Learning Problems [6]*

A metric between time-series distributions is proposed that can be evaluated using binary classification methods, which were originally developed to work on i.i.d. data. It is shown how this metric can be used for solving statistical problems that are seemingly unrelated to classification and concern highly dependent time series. Specifically, the problems of time-series clustering, homogeneity testing and the three-sample problem are addressed. Universal consistency of the resulting algorithms is proven under most general assumptions. The theoretical results are illustrated with experiments on synthetic and real-world data.

### 6.2.3. Semi-Supervised and Unsupervised Learning

#### *Learning from a Single Labeled Face and a Stream of Unlabeled Data [19]*

Face recognition from a single image per person is a challenging problem because the training sample is extremely small. We consider a variation of this problem. In our problem, we recognize only one person, and there are no labeled data for any other person. This setting naturally arises in authentication on personal computers and mobile devices, and poses additional challenges because it lacks negative examples. We formalize our problem as one-class classification, and propose and analyze an algorithm that learns a non-parametric model of the face from a single labeled image and a stream of unlabeled data. In many domains, for instance when a person interacts with a computer with a camera, unlabeled data are abundant and easy to utilize. This is the first paper that investigates how these data can help in learning better models in the single-image-per-person setting. Our method is evaluated on a dataset of 43 people and we show that these people can be recognized 90% of time at nearly zero false positives. This recall is 25+% higher than the recall of our best performing baseline. Finally, we conduct a comprehensive sensitivity analysis of our algorithm and provide a guideline for setting its parameters in practice.

#### *Unsupervised model-free representation learning [23]*

Numerous control and learning problems face the situation where sequences of high-dimensional highly dependent data are available, but no or little feedback is provided to the learner. In such situations it may be useful to find a concise representation of the input signal, that would preserve as much as possible of the relevant information. In this work we are interested in the problems where the relevant information is in the time-series dependence. Thus, the problem can be formalized as follows. Given a series of observations  $X_0, \dots, X_n$  coming from a large (high-dimensional) space  $\mathcal{X}$ , find a representation function  $f$  mapping  $\mathcal{X}$  to a finite space  $\mathcal{Y}$  such that the series  $f(X_0), \dots, f(X_n)$  preserve as much information as possible about the original time-series dependence in  $X_0, \dots, X_n$ . For stationary time series, the function  $f$  can be selected as the one maximizing the time-series information  $I_\infty(f) = h_0(f(X)) - h_\infty(f(X))$  where  $h_0(f(X))$  is the Shannon entropy of  $f(X_0)$  and  $h_\infty(f(X))$  is the entropy rate of the time series

$f(X_0), \dots, f(X_n), \dots$ . In this paper we study the functional  $I_\infty(f)$  from the learning-theoretic point of view. Specifically, we provide some uniform approximation results, and study the behaviour of  $I_\infty(f)$  in the problem of optimal control.

#### ***Time-series information and learning [22]***

Given a time series  $X_1, \dots, X_n, \dots$  taking values in a large (high-dimensional) space  $\mathcal{X}$ , we would like to find a function  $f$  from  $\mathcal{X}$  to a small (low-dimensional or finite) space  $\mathcal{Y}$  such that the time series  $f(X_1), \dots, f(X_n), \dots$  retains all the information about the time-series dependence in the original sequence, or as much as possible thereof. This goal is formalized in this work, and it is shown that the target function  $f$  can be found as the one that maximizes a certain quantity that can be expressed in terms of entropies of the series  $(f(X_i))_{i \in \mathcal{N}}$ . This quantity can be estimated empirically, and does not involve estimating the distribution on the original time series  $(X_i)_{i \in \mathcal{N}}$ .

## **6.3. Statistical Learning and Bayesian Analysis**

### **6.3.1. Dictionary learning**

#### ***Learning a common dictionary over a sensor network [10]***

We consider the problem of distributed dictionary learning, where a set of nodes is required to collectively learn a common dictionary from noisy measurements. This approach may be useful in several contexts including sensor networks. Diffusion cooperation schemes have been proposed to solve the distributed linear regression problem. In this work we focus on a diffusion-based adaptive dictionary learning strategy: each node records independent observations and cooperates with its neighbors by sharing its local dictionary. The resulting algorithm corresponds to a distributed alternate optimization. Beyond dictionary learning, this strategy could be adapted to many matrix factorization problems in various settings. We illustrate its efficiency on some numerical experiments.

#### ***Distributed dictionary learning over a sensor network [29]***

We consider the problem of distributed dictionary learning, where a set of nodes is required to collectively learn a common dictionary from noisy measurements. This approach may be useful in several contexts including sensor networks. Diffusion cooperation schemes have been proposed to solve the distributed linear regression problem. In this work we focus on a diffusion-based adaptive dictionary learning strategy: each node records observations and cooperates with its neighbors by sharing its local dictionary. The resulting algorithm corresponds to a distributed block coordinate descent (alternate optimization). Beyond dictionary learning, this strategy could be adapted to many matrix factorization problems and generalized to various settings. This article presents our approach and illustrates its efficiency on some numerical examples.

## **6.4. Applications**

### **6.4.1. Medical Applications**

#### ***Outlier detection for patient monitoring and alerting [5]***

We develop and evaluate a data-driven approach for detecting unusual (anomalous) patient-management decisions using past patient cases stored in electronic health records (EHRs). Our hypothesis is that a patient-management decision that is unusual with respect to past patient care may be due to an error and that it is worthwhile to generate an alert if such a decision is encountered. We evaluate this hypothesis using data obtained from EHRs of 4486 post-cardiac surgical patients and a subset of 222 alerts generated from the data. We base the evaluation on the opinions of a panel of experts. The results of the study support our hypothesis that the outlier-based alerting can lead to promising true alert rates. We observed true alert rates that ranged from 25% to 66% for a variety of patient-management actions, with 66% corresponding to the strongest outliers.

## 6.5. Miscellaneous

### 6.5.1. Miscellaneous

#### *A confidence-set approach to signal denoising [7]*

The problem of filtering of finite-alphabet stationary ergodic time series is considered. A method for constructing a confidence set for the (unknown) signal is proposed, such that the resulting set has the following properties. First, it includes the unknown signal with probability  $\gamma$ , where  $\gamma$  is a parameter supplied to the filter. Second, the size of the confidence sets grows exponentially with a rate that is asymptotically equal to the conditional entropy of the signal given the data. Moreover, it is shown that this rate is optimal. We also show that the described construction of the confidence set can be applied to the case where the signal is corrupted by an erasure channel with unknown statistics.

#### *Quantification adaptative pour la stéganalyse d'images texturées [28]*

Nous cherchons à améliorer les performances d'un schéma de stéganalyse (i.e. la détection de messages cachés) pour des images texturées. Le schéma de stéganographie étudié consiste à modifier certains pixels de l'image par une perturbation  $\pm 1$ , et le schéma de stéganalyse utilise les caractéristiques construites à partir de la probabilité conditionnelle empirique de différences de 4 pixels voisins. Dans sa version originale, la stéganalyse n'est pas très efficace sur des images texturées et ce travail vise à explorer plusieurs techniques de quantification en utilisant d'abord un pas de quantification plus important puis une quantification adaptative scalaire ou vectorielle. Les cellules de la quantification adaptative sont générées en utilisant un K-means ou un K-means "équilibré" de manière à ce que chaque cellule quantifie approximativement le même nombre d'échantillon. Nous obtenons un gain maximal de classification de 3% pour un pas de quantification uniforme de 3. En utilisant l'algorithme K-means équilibré sur  $[-18,18]$ , le gain par rapport à la version de base est de 4.7%.

#### *Cost-sensitive Multiclass Classification Risk Bounds [8]*

A commonly used approach to multiclass classification is to replace the 0-1 loss with a convex surrogate so as to make empirical risk minimization computationally tractable. Previous work has uncovered sufficient and necessary conditions for the consistency of the resulting procedures. In this paper, we strengthen these results by showing how the 0-1 excess loss of a predictor can be upper bounded as a function of the excess loss of the predictor measured using the convex surrogate. The bound is developed for the case of cost-sensitive multiclass classification and a convex surrogate loss that goes back to the work of Lee, Lin and Wahba. The bounds are as easy to calculate as in binary classification. Furthermore, we also show that our analysis extends to the analysis of the recently introduced "Simplex Coding" scheme.

#### *Approximate Dynamic Programming Finally Performs Well in the Game of Tetris [12]*

Tetris is a video game that has been widely used as a benchmark for various optimization techniques including approximate dynamic programming (ADP) algorithms. A look at the literature of this game shows that while ADP algorithms that have been (almost) entirely based on approximating the value function (value function based) have performed poorly in Tetris, the methods that search directly in the space of policies by learning the policy parameters using an optimization black box, such as the cross entropy (CE) method, have achieved the best reported results. This makes us conjecture that Tetris is a game in which good policies are easier to represent, and thus, learn than their corresponding value functions. So, in order to obtain a good performance with ADP, we should use ADP algorithms that search in a policy space, instead of the more traditional ones that search in a value function space. In this paper, we put our conjecture to test by applying such an ADP algorithm, called classification-based modified policy iteration (CBMPI), to the game of Tetris. Our experimental results show that for the first time an ADP algorithm, namely CBMPI, obtains the best results reported in the literature for Tetris in both small  $10 \times 10$  and large  $10 \times 20$  boards. Although the CBMPI's results are similar to those of the CE method in the large board, CBMPI uses considerably fewer (almost 1/6) samples (calls to the generative model) than CE.



### ***A Generalized Kernel Approach to Structured Output Learning [15]***

We study the problem of structured output learning from a regression perspective. We first provide a general formulation of the kernel dependency estimation (KDE) problem using operator-valued kernels. We show that some of the existing formulations of this problem are special cases of our framework. We then propose a covariance-based operator-valued kernel that allows us to take into account the structure of the kernel feature space. This kernel operates on the output space and encodes the interactions between the outputs without any reference to the input space. To address this issue, we introduce a variant of our KDE method based on the conditional covariance operator that in addition to the correlation between the outputs takes into account the effects of the input variables. Finally, we evaluate the performance of our KDE approach using both covariance and conditional covariance kernels on two structured output problems, and compare it to the state-of-the-art kernel-based structured output regression methods.

### ***Gossip-based distributed stochastic bandit algorithms [24]***

The multi-armed bandit problem has attracted remarkable attention in the machine learning community and many efficient algorithms have been proposed to handle the so-called exploitation-exploration dilemma in various bandit setups. At the same time, significantly less effort has been devoted to adapting bandit algorithms to particular architectures, such as sensor networks, multi-core machines, or peer-to-peer (P2P) environments, which could potentially speed up their convergence. Our goal is to adapt stochastic bandit algorithms to P2P networks. In our setup, the same set of arms is available in each peer. In every iteration each peer can pull one arm independently of the other peers, and then some limited communication is possible with a few random other peers. As our main result, we show that our adaptation achieves a linear speedup in terms of the number of peers participating in the network. More precisely, we show that the probability of playing a suboptimal arm at a peer in iteration  $t = \Omega(\log N)$  is proportional to  $1/(Nt)$  where  $N$  denotes the number of peers. The theoretical results are supported by simulation experiments showing that our algorithm scales gracefully with the size of network.

### ***Sur quelques problèmes non-supervisés impliquant des séries temporelles hautement dépendantes [1]***

Cette thèse est consacrée à l'analyse théorique de problèmes non supervisés impliquant des séries temporelles hautement dépendantes. Plus particulièrement, nous abordons les deux problèmes fondamentaux que sont le problème d'estimation des points de rupture et le partitionnement de séries temporelles. Ces problèmes sont abordés dans un cadre extrêmement général où les données sont générées par des processus stochastiques ergodiques stationnaires. Il s'agit de l'une des hypothèses les plus faibles en statistiques, comprenant non seulement, les hypothèses de modèles et les hypothèses paramétriques habituelles dans la littérature scientifique, mais aussi des hypothèses classiques d'indépendance, de contraintes sur l'espace mémoire ou encore des hypothèses de mélange. En particulier, aucune restriction n'est faite sur la forme ou la nature des dépendances, de telles sortes que les échantillons peuvent être arbitrairement dépendants. Pour chaque problème abordé, nous proposons de nouvelles méthodes non paramétriques et nous prouvons de plus qu'elles sont, dans ce cadre, asymptotiquement consistantes. Pour l'estimation de points de rupture, la consistance asymptotique se rapporte à la capacité de l'algorithme à produire des estimations des points de rupture qui sont asymptotiquement arbitrairement proches des vrais points de rupture. D'autre part, un algorithme de partitionnement est asymptotiquement consistant si le partitionnement qu'il produit, restreint à chaque lot de séquences, coïncides, à partir d'un certain temps et de manière consistante, avec le partitionnement cible. Nous montrons que les algorithmes proposés sont implémentables efficacement, et nous accompagnons nos résultats théoriques par des évaluations expérimentales. L'analyse statistique dans le cadre stationnaire ergodique est extrêmement difficile. De manière générale, il est prouvé que les vitesses de convergence sont impossibles à obtenir. Dès lors, pour deux échantillons générés indépendamment par des processus ergodiques stationnaires, il est prouvé qu'il est impossible de distinguer le cas où les échantillons sont générés par le même processus de celui où ils sont générés par des processus différents. Ceci implique que des problèmes tels le partitionnement de séries temporelles sans la connaissance du nombre de partitions ou du nombre de points de rupture ne peut admettre de solutions consistantes. En conséquence, une tâche difficile est de découvrir les formulations du problème qui en permettent une résolution dans ce cadre général. La principale contribution de cette thèse est de démon-

trer (par construction) que malgré ces résultats d'impossibilités théoriques, des formulations naturelles des problèmes considérés existent et admettent des solutions consistantes dans ce cadre général. Ceci inclut la démonstration du fait que le nombre de points de rupture corrects peut être trouvé, sans recourir à des hypothèses plus fortes sur les processus stochastiques. Il en résulte que, dans cette formulation, le problème des points de rupture peut être réduit à du partitionnement de séries temporelles. Les résultats présentés dans ce travail forment les fondations théoriques pour l'analyse des données séquentielles dans un espace d'applications bien plus large.

#### ***Actor-Critic Algorithms for Risk-Sensitive MDPs [32]***

In many sequential decision-making problems we may want to manage risk by minimizing some measure of variability in rewards in addition to maximizing a standard criterion. Variance-related risk measures are among the most common risk-sensitive criteria in finance and operations research. However, optimizing many such criteria is known to be a hard problem. In this paper, we consider both discounted and average reward Markov decision processes. For each formulation, we first define a measure of variability for a policy, which in turn gives us a set of risk-sensitive criteria to optimize. For each of these criteria, we derive a formula for computing its gradient. We then devise actor-critic algorithms for estimating the gradient and updating the policy parameters in the ascent direction. We establish the convergence of our algorithms to locally risk-sensitive optimal policies. Finally, we demonstrate the usefulness of our algorithms in a traffic signal control application.

#### ***Bayesian Policy Gradient and Actor-Critic Algorithms [31]***

Policy gradient methods are reinforcement learning algorithms that adapt a parameterized policy by following a performance gradient estimate. Many conventional policy gradient methods use Monte-Carlo techniques to estimate this gradient. The policy is improved by adjusting the parameters in the direction of the gradient estimate. Since Monte-Carlo methods tend to have high variance, a large number of samples is required to attain accurate estimates, resulting in slow convergence. In this paper, we first propose a Bayesian framework for policy gradient, based on modeling the policy gradient as a Gaussian process. This reduces the number of samples needed to obtain accurate gradient estimates. Moreover, estimates of the natural gradient as well as a measure of the uncertainty in the gradient estimates, namely, the gradient covariance, are provided at little extra cost. Since the proposed Bayesian framework considers system trajectories as its basic observable unit, it does not require the dynamic within each trajectory to be of any special form, and thus, can be easily extended to partially observable problems. On the downside, it cannot take advantage of the Markov property when the system is Markovian. To address this issue, we then extend our Bayesian policy gradient framework to actor-critic algorithms and present a new actor-critic learning model in which a Bayesian class of non-parametric critics, based on Gaussian process temporal difference learning, is used. Such critics model the action-value function as a Gaussian process, allowing Bayes' rule to be used in computing the posterior distribution over action-value functions, conditioned on the observed data. Appropriate choices of the policy parameterization and of the prior covariance (kernel) between action-values allow us to obtain closed-form expressions for the posterior distribution of the gradient of the expected return with respect to the policy parameters. We perform detailed experimental comparisons of the proposed Bayesian policy gradient and actor-critic algorithms with classic Monte-Carlo based policy gradient methods, as well as with each other, on a number of reinforcement learning problems.

## SIMPAF Project-Team

## 6. New Results

### 6.1. Quantitative homogenization theory

In collaboration with S. Neukamm and F. Otto, A. Gloria developed in [46] and [45] a quantitative approach of the stochastic homogenization of discrete elliptic equations. There are two main achievements. In [46] we developed a general theory which quantifies optimally in time the decay of the non-constant coefficients semi-group associated with discrete random diffusion coefficients satisfying a spectral gap assumption (namely, the environment seen from the particle). Combined with spectral theory this allowed us to make a sharp numerical analysis of the popular periodization method to approximate homogenized coefficients. In [45], we obtained a quantitative two-scale expansion result, and essentially proved that the difference between the solution of a (discrete) elliptic equation with random coefficients on the torus and the first two terms of the two-scale expansion scales as in the periodic case (except in dimension 2, for which there is a logarithmic correction).

### 6.2. Corrosion

The Diffusion Poisson Coupled Model [32] is a model of iron based alloy in a nuclear waste repository. It describes the growth of an oxide layer in this framework. The system is made of a Poisson equation on the electrostatic potential and convection-diffusion equations on the densities of charge carriers (electrons, ferric cations and oxygen vacancies), supplemented with coupled Robin boundary conditions. The DPCM model also takes into account the growth of the oxide host lattice and its dissolution, leading to moving boundary equations. In [44], C. Chainais-Hillairet and I. Lacroix-Violet consider a simplified version of this model, where only two charge carriers are taken into account and where there is no evolution of the layer thickness. They prove the existence of a steady-state solution to this model. More recently, C. Chainais-Hillairet and I. Lacroix-Violet have also obtained an existence result for the time-dependent simplified model. This result is submitted for publication [47].

P.-L. Colin, C. Chainais-Hillairet and I. Lacroix-Violet have recently performed the numerical analysis of the numerical scheme presented in [31]. The scheme is a Euler implicit in time scheme with Scharfetter-Gummel approximation of the convection-diffusion fluxes. They prove existence of a solution to the scheme, a priori estimates satisfied by the solution and convergence of the numerical scheme to a weak solution of the corrosion model.

Numerical experiments done for the simulation of the full DPCM model with moving boundaries shows the convergence in time towards a pseudo-steady-state. T. Gallouët has proposed a new scheme in order to compute directly this pseudo-steady-state. This scheme has been implemented in the code CALIPSO (ANDRA). Validation is in progress, as the numerical analysis of the scheme.

### 6.3. New results on finite volume schemes

In [5], C. Chainais-Hillairet, S. Krell and A. Mouton develop Discrete Duality Finite Volume methods for the finite volume approximation of a system describing miscible displacement in porous media (Peaceman model). They establish relevant a priori estimates satisfied by the numerical solution and prove existence and uniqueness of the solution to the scheme. They show the efficiency of the schemes through numerical experiments. Recently, they also proved the convergence of the DDFV scheme for the Peaceman model. This work will be soon submitted for publication.

In [35], M. Bessemoulin-Chatard, C. Chainais-Hillairet and F. Filbet prove several discrete Gagliardo-Nirenberg-Sobolev and Poincaré-Sobolev inequalities for some approximations with arbitrary boundary values on finite volume meshes. The keypoint of their approach is to use the continuous embedding of the space  $BV(\Omega)$  into  $L^{N/(N-1)}(\Omega)$  for a Lipschitz domain  $\Omega \subset \mathbb{R}^N$ , with  $N \geq 2$ . Finally, they give several applications to discrete duality finite volume (DDFV) schemes which are used for the approximation of nonlinear and on isotropic elliptic and parabolic problems.

In [22], M. Bessemoulin-Chatard, C. Chainais-Hillairet and M.-H. Vignal consider the numerical approximation of the classical time-dependent drift-diffusion system near quasi-neutrality by a fully implicit in time and finite volume in space scheme, where the convection-diffusion fluxes are approximated by Scharfetter-Gummel fluxes. They establish that all the a priori estimates needed to prove the convergence of the scheme does not depend on the Debye length  $\lambda$ . This proves that the scheme is asymptotic preserving in the quasi-neutral limit  $\lambda \rightarrow 0$ .

In [24], C. Chainais-Hillairet, A. Jüngel and S. Schuchnigg prove the time decay of fully discrete finite-volume approximations of porous-medium and fast-diffusion equations with Neumann or periodic boundary conditions in the entropy sense. The algebraic or exponential decay rates are computed explicitly. In particular, the numerical scheme dissipates all zeroth-order entropies which are dissipated by the continuous equation. The proofs are based on novel continuous and discrete generalized Beckner inequalities.

## 6.4. New results in numerical fluid dynamics

In the case of compressible models, as the Euler equations, a careful analysis of sharp and practical stability conditions to ensure the positivity of both density and pressure variables was performed[4]. We are also concerned with the numerical simulation of certain multi-fluids flows, which in particular arises in the modeling of powder/snow avalanches. The hybrid scheme works on unstructured meshes and can be advantageously coupled to mesh refinements strategies in order to follow fronts of high density variation [42]. In particular, we investigate the influence of the characteristics Froude number, Schmidt number and Reynolds number on the front progression. In the context of the PhD thesis of Meriem Ezzoug (University of Monastir, Tunisia), co-advised by C. Calgaro and E. Zahrouni (University of Monastir, Tunisia), we investigate theoretically and numerically the influence of a specific stress tensor, introduced for the first time by Korteweg, in some diffuse interface models which allow to describe some phase transition phenomena, such as surface tension force formulation for multiphase fluid flows. In order to answer these questions, we have developed respectively a Fortran code, a C++ code (NS2DDV-C++, see the softwares section) and a MATLAB code (NS2DDV-M, see the softwares section).

## 6.5. New results on a posteriori estimates

Some residual-type a posteriori error estimators were developed in the context of magnetostatic and magnetodynamic Maxwell equations, given in their potential and harmonic formulations. Here, the task was to find a relevant decomposition of the error in order to obtain the reliability of the estimator, with the use of ad-hoc interpolations. This work was realized in collaboration with the L2EP Laboratory (Laboratoire d'Electrotechnique et d'Electronique de Puissance de Lille, Lille 1 University), and gave rise to several contributions [7], [18], [19], [20], [21], obtained in the context of the Ph-D thesis and of the Post-doc position of Zuqi Tang. Then, other results about a posteriori error estimators were obtained in other contexts [6], [8]. Recently, we started working on space/time error estimators for finite element methods, arising in the context of low-frequency Maxwell equations (PhD of R. Tittarelli, CIFRE EDF R&D, see [25]).

## 6.6. New results in control in fluid mechanics

Recently, we studied more particularly passive control techniques using porous media for incompressible aerodynamics on several bodies, with the use of the penalisation method [3].

## BONSAI Project-Team

### 6. New Results

#### 6.1. High-throughput sequence processing

- Within our collaboration with Montpellier (IRB and LIRMM) we published a paper on CRAC, a software for analysing short RNA sequences and detecting variations among them [5].
- We have been invited to contribute an invited book chapter on metatranscriptomic data analysis (*Methods in Molecular Biology*, in press). This chapter covers the complete bioinformatic analysis from raw reads to taxonomic assignation, and introduces our software SortMeRNA (see Paragraph 5.6). This is a joint work with team LABIS in Genoscope.
- Evguenia Kopylova defended her thesis on December, the 11th ("*New algorithmic and bioinformatic approaches for the analysis of data from high throughput sequencing*", [1]). The second part of her work deals with a new read mapper for metagenomic sequence data.
- Within our collaboration with the Lille hospital, we developed a seed-based heuristics for the detection of lymphocyte rearrangements from high-throughput data. This method is implemented in the software Vidjil (see Section 5.7). Our results were presented at the Jobim conference [8], and a journal article was submitted.

#### 6.2. RNA algorithms

- We have started a new collaborative project with Bielefeld Universität on an extension of *Algebraic Dynamic Programming*. We introduced a generic specification framework, called *inverted coupled rewrite systems* [9], that can deal with optimization problems on strings, trees, and arc-annotated sequences. It is based on the following ideas: the solutions of combinatorial optimization problems are the inverse image of a term rewrite relation that reduces problem solutions to problem inputs. A tree grammar is used to further refine the search space, and optimization objectives are specified as interpretations of these terms. All these constituents provide a mathematically precise and complete problem specification, leading to concise yet translucent specifications of dynamic programming algorithms.

#### 6.3. Genomic rearrangements

- Within a collaboration with LIAFA (CNRS UMR 7089, and University Paris 7) we published a method for the assembling of ancestral gene orders from contiguous ancestral fragments [4].

#### 6.4. Nonribosomal peptides

- Yoann Dufresne is starting a PhD thesis on computational biology for nonribosomal peptides (NRPs) under the supervision of Maude Pupin and Laurent Noe, after doing his master thesis with them. He already worked on the translation of the chemical structure of the NRPs into their monomeric structure. NRPs can be represented by their chemical structure that is a graph where the atoms are represented by nodes and the chemical bonds by arcs; or by their monomeric structure that is a graph where the monomers are represented by nodes and the chemical bonds between monomers by arcs. We designed a novel algorithm capable of localizing the monomers from a reference list in the chemical structures of peptides [7]. It is based on a heuristic that utilizes chemical information of NRPs. The preliminary results are encouraging, and should lead to further studies.

## SHACRA Project-Team

### 6. New Results

#### 6.1. Electrophysiology

Cardiac arrhythmia is a very frequent pathology that comes from an abnormal electrical activity in the myocardium. This work aims at developing a training simulator for interventional radiology and thermo-ablation of these arrhythmias. After tackling the issue of fast electrophysiology, a first version of our training simulator was proposed.



*Figure 3. Cardiac electrophysiology computed on a patient-specific geometry*

The first main contribution of this work is the interactive catheter navigation inside a moving venous system and a beating heart. The virtual catheterization reproduces navigation issues that can be solved using a bending catheter. Second, our real-time GPU electrophysiology model allows interactions during the simulation such as extra-cellular potential measurement, RF ablation, and electrical stimulation. An innovative management of the computational units based on multithreading offers performances close to real-time. This framework is therefore a substantial step towards realistic and highly efficient virtual training systems in cardiology. As future work, we intend to use patient-specific data in our framework so that cardiologists could quantitatively assess the realism of our virtual training.

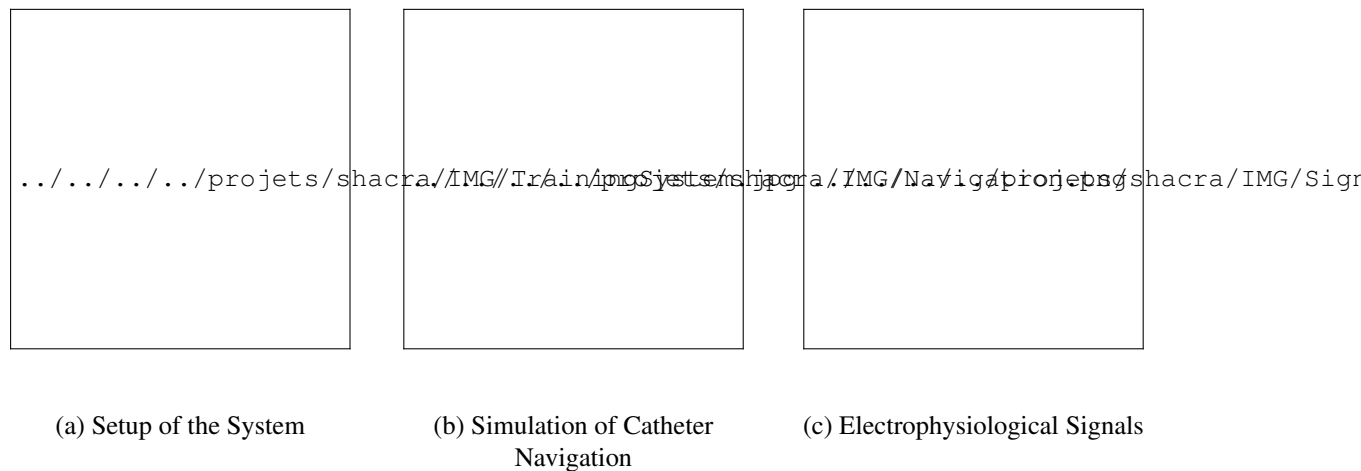


Figure 4. The first simulation dedicated to electrocardiology training

## 6.2. Cryoablation

A new project started this year around cryotherapy. This technique consists in inserting needles that freezing the surrounding tissues, thus immediately leading to cellular death of the tissues. Cryoablation procedure is used in many medical fields for tumor ablation, and even starts being used in cardiology. In this scope, we build a simulator able to place the cryoprobes and run a simulation representing the evolution of iceballs in living tissues.

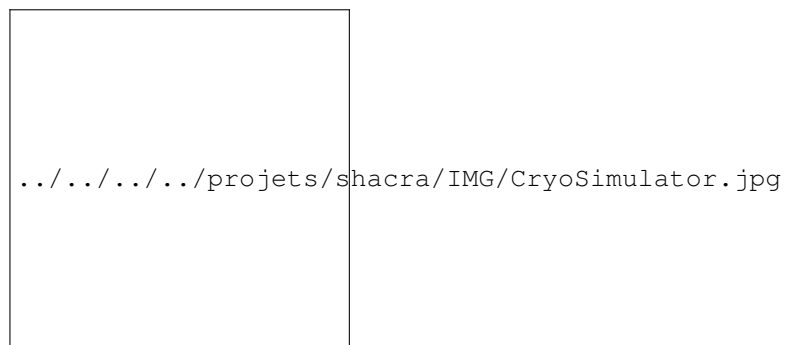


Figure 5. Simulation framework for cryoablation planning

## 6.3. Stapedotomy

Stapedotomy is a challenging procedure of the middle ear microsurgery, since the surgeons is in direct contact with sensitive structures such as the ossicular chain. This procedures is taught and performed in the last phase of the surgical apprenticeship. To improve surgical teaching, we propose to use a virtual surgical simulator based on a finite element model of the middle ear. The static and dynamic behavior of the developed finite element model was successfully compared to published data on human temporal bones specimens. A semi-automatic algorithm was developed to perform a quick and accurate registration of our validated mechanical

atlas to match the patient dataset. This method avoids a time-consuming work of manual segmentation, parameterization, and evaluation. A registration is obtained in less than 260 seconds with an accuracy close to a manual process and within the imagery resolution. The computation algorithms, allowing carving, deformation of soft and hard tissues, and collision response, are compatible with a real-time interactive simulation of a middle ear procedure. As a future work, we propose to investigate new robotized procedures of the middle ear surgery in order to develop new applications for the RobOtol device and to provide a training tool for the surgeons.



Figure 6. Simulation of the stapedotomy procedure.

## 6.4. Radiotherapy planning

The main challenge of radiotherapy treatment is to irradiate the tumor while sparing the surrounding healthy tissues. In the case of throat cancer, the complexity of the therapy treatment is due to the proximity of organs at risk such as the two parotid glands. The parotid glands are the main salivary glands. An overdose of radiation in these glands may cause xerostomia, which is a medical term for the symptom of dryness in the mouth, or in other words, a lack of saliva. This disease affects significantly the life of the patient: difficulty talking, tasting, chewing, swallowing, excessive thirst, constant pain in the throat etcetera. A radiation therapy treatment of throat cancer takes from 5 to 7 weeks. The treatment is planned several days before the therapy. The planning consists in contouring each organ of the area on CT-scan images and defining the dose of radiation to deliver to each of these organs. This stage is lengthy and takes around two hours per patient. Yet, some anatomical variations occur in the course of the treatment, mainly due to the weight loss of the patients. These variations compromise the safety of the healthy tissues, because the planned treatments are no more up to date. For now, the physicians have no solution good enough to handle these changes. Xerostomia affects around 20 per-cent of the patients suffering from throat cancer.

The main idea of this work is to create an interface that the physicians could use to redo the planning when it is needed, when the anatomical changes are significant. The purpose is to give to them the possibility to use what they see on images, to recreate the right shape of the contours without recontouring each image, and in a reasonable time. This interface will use their knowledge to determine the new shape of the organs. The work does not aim at providing a fully automatic method because it would reduce its acceptance by the physicians. As the method is based on the input of the physician, they can control the deformation based on images but also on their knowledge.

## 6.5. Image-based diagnoses

In the context of the female pelvic medicine, image-based diagnoses of pelvic floor disorders like prolapse or endometriosis rely on mechanical indicators, such as mobilities of organs and shear displacements between



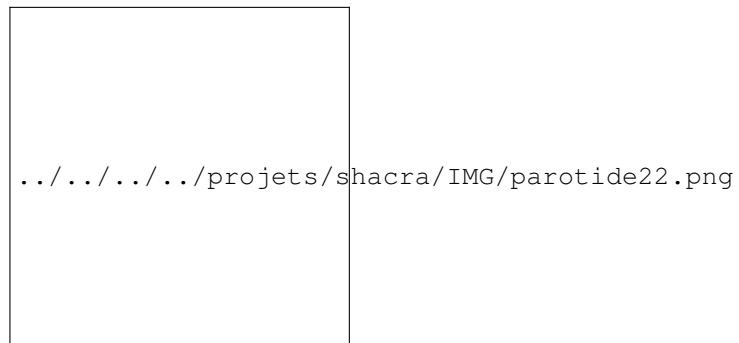


Figure 7. Screenshot of our radiotherapy planning tool.

organs. This information would be useful for both precise diagnoses and planning of surgical procedure. Involving numerical tools for diagnoses and surgery planing becomes increasingly interesting for physicians in clinical uses. The advantages of numerical models are not only in visualization, but also in quantitative measurements on a group of organs, such as their shapes and their relative movements. The processing pipeline includes patient data retrieval, image analysis, patient-specific modeling and biomechanical simulation. Our work consists in proposing new methods and algorithms for modeling the 3D anatomy of specific patients based on image data. This model should be compatible with the requirements of a biomechanical simulation. Moreover, we aim at developing new image processing tools for analyzing 2D dynamic MRI (to assess the mobilities of the pelvic system by extracting certain mechanical indicators from images) and for comparison with simulations.

Registration between geometric models and images remains a major challenge in these applications. We proposed a new model-to-image registration approach which was developed and tested for segmentation of organs in 2D images and for tracking the motion of pelvic organs from 2D dynamic MRI. Thanks to this technique, evaluation of the level of shear strain that is encountered by the fascias (connective tissues between organs) during the motion became possible. This tool could help in early diagnostic of prolapse. In the next step, our objective is to extend this method for adapting it to 3D reconstruction (with 3D geometric models and 3D MR images) and for the comparison of 3D simulations with deformable images.

## 6.6. Dynamic Deformations Simulated at Different Frequencies

The dynamic response of deformable bodies varies significantly in dependence on mechanical properties of the objects: while the dynamics of a stiff and light object (e. g. wire or needle) involves high-frequency phenomena such as vibrations, much lower frequencies are sufficient for capturing dynamic response of an object composed of a soft tissue. Yet, when simulating mechanical interactions between soft and stiff deformable models, a single time-step is usually employed to compute the time integration of dynamics of both objects. However, this can be a serious issue when haptic rendering of complex scenes composed of various bodies is considered. In this work, we present a novel method allowing for dynamic simulation of a scene composed of colliding objects modelled at different frequencies: typically, the dynamics of soft objects are calculated at frequency about 50 Hz, while the dynamics of stiff object is modeled at 1 kHz, being directly connected to the computation of haptic force feedback. The collision response is performed at both low and high frequencies employing data structures which describe the actual constraints and are shared between the high and low frequency loops. During the simulation, the realistic behaviour of the objects according to the mechanical principles (such as non-interpenetration and action-reaction principle) is guaranteed. We have shown several scenarios involving different bodies in interaction, demonstrating the benefits of the proposed method. This research has been published at IROS 2013.

## **6.7. Simulation of Lipofilling Reconstructive Surgery**

We have developed a method to simulate the outcome of reconstructive facial surgery based on fat-filling. Facial anatomy is complex: the fat is constrained between layers of tissues which behave as walls along the face; in addition, connective tissues that are present between these different layers also influence the fat-filling procedure. To simulate the end result, we have proposed a method which couples a 2.5D Eulerian fluid model for the fat and a finite element model for the soft tissues. The two models are coupled using the computation of the mechanical compliance matrix. We had two contributions: a solver for fluids which couples properties of solid tissues and fluid pressure, and an application of this solver to fat-filling surgery procedure simulation. This research has been published at MICCAI 2013.

## **6.8. Real-time simulation of contact and cutting of heterogeneous soft-tissues**

We have developed a new numerical method for interactive (real-time) simulations, which considerably improves the accuracy of the response of heterogeneous soft-tissue models undergoing contact, cutting and other topological changes. It provides an integrated methodology able to deal both with the ill-conditioning issues associated with material heterogeneities, contact boundary conditions which are one of the main sources of inaccuracies, and cutting which is one of the most challenging issues in interactive simulations. Our approach is based on an implicit time integration of a non-linear finite element model. To enable real-time computations, we propose a new preconditioning technique, based on an asynchronous update at low frequency. The preconditioner is not only used to improve the computation of the deformation of the tissues, but also to simulate the contact response of homogeneous and heterogeneous bodies with the same accuracy. We also address the problem of cutting the heterogeneous structures and propose a method to update the preconditioner according to the topological modifications. Finally, we have applied our approach to three challenging demonstrators: i) a simulation of cataract surgery ii) a simulation of laparoscopic hepatectomy iii) a brain tumor surgery. This research was done in collaboration with the University of Cardiff and has been published in the journal *Media* this year.

## **6.9. Control of Elastic Soft Robots**

In this work, we present a new method for the control of soft robots with elastic behavior, piloted by several actuators. The central contribution of this work is the use of the Finite Element Method (FEM), computed in real-time, in the control algorithm. The FEM based simulation computes the nonlinear deformations of the robots at interactive rates. The model is completed by Lagrange multipliers at the actuation zones and at the end-effector position. A reduced compliance matrix is built in order to deal with the necessary inversion of the model. Then, an iterative algorithm uses this compliance matrix to find the contribution of the actuators (force and/or position) that will deform the structure so that the terminal end of the robot follows a given position. Additional constraints, like rigid or deformable obstacles, or the internal characteristics of the actuators are integrated in the control algorithm. We illustrate our method using simulated examples of both serial and parallel structures and we validate it on a real 3D soft robot made of silicone.

## ADAM Project-Team

# 6. New Results

## 6.1. Self-Adaptive Software Systems

**Participants:** Russel Nzekwa, Romain Rouvoy [correspondant], Lionel Seinturier.

The design of self-adaptive and autonomic software systems raises many challenges. In his PhD thesis, Russel Nzekwa [12] proposes a new result with the CORONA framework that enables to build flexible autonomic systems. CORONA relies on an architectural description language which reifies the structure of the control system architecture. CORONA enables the flexible integration of non-functional-properties during the design of autonomic systems. It also provides tools for checking conflicts in the architecture of autonomic systems. Finally, the traceability between the design and the runtime implementation is carried out through the code generation of skeletons from architectural descriptions of control systems. The work on CORONA goes toward the long term objective of setting up an integrated design and programming solution for self-adaptive systems, where feedback control loops play the central role as first class elements.

## 6.2. Energy Management in Software Systems

**Participants:** Rémi Druilhe, Laurence Duchien, Lionel Seinturier [correspondant].

Energy management and saving is a concern that spans the entire domain of information and communication technologies and sciences. Recently it has been recognized that to improve its efficiency, energy has to be managed, not only at the hardware level, but also at the level of software systems, especially in distributed environments. In his PhD thesis, Rémi Druilhe [11] proposes a new result with the HOMENAP system for networked digital home environments. This work is the result of a collaboration with Orange Labs. HOMENAP takes into account three main properties: heterogeneity, dynamicity and quality of service. HOMENAP proposes an autonomic decision-making system to deal with the placement of digital services on networked devices. Based on the observation of relevant events, the system takes the decision to modify the distribution of digital services on devices in order to preserve a defined tradeoff between energy efficiency and quality of service. HOMENAP participates to the long term objective of dealing with energy as a main steering factor for self-optimizing software systems.

## 6.3. Automated Software Repair

**Participant:** Martin Monperrus [correspondant].

Automated software repair aims at assisting developers in order to improve the quality of software systems, for example by recommending some repair actions to fix bugs. In [15], we present some major results in this direction by mining fix transactions of existing software repositories. From the empirical study of 14 software repositories containing 89,993 versioning transactions, we show that we can learn a probability distribution of repair actions. We show that certain distributions over repair actions can result in an infinite time (in average) to find a repair shape while other fine-tuned distributions enable to find a repair shape in hundreds of repair attempts. We now aim at going beyond this empirical study and theoretical analysis by exploring how to use this learned knowledge for new software systems.

## FUN Project-Team

# 5. New Results

## 5.1. Routing in FUN

**Participants:** Thierry Delot, Tony Ducrocq, Nicolas Gouvy, Nathalie Mitton, Enrico Natalizio, David Simplot-Ryl, Tahiry Razafindralambo, Dimitrios Zormpas.

Wireless sensor and actuator/robot networks need some routing mechanisms to ensure that data travel the network to the sink with some guarantees. The FUN research group has investigated different geographic routing paradigms. Georouting assumes that every node is aware of its location, the one of its neighbors and of the destination(s).

In this context, we first propose the first  $k$ -anycasting georouting protocol, ie in which a node wishes to send a message to  $k$  sinks in the network [13]. Then, we tried to relax some of the assumptions. For instance in [12], we introduce HECTOR which is the first position based routing protocol which relies on virtual positions, is energy-aware and guarantees the data delivery.

In [46], [21], we assume that only a part of nodes is aware of its position and proposes a hybrid approach between position-based greedy approach and traditional on-demand routing. Indeed, geographic routing protocols show good properties for WSNs. They are stateless, local and scalable. However they require that each node of the network is aware of its own position. While it may be possible to equip each node with GPS receiver, even if it is costly, there are some issues and receiving a usable GPS signal may be difficult in some situations. For these reasons, we propose a geographic routing algorithm, called HGA, able to take advantages of position informations of nodes when available but also able to continue the routing in a more traditional way if position information is not available. We show with simulations that our algorithm offers an alternative solution to classical routing algorithm (non-geographic) and offers better performances for network with a density above 25 and more than 5% of nodes are aware of their position. [46] analyses the impact of nodes topology on network performances. We show that different topologies can lead to a difference of up to 25% on delivery ratio and average route length and more than 100% on overall cost of transmissions.

In [24], [25], [26], [3], we consider that nodes are able to move by themselves and we try to take advantage of this feature to improve the network performance. In sensor networks, there is often more than one sensor which reports an event to the sink in WSN. In existing solutions, this leads to oscillation of nodes which belong to different routes and their premature death. Experiments show that the need of a routing path merge solution is high. As a response, [24], [25] introduce the first routing protocol which locates and uses paths crossing to adapt the topology to the network traffic in a fully localized way while still optimizing energy efficiency. Furthermore the protocol makes the intersection to move away from the destination, getting closer to the sources, allowing higher data aggregation and energy saving. Our approach outperforms existing solutions and extends network lifetime up to 37%.

Using nodes location, position-based routing protocols generally apply a greedy routing that makes a sensor forward data to route to one of its neighbors in the forwarding direction of the destination. If this greedy step fails, the routing protocol triggers a recovery mechanism. Such recovery mechanisms are mainly based on graph planarization and face traversal or on a tree construction. Nevertheless real-world network planarization is very difficult due to the dynamic nature of wireless links and trees are not so robust in such dynamic environments. Recovery steps generally provoke huge energy overhead with possibly long inefficient paths. In [26], we propose to take advantage of the introduction of controlled mobility to reduce the triggering of a recovery process. We propose Greedy Routing Recovery (GRR) routing protocol. GRR enhances greedy routing energy efficiency as it adapts network topology to the network activity. Furthermore GRR uses controlled mobility to relocate nodes in order to restore greedy and reduce energy consuming recovery step triggering. Simulations demonstrate that GRR successfully bypasses topology holes in more than 72% of network topologies avoiding calling to expensive recovery steps and reducing energy consumption while preserving network connectivity.

[31] relaxes the assumption that nodes are aware of their neighbors and considers that dynamic energy sources could be available. It introduces MEGAN (Mobility assisted Energy efficient Georouting in energy harvesting Actuator and sensor Networks), a beacon-less protocol that uses controlled mobility, and takes account of the energy consumption and the energy harvesting to select next hop. MEGAN aims at prolonging the overall network lifetime rather than reducing the energy consumption over a single path. When node  $s$  needs to send a message to the sink  $d$ , it first computes the ideal position of the forwarder node based on available and needed energy, and then broadcasts this data. Every node within the transmission range of  $s$  in the forward direction toward  $d$  will start a backoff timer. The backoff time is based on its available energy and on its distance from the ideal position. The first node whose backoff timer goes off is the forwarder node. This node informs its neighborhood and then moves toward the ideal position. If, on its route, it finds a good spot for energy harvesting, it will actually stop its movement and forward the original message by using MEGAN, which will run on all the intermediate nodes until the destination is reached. Simulations show that MEGAN reduces energy consumption up to 50% compared to algorithms where mobility and harvesting capabilities are not exploited.

Additionally, according to a wide range of studies, (Informatics Technologies) IT should become a key facilitator in establishing primary education, reducing mortality and supporting commercial initiatives in Least Developed Countries (LDCs). The main barrier to the development of IT services in these regions is not only the lack of communication facilities, but also the lack of consistent information systems, security procedures, economic and legal support, as well as political commitment. In [18], we propose the vision of an infrastructureless data platform well suited for the development of innovative IT services in LDCs. We propose a participatory approach, where each individual implements a small subset of a complete information system thanks to highly secure, portable and low-cost personal devices as well as opportunistic networking, without the need of any form of infrastructure. [18] reviews the technical challenges that are specific to this approach.

## 5.2. Self-organization

**Participants:** Tony Ducrocq, Nathalie Mitton, David Simplot-Ryl, Isabelle Simplot-Ryl.

Self-organization encompasses several mechanisms. This year, the FUN research group contributes to some of them such as neighbor discovery, localization, clustering and topology control in FUN.

### 5.2.1. Neighbor discovery

HELLO protocol or neighborhood discovery is essential in wireless ad hoc networks. It makes the rules for nodes to claim their existence/aliveness. In the presence of node mobility, no x optimal HELLO frequency and optimal transmission range exist to maintain accurate neighborhood tables while reducing the energy consumption and bandwidth occupation. Thus a Turnover based Frequency and transmission Power Adaptation algorithm (TFPA) is presented in [27]. The method enables nodes in mobile networks to dynamically adjust both their HELLO frequency and transmission range depending on the relative speed. In TFPA, each node monitors its neighborhood table to count new neighbors and calculate the turnover ratio. The relationship between relative speed and turnover ratio is formulated and optimal transmission range is derived according to battery consumption model to minimize the overall transmission energy. By taking advantage of the theoretical analysis, the HELLO frequency is adapted dynamically in conjunction with the transmission range to maintain accurate neighborhood table and to allow important energy savings. The algorithm is simulated and compared to other state-of-the-art algorithms. The experimental results demonstrate that the TFPA algorithm obtains high neighborhood accuracy with low HELLO frequency (at least 11% average reduction) and with the lowest energy consumption. Besides, the TFPA algorithm does not require any additional GPS-like device to estimate the relative speed for each node, hence the hardware cost is reduced.

### 5.2.2. Topology control

Topology control is a tool for self-organizing wireless networks locally. It allows a node to consider only a subset of links/neighbors in order to later reduce computing and memory complexity. Topology control in wireless sensor networks is an important issue for scalability and energy efficiency. It is often based on graph reduction performed through the use of Gabriel Graph or Relative Neighborhood Graph. This graph reduction is usually based on geometric values.

In [11], we propose a radically new family of geometric graphs, i.e., Hypocomb, Reduced Hypocomb and Local Hypocomb for topology control. The first two are extracted from a complete graph; the last is extracted from a Unit Disk Graph (UDG). We analytically study their properties including connectivity, planarity and degree bound. All these graphs are connected (provided the original graph is connected) planar. Hypocomb has unbounded degree while Reduced Hypocomb and Local Hypocomb have maximum degree 6 and 8, respectively. To the best of our knowledge, Local Hypocomb is the first strictly-localized, degree-bounded planar graph computed using merely 1-hop neighbor position information. We present a construction algorithm for these graphs and analyze its time complexity. Hypocomb family graphs are promising for wireless ad hoc networking. We report our numerical results on their average degree and their impact on FACE [49] routing. We discuss their potential applications and some open problems.

### 5.2.3. Clustering

Clustering in wireless sensor networks is an efficient way to structure and organize the network. It aims to identify a subset of nodes within the network and bind it a leader (i.e. cluster-head). This latter becomes in charge of specific additional tasks like gathering data from all nodes in its cluster and sending them by using a longer range communication to a sink.

As a consequence, a cluster-head exhausts its battery more quickly than regular nodes. In [8], [22], [1], we present BLAC, a novel Battery-Level Aware Clustering family of schemes. BLAC considers the battery-level combined with another metric to elect the cluster-head. It comes in four variants. The cluster-head role is taken alternately by each node to balance energy consumption. Due to the local nature of the algorithms, keeping the network stable is easier. BLAC aims to maximize the time with all nodes alive to satisfy application requirements. Simulation results show that BLAC improves the full network lifetime 3-time more than traditional clustering schemes by balancing energy consumption over nodes and still delivering high data percentage.

On another approach, [34] considers the Slepian-Wolf coding based data aggregation problem and the corresponding dependable clustering problem in WSN. A dependable Slepian-Wolf coding based clustering (DSWC) algorithm is proposed to provide dependable clustering against cluster-head failures. The proposed D-SWC algorithm attempts to elect a primary cluster head and a backup cluster head for each cluster member during clustering so that once a failure occurs to the primary cluster head the cluster members within the failed cluster can promptly switchover to the backup cluster head and thus recover the connectivity of the failed cluster to the data sink without waiting for the next-round clustering to be performed. Simulation results show that the DSWC algorithm can effectively increase the amount of data transmitted to the data sink as compared with an existing nondependable clustering algorithm for Slepian-Wolf coding based data aggregation in WSNs.

## 5.3. Controlled mobility

**Participants:** Milan Erdelj, Valeria Loscri, Kalypso Magklara, Karen Miranda, Enrico Natalizio, Jean Razafimandimby Anjalalaina, Tahiry Razafindralambo, David Simplot-Ryl, Dimitrios Zormpas.

Controlled mobility [5] is a new paradigm that leads to a set of great new challenges.

### 5.3.1. Target coverage

One of the main operations in wireless sensor networks is the surveillance of a set of events (targets) that occur in the field. In practice, a node monitors an event accurately when it is located closer to it, while the opposite happens when the node is moving away from the target. This detection accuracy can be represented by a probabilistic distribution. Since the network nodes are usually randomly deployed, some of the events are monitored by a few nodes and others by many nodes. In applications where there is a need of a full coverage and of a minimum allowed detection accuracy, a single node may not be able to sufficiently cover an event by itself. In this case, two or more nodes are needed to collaborate and to cover a single target. Moreover, all the nodes must be connected with a base station that collects the monitoring data.

In [15], we describe the problem of the minimum sampling quality, where an event must be sufficiently detected by the maximum possible amount of time. Since the probability of detecting a single target using randomly deployed static nodes is quite low, we present a localized algorithm based on mobile nodes. Our algorithm sacrifices a part of the energy of the nodes by moving them to a new location in order to satisfy the desired detection accuracy. It divides the monitoring process in rounds to extend the network lifetime, while it ensures connectivity with the base station. Furthermore, since the network lifetime is strongly related to the number of rounds, we propose two redeployment schemes that enhance the performance of our approach by balancing the number of sensors between densely covered areas and areas that are poorly covered. Finally, our evaluation results show an over 10 times improvement on the network lifetime compared to the case where the sensors are static. Our approaches, also, outperform a virtual forces algorithm when connectivity with the base station is required. The redeployment schemes present a good balance between network lifetime and convergence time.

[47], [28] assume that these targets to cover are dynamic. We assume that no knowledge about either event position or duration is given a priori. Nonetheless, the events need to be monitored and covered thanks to mobile wireless sensors. Thus, mobile sensors have to discover the events and move towards a new Zone of Interest (ZoI) when the previous monitored event is over. An efficient, distributed and localized solution of this problem would be immediately exploitable by several applications domains, such as environmental, civil, etc. We propose two novel approaches to deal with dynamic event coverage. The first one is a modified version of the PSO, where particles (mobile sensors, nodes or devices in the following) update their velocity by using only local information coming from their neighbors. In practice, the velocity update is performed by considering neighbors' sensed events. Our distributed version of PSO is integrated with a distributed version of the Virtual Force Algorithm (VFA). Virtual Force approach has the ability to "position" nodes with no overlap, by using attractive and repulsive forces based on the distance between nodes. The other proposed algorithm is a distributed implementation of the VFA by itself. Both techniques are able to reach high levels of coverage and show a satisfying reactivity when the ZoI changes. This output parameter is measured as the capability for the sensors to "follow" a sequence of events happening in different ZoIs. The effectiveness of our techniques is shown through a series of simulations and comparisons with the classical centralized VFA.

On another approach consists in using flying drone to cover this set of targets. [39] focuses on the energy efficiency problem where camera equipped flying drones are able to detect and follow mobile events that happen on the ground. We give a mathematical formulation of the problem of minimizing the total energy consumption of a fleet of drones when coverage of all events is required. Due to the extremely high complexity of the binary optimization problem, the optimum solution cannot be obtained even for small instances. On the contrary, we present LAS, a localized solution for the aforementioned problem which takes into account the ability of the drones to fly at lower altitudes in order to conserve energy. We simulate LAS and we compare its performance to a centralized algorithm and to an approach that uses static drones to cover all the terrain. Our findings show that LAS performs similar to the centralized algorithm, while it outperforms the static approach by up to 150% in terms of consumed energy. Finally, the simulation results show that LAS is very sustainable in presence of communication errors.

### 5.3.2. Multiple Point of Interest coverage

The coverage of Points of Interest (PoI) is a classical requirement in mobile wireless sensor applications. Optimizing the sensors self-deployment over a PoI while maintaining the connectivity between the sensors and the base station is thus a fundamental issue.

The problems of multiple PoI discovery, coverage and data report are still solved separately and there are no works that combine the aforementioned problems into a single deployment scheme. In [9], [2], we present a novel approach for mobile sensor deployment, where we combine multiple PoI discovery and coverage with network connectivity preservation in order to capture the dynamics of the monitored area. Furthermore, we derive analytical expressions for circular movement parameters and examine the performance of our approach through extensive simulation campaigns.

[10] addresses the problem of autonomous deployment of mobile sensors that need to cover a predefined PoI with a connectivity constraint. In our algorithm, each sensor moves toward a PoI but has also to maintain the connectivity with a subset of its neighboring sensors that are part of the Relative Neighborhood Graph (RNG). The Relative Neighborhood Graph reduction is chosen so that global connectivity can be provided locally. Our deployment scheme minimizes the number of sensors used for connectivity thus increasing the number of monitoring sensors. Analytical results, simulation results and practical implementation are provided to show the efficiency of our algorithm.

### 5.3.3. Robot cooperation

The concept of autonomous mobile agents gets a lot of attention in the domain of WSN or wireless sensor and actuator networks (WSAN). Multiple robots that coordinate or cooperate with other sensors, robots or human operator, allow the WSN/WSAN to perform tasks that are far beyond the scope of single robot unit. In [23], we describe the robot middleware architecture that allows networked multi-robot control and data acquisition in the context of wireless sensor networks. Furthermore, we present three examples of robot network deployment and illustrate the proposed architecture usability: the robotic network deployment with the goal of covering the Point of Interest, adaptable multi-hop video transmission scenario, and the case of obtaining the energy consumption during the deployment.

### 5.3.4. Substitution networks

A substitution network [4] is a temporary network that will be deployed to support a base network in trouble and help it to provide the best service.

WSN are widely deployed nowadays on a large variety of applications. The major goal of a WSN is to collect information about a set of phenomena. Such process is non trivial since batteries' life is limited and thus wireless transmissions as well as computing operations must be minimized. A common task in WSNs is to estimate the sensed data and to spread the estimated samples over the network. Thus, time series estimation mechanisms are vital on this type of processes so as to reduce data transmission. In [30], we assume a single-hop clustering mechanism in which sensor nodes are grouped into clusters and communicate with a sink through a single hop. We propose a couple of autoregressive mechanisms to predict local sensed samples in order to reduce wireless data communication. We compare our proposal with a model called EEE that has been previously proposed in the literature. We prove the efficiency of our algorithms with real samples publicly available and show that they outperform the EEE mechanism.

In [32], we propose an algorithm to efficiently (re-)deploy the wireless mobile routers of a substitution network by considering the energy consumption, a fast deployment scheme and a mix of the network metric. We consider a scenario where we have two routers in a fixed network and where connectivity must be restored between those two routers with a wireless mobile router. The main objective of the wireless mobile router is to increase the communication performance such as the throughput by acting as relay node between the two routers of the fixed network. We present a fast, adaptive and localized approach which takes into account different network metrics such as Received Signal Strength (RSS), Round-Trip Time (RTT) and the Transmission Rate, between the wireless mobile router and the two routers of the fixed network. Our method ameliorates the performance of our previous approach from the literature by shortening the deployment time, increasing the throughput, and consuming less energy in some specific cases.

## 5.4. Security

**Participants:** Nathalie Mitton, Enrico Natalizio.

[19] deals with the energy efficient issue of cryptographic mechanisms used for secure communication between devices in wireless sensor networks. Since these devices are mainly targeted for low power consumption appliances, there is an effort for optimization of any aspects needed for regular sensor operation. On a basis of utilization of hardware cryptographic accelerators integrated in microcontrollers, this article provides the comparison between software and hardware solutions. Proposed work examines the problems and solutions for implementation of security algorithms for WSN devices. Because the speed of hardware accelerator should



be much higher than the software implementation, there are examination tests of energy consumption and validation of performance of this feature. Main contribution of the article is real testbed evaluation of the time latency and energy requirements needed for securing the communication. In addition, global evaluation for all important network communication parameters like throughput, delay and delivery ratio are also provided.

The Internet of Things (IoT) will enable objects to become active participants of everyday activities. Introducing objects into the control processes of complex systems makes IoT security very difficult to address. Indeed, the Internet of Things is a complex paradigm in which people interact with the technological ecosystem based on smart objects through complex processes. The interactions of these four IoT components, person, intelligent object, technological ecosystem, and process highlight a systemic and cognitive dimension within security of the IoT. The interaction of people with the technological ecosystem requires the protection of their privacy. Similarly, their interaction with control processes requires the guarantee of their safety. Processes must ensure their reliability and realize the objectives for which they are designed. We believe that the move towards a greater autonomy for objects will bring the security of technologies and processes and the privacy of individuals into sharper focus. Furthermore, in parallel with the increasing autonomy of objects to perceive and act on the environment, IoT security should move towards a greater autonomy in perceiving threats and reacting to attacks, based on a cognitive and systemic approach. In [33], we will analyze the role of each of the mentioned actors in IoT security and their relationships, in order to highlight the research challenges and present our approach to these issues based on a holistic vision of IoT security.

## 5.5. RFID

**Participants:** Ibrahim Amadou, Nathalie Mitton.

Mitigating reader-to-reader collisions is one of the principal challenges in a large-scale dynamic RFID system with a number of readers deployed in order to maximize the system performance (i.e., throughput, fairness and latency). In prior works, contention-based and activity scheduling medium access control (MAC) protocols are commonly used approaches to reduce such problems. Existing protocols typically perform worse in a large-scale RFID dynamic system and require more additional components or are based on unrealistic assumptions. So far, many research efforts have been made to improve the performance or the reliability of Carrier Sense Multiple Access (CSMA) techniques for Mobile Ad-Hoc Networks (MANETs) by using an adaptive Backoff schemes. In [17], we look at these well known solutions that proved their efficiency in high congestion wireless networks. We evaluate the performance and characterize these solutions when they are used to reserve the wireless channel through broadcasting message for reader-to-tag communication. Based on the application requirements, we study their capacity to mitigate collisions, the channel access latency, the average number of successful requests sent per reader and the fairness index in the context of RFID networks.

## 5.6. Data collection and aggregation

**Participant:** Nathalie Mitton.

Named Data Networking (NDN) is a new promising paradigm for content retrieval and distribution in the future Internet. NDN communication is driven by data consumers that broadcast Interest packets to require named contents. The requests are forwarded towards the source(s) by directly using content names (instead of IP addresses), while in-network caching is used to improve delivery performance. NDN shows many similarities with data-centric models defined for wireless sensor networks (WSNs), e.g., directed diffusion. In addition, NDN defines a new complete communication framework with innovative naming and security schemes and novel routing and transport strategies. This clearly opens new perspectives in the design and development of sensor networks, which can benefit of the NDN framework to better support different kinds of applications and services. In [16], we explore the potentialities of NDN applied to WSNs and propose enhanced delivery strategies inspired by the directed diffusion scheme to be deployed in the NDN framework. Performance of a plain NDN scheme and of our enhanced solution is evaluated through the ndnSIM simulator. Achieved results confirm the viability of a NDN-like approach over WSNs and the better efficiency and effectiveness of the proposed solution compared to a plain NDN.

[38] considers the Slepian-Wolf coding based energy-minimization rate allocation problem in a WSN and propose a distributed rate allocation algorithm to solve the problem. The proposed distributed algorithm is based on an existing centralized rate allocation algorithm which has a high computational complexity. To reduce the computational complexity of the centralized algorithm and make the rate allocation performable in a distributed manner, we make necessary modifications to the centralized algorithm by reducing the number of sets in calculating the average energy consumption cost and limiting the number of conditional nodes that a set can use. Simulation results show that the proposed distributed algorithm can significantly reduce the computational time when compared with the existing centralized algorithm at the cost of the overall energy consumption for data transmission and the total amount of data transmitted in the network.

## 5.7. VANET

**Participant:** Nathalie Mitton.

Routing is a critical issue in vehicular ad hoc networks (VANETs). This paper considers the routing issue in both vehicle to vehicle (V2V) and vehicle to infrastructure (V2I) communications in VANETs, and proposes a Moving dirEction and DestinAtion Location based routing (MEDAL) algorithm for supporting V2V and V2I communications. MEDAL [36] takes advantage of both the moving directions of vehicles and the destination location to select a neighbor vehicle as the next hop for forwarding data. Unlike most existing routing algorithms, it only uses a HELLO message to obtain or update routing information without using other control messages, which largely reduces the number of control messages used in routing. Simulation results show that MEDAL can significantly improve the packet delivery ratio of the network as compared with the well-known Ad hoc On-demand Distance Vector Routing (AODV) algorithm.

## 5.8. Industrial Applications

**Participants:** Milan Erdelj, Nathalie Mitton, Enrico Natalizio.

The collaborative nature of industrial wireless sensor networks (IWSNs) brings several advantages over traditional wired industrial monitoring and control systems, including self-organization, rapid deployment, flexibility, and inherent intelligent processing. In this regard, IWSNs play a vital role in creating more reliable, efficient, and productive industrial systems, thus improving companies' competitiveness in the marketplace. Industrial Wireless Sensor Networks: Applications, Protocols, and Standards [42] examines the current state of the art in industrial wireless sensor networks and outlines future directions for research.

## RMOD Project-Team

## 6. New Results

### 6.1. Tools for understanding applications: IDEs and Visualization

**Performance Evolution Blueprint: Understanding the Impact of Software Evolution on Performance.** Understanding the root of a performance drop or improvement requires analyzing different program executions at a fine grain level. Such an analysis involves dedicated profiling and representation techniques. JProfiler and YourKit, two recognized code profilers fail, on both providing adequate metrics and visual representations, conveying a false sense of the performance variation root. We propose performance evolution blueprint, a visual support to precisely compare multiple software executions. Our blueprint is offered by Rizel, a code profiler to efficiently explore performance of a set of benchmarks against multiple software revisions. [31]

**Seamless Composition and Reuse of Customizable User Interfaces with Spec** Implementing UIs is often a tedious task. To address this, UI Builders have been proposed to support the description of widgets, their location, and their logic. A missing aspect of UI Builders is however the ability to reuse and compose widget logic. In our experience, this leads to a significant amount of duplication in UI code. To address this issue, we built Spec: a UIBuilder for Pharo with a focus on reuse. With Spec, widget properties are defined declaratively and attached to specific classes known as composable classes. A composable class defines its own widget description as well as the model-widget bridge and widget interaction logic. Spec enables seamless reuse of widgets, its use in Pharo 2.0 has cut in half the amount of lines of code of six of its tools, mostly through reuse. This shows that Spec meets its goals of allowing reuse and composition of widget logic. [17]

**Pragmatic Visualizations for Roassal: a Florilegium** Software analysis and in particular reverse engineering often involves a large amount of structured data. This data should be presented in a meaningful form so that it can be used to improve software artefacts. The software analysis community has produced numerous visual tools to help understand different software elements. However, most of the visualization techniques, when applied to software elements, produce results that are difficult to interpret and comprehend. We present five graph layouts that are both expressive for polymetric views and agnostic to the visualization engine. These layouts favor spatial space reduction while emphasizing on clarity. Our layouts have been implemented in the Roassal visualization engine and are available under the MIT License. [23]

### 6.2. Software Quality: Bugs and Debuggers

**BugMaps-Granger: A Tool for Causality Analysis between Source Code Metrics and Bugs.** Despite the increasing number of bug analysis tools for exploring bugs in software systems, there are no tools supporting the investigation of causality relationships between internal quality metrics and bugs. We propose an extension of the BugMaps tool called BugMaps-Granger that allows the analysis of source code properties that caused bugs. For this purpose, we relied on Granger Causality Test to evaluate whether past changes to a given time series of source code metrics can be used to forecast changes in a time series of defects. Our tool extracts source code versions from version control platforms, generates source code metrics and defects time series, computes Granger, and provides interactive visualizations for causal analysis of bugs. We also provide a case study in order to evaluate the tool. [22]

#### **Mining Architectural Patterns Using Association Rules**

Software systems usually follow many programming rules prescribed in an architectural model. However, developers frequently violate these rules, introducing architectural drifts in the source code. We present a data mining approach for architecture conformance based on a combination of static and historical software analysis. For this purpose, the proposed approach relies on data mining techniques to extract structural and historical architectural patterns. In addition, we propose a methodology that uses the extracted patterns to detect both absences and divergences in source-code based architectures. We applied the proposed approach in an industrial strength system. As a result we detected 137 architectural violations, with an overall precision of 41.02%. [27]

### Heuristics for Discovering Architectural Violations

Software architecture conformance is a key software quality control activity that aims to reveal the progressive gap normally observed between concrete and planned software architecture. We present ArchLint, a lightweight approach for architecture conformance based on a combination of static and historical source code analysis. For this purpose, ArchLint relies on four heuristics for detecting both absences and divergences in source code based architectures. We applied ArchLint in an industrial-strength system and as a result we detected 119 architectural violations, with an overall precision of 46.7% and a recall of 96.2%, for divergences. We also evaluated ArchLint with four open-source systems, used in an independent study on reflexion models. In this second study, ArchLint achieved precision results ranging from 57.1% to 89.4%. [26]

## 6.3. Software Quality: History and Changes

**Representing Code History with Development Environment Events.** Modern development environments handle information about the intent of the programmer: for example, they use abstract syntax trees for providing high-level code manipulation such as refactorings; nevertheless, they do not keep track of this information in a way that would simplify code sharing and change understanding. In most Smalltalk systems, source code modifications are immediately registered in a transaction log often called a ChangeSet. Such mechanism has proven reliability, but it has several limitations. We analyse such limitations and describe scenarios and requirements for tracking fine-grained code history with a semantic representation. We want to enrich code sharing with extra information from the IDE, which will help understanding the intention of the changes and let a new generation of tools act in consequence. [24]

**Mining System Specific Rules from Change Patterns** A significant percentage of warnings reported by tools to detect coding standard violations are false positives. Thus, there are some works dedicated to provide better rules by mining them from source code history, analyzing bug-fixes or changes between system releases. However, software evolves over time, and during development not only bugs are fixed, but also features are added, and code is refactored. In such cases, changes must be consistently applied in source code to avoid maintenance problems. We propose to extract system specific rules by mining systematic changes over source code history, i.e., not just from bug-fixes or system releases, to ensure that changes are consistently applied over source code. We focus on structural changes done to support API modification or evolution with the goal of providing better rules to developers. Also, rules are mined from predefined rule patterns that ensure their quality. In order to assess the precision of such specific rules to detect real violations, we compare them with generic rules provided by tools to detect coding standard violations on four real world systems covering two programming languages. The results show that specific rules are more precise in identifying real violations in source code than generic ones, and thus can complement them. [25]

## 6.4. Reconciling Dynamic Languages and Isolation

**Virtual Smalltalk Images: Model and Applications.** Reflective architectures are a powerful solution for code browsing, debugging or in-language process handling. However, these reflective architectures show some limitations in edge cases of self-modification and self-monitoring. Modifying the modifier process or monitoring the monitor process in a reflective system alters the system itself, leading to the impossibility to perform some of those tasks properly. We analyze the problems of reflective architectures in the context of image based object-oriented languages and solve them by providing a first-class representation of an image: a virtualized image. We present Oz, our virtual image solution. In Oz, a virtual image is represented by an object space. Through an object space, an image can manipulate the internal structure and control the execution of other images. An Oz object space allows one to introspect and modify execution information such as processes, contexts, existing classes and objects. We show how Oz solves the edge cases of reflective architectures by adding a third participant, and thus, removing the self modification and self-observation constraints. [30]

**Bootstrapping Reflective Systems: The Case of Pharo.** Bootstrapping is a technique commonly known by its usage in language definition by the introduction of a compiler written in the same language it compiles. This process is important to understand and modify the definition of a given language using the same language,

taking benefit of the abstractions and expression power it provides. A bootstrap, then, supports the evolution of a language. However, the infrastructure of reflective systems like Smalltalk includes, in addition to a compiler, an environment with several self-references. A reflective system bootstrap should consider all its infrastructural components. We propose a definition of bootstrap for object-oriented reflective systems, we describe the architecture and components it should contain and we analyze the challenges it has to overcome. Finally, we present a reference bootstrap process for a reflective system and Hazelnut, its implementation for bootstrapping the Pharo Smalltalk-inspired system. [15]

**Object Graph Isolation with Proxies** More and more software systems are now made of multiple collaborating third-party components. Enabling fine-grained control over the communication between components becomes a major requirement. While software isolation has been studied for a long time in operating systems (OS), most programming languages lack support for isolation. In this context we explore the notion of proxy. A proxy is a surrogate for another object that controls access to this object. We are particularly interested in generic proxy implementations based on language-level reflection. We present an analysis that shows how these reflective proxies can propagate a security policy thanks to the transitive wrapping mechanism. We present a prototype implementation that supports transitive wrapping and allows a fine-grained control over an isolated object graph. [33]

## 6.5. Dynamic Languages: Compilers

**Towards a flexible Pharo Compiler** The Pharo Smalltalk-inspired language and environment started its development with a codebase that can be traced back to the original Smalltalk-80 release from 1983. Over the last years, Pharo has been used as the basis of many research projects. Often these experiments needed changes related to the compiler infrastructure. However, they did not use the existing compiler and instead implemented their own experimental solutions. This shows that despite being an impressive achievement considering its age of over 35 years, the compiler infrastructure needs to be improved. We identify three problems: (i) The architecture is not reusable, (ii) compiler can not be parametrized and (iii) the mapping between source code and bytecode is overly complex. Solving these problems will not only help researchers to develop new language features, but also the enhanced power of the infrastructure allows many tools and frameworks to be built that are important even for day-to-day development, such as debuggers and code transformation tools. [20]

**Gradual Typing for Smalltalk** Being able to combine static and dynamic typing within the same language has clear benefits in order to support the evolution of prototypes or scripts into mature robust programs. While being an emblematic dynamic object-oriented language, Smalltalk is lagging behind in this regard. We report on the design, implementation and application of Gradualtalk, a gradually-typed Smalltalk meant to enable incremental typing of existing programs. The main design goal of the type system is to support the features of the Smalltalk language, like metaclasses and blocks, live programming, and to accommodate the programming idioms used in practice. We studied a number of existing projects in order to determine the features to include in the type system. As a result, Gradualtalk is a practical approach to gradual types in Smalltalk, with a novel blend of type system features that accommodate most programming idioms. [13]

## LINKS Team

# 5. New Results

## 5.1. Querying Heterogeneous Linked Data

**Participants:** Guillaume Bagan, Iovka Boneva, Angela Bonifati, Pierre Bourhis, Radu Ciucanu, Tom Sebastian, Slawomir Staworko, Sophie Tison.

Staworko, Ciucanu and Boneva presented a new class of schemas for unordered XML trees, which are based on unordered regular expressions, also called multiplicity schemas. They show that many static analysis problems become feasible when removing disjunctions there [6].

Ciucanu and Staworko [8] investigated the case of unordered XML, where the relative order among siblings is ignored, and focused on the problem of learning schemas from examples given by the user. They considered disjunctive multiplicity schemas (DMS) and their restrictions, disjunction-free multiplicity schemas (MS). For both DMS and MS, they prove the learnable cases.

Regular path queries in graphs have found much recent interest in the context of SPARQL queries for linked open data in the RDF format. Bagan, Bonifati and Groz (former PhD student of Mostrare, now PostDoc at Tel-Aviv University) have obtained a precise characterization of those regular path queries that can be answered with polynomial data complexity [5] leading to a trichotomy (AC0, NL-complete, or else NP-complete). Thereby, they have solved an open question (raised by W. Martens in PODS'12).

XPath query evaluation over compressed trees has been studied in [12]. They focused on a fragment of XPath, which is the downward, navigational XPath and presented precise bounds on the time complexity of XPath query execution over grammar-compressed trees. In particular, they focused on counting the nodes selected by an XPath expression, extracting and materializing their pre-order numbers and serializing the obtained subtrees.

In [2], Groz, Staworko, Caron, Roos and Tison studied query rewriting with views when the classes used to define queries and views are Regular XPath and MSO. Next, they investigated problems of static analysis of security access specifications (SAS) by introducing the novel class of interval-bounded SAS and they defined three different manners to compare views (i.e. queries), with a security point of view. Finally, they provided a systematic study of the complexity for deciding these three comparisons.

## 5.2. Managing Dynamic Linked Data

**Participants:** Angela Bonifati, Denis Debarbieux, Joachim Niehren, Tom Sebastian.

Bonifati, Goodfellow (former PhD student at the University of Strathclyde, UK), Manolescu and Sileo (former PhD student at the University of Basilicata, Italy, directed by Bonifati) studied XML view maintenance in the presence of updates [1]. Their approach relies on algebraic operators for propagating source updates to the target XML view, e.g. in a typical scenario of GAV (global-as-view) schema mappings. Their algebraic approach is set-oriented as opposed to tuple-oriented methods presented in the literature. Moreover, it leverages structural identifiers and structural join algorithms. As such, it proved to be more efficient than existing methods for updating materialized XML views.

Debarbieux, Gauwin (former PhD student in the team, now Assistant Professor at the University of Bordeaux), Niehren, Sebastian and Zergaoui (CEO at Innovimax) focused on using early nested word automata in order to approximate earliest query answering algorithms for nested word automata in a highly efficient manner [9]. This approximation can be made tight in practice for automata obtained from XPath expressions. An XPath streaming algorithm based on early nested word automata has been implemented in the FXP tool. FXP outperforms most previous tools in efficiency, while covering more queries of the XPathMark benchmark.

### 5.3. Linking Data Graphs

**Participants:** Angela Bonifati, Radu Ciucanu, Joachim Niehren, Aurélien Lemay, Grégoire Laurence, Antoine Ndione, Slawomir Staworko.

In [7], Bonifati, Ciucanu and Staworko investigate the problem of inferring arbitrary n-ary join predicates across two relations via user interactions. The relations can be found on the Web, thus they lack integrity constraints. In such a scenario, the user is asked to label as positive or negative a few tuples depending on whether she would like them in the join result or not. Deciding whether the remaining tuples are uninformative, i.e. do not allow to infer the query goal, can be done in polynomial time.

The PhD thesis of Ndione focuses on probabilistic algorithms to decide approximate membership of words in a language by using property testing. In [3], Ndione, Lemay and Niehren presented an algorithm that tests the membership modulo the edit distance. Their algorithm runs in polynomial time, as opposed to other property testing algorithms, leveraging the Hamming distance or the edit distance with moves, that are exponential.

In [11], Laurence, Lemay, Niehren, Staworko and Tommasi (project leader of the Magnet team) studied the problem of learning sequential top-down tree-to-word transducers (STWs). They present a Myhill-Nerode characterization of the corresponding class of sequential tree-to-word transformations (STW). Next, they investigate what learning of STWs means, identify fundamental obstacles, and propose a learning model with abstain. Finally, they present a polynomial learning algorithm.

In [4], Niehren, Champavère (former PhD student in the team), Gilleron and Lemay addressed the problem of learnability of regular queries in unranked trees. The idea is that tree pruning strategies and the schemas (DTD in the specific case) can guide the learning process and lead to a class of queries that are learnable according to those. The obtained learning algorithm adds pruning heuristics to the traditional learning algorithm based on tree automata and exploiting positive and negative examples.

## MAGNET Team

# 6. New Results

## 6.1. Probabilistic models for large graph

We have developed new approaches for the statistical analysis of large-scale undirected graphs. The main insight is to exploit the spectral decomposition of subgraph samples, and in particular their Fiedler eigenvalues, as basic features for density estimation and probabilistic inference. Our contributions are twofold. First, we develop a conditional random graph model for learning to predict links in information networks (such as scientific coauthorship and email communication). Second, we propose to apply the resulting model to graph generation and link prediction. This work is published in the *Journal of Machine Learning Research*, the top journal in the field of machine learning.

## 6.2. Learning in hypergraphs

In this work, we focus on the problem of learning from several sources of heterogeneous data represented as input graphs that encode different relations over the same set of nodes. Our goal is to merge those input graphs by embedding them into an Euclidean space related to the commute time distance in the original graphs. Our algorithm outputs a combined kernel that can be used for different graph learning tasks. This work has been published in [5].

The approach designed in that paper has raised a new definition of undirected hypergraphs with bipartite hyperedges. A bipartite hyperedge is a pair of disjoint sets of nodes in which every node is associated with a weight. A bipartite hyperedge can be viewed as a relation between two teams of nodes in which every node has a weighted contribution to its team. Undirected hypergraphs generalize over undirected graphs. Consistently with the case of graphs, we have studied the hypergraph spectral framework. We have defined the notions of hypergraph gradient, hypergraph Laplacian, and hypergraph kernel as the Moore-Penrose pseudoinverse of a hypergraph Laplacian. Therefore, smooth labeling of (teams of) nodes and hypergraph regularization methods can be performed. Contrary to the graph case, we show that the class of hypergraph Laplacians is closed by the pseudoinverse operation (thus it is also the class of hypergraphs kernels), and is closed by convex linear combination. Closure properties allow us to define (hyper)graph combinations and operations while keeping a hypergraph interpretation of the result. We exhibit a subclass of signed graphs that can be associated with hypergraphs in a constructive way. A hypergraph and its associated signed graph have the same Laplacian. This property allows us to define a distance between nodes in undirected hypergraphs as well as in the subclass of signed graphs. The distance coincides with the usual definition of commute-time distance when the equivalent signed graph turns out to be a graph. We claim that undirected hypergraphs open the way to solve new learning tasks and model new problems based on set similarity or dominance. We are currently exploring applications for modeling games between teams and for graph summarization. This work [8] has been submitted to *Journal of Machine Learning Research*.

## 6.3. Natural Language Processing

In [7] and [3], we develop a new algorithm for drastically improving a pairwise coreference classification system. Specifically, this algorithm works by learning the best partition over mention type pairs by training different pairwise coreference models for each pair type. In effect, our algorithm finds the optimal feature space (from a base feature set and set of types) for separating coreferential mention pairs, but it remains tractable by exploiting the structure of the hierarchies built from the pair types. In [6], we propose a new approach for the automatic identification of so-called implicit discourse relations. Our system combines hand-labeled examples and automatically annotated examples (based on explicit relations) using different methods inspired by work on domain adaptation. Our system is evaluated empirically and yields important performance gains compared to only using hand-labeled data. This paper has received the best paper award at the *TALN 2013* conference, the national NLP conference.



## 6.4. Query Induction

We have proposed a new algorithm for query learning that combines schema-guided pruning heuristics with the traditional learning algorithm for tree automata from positive and negative examples. We show that this algorithm is justified by a formal learning model, and that for stable queries it performs very well in practice of XML information extraction. This work [1] has also been published in *JMLR*.

## 6.5. Learning Transducers

We have pursued the work on learning finite state tree-to-word transducers. Tree-to-word transformations are ubiquitous in computer science. They are the core of many computation paradigms from the evaluation of abstract syntactic trees to modern programming languages *XSLT*. We have extended the results obtained last year on the study of a class of sequential top-down tree-to-word transducers, called *STWs*. Transducers in *STWs* are capable of: concatenation in the output, producing arbitrary context-free languages, deleting inner nodes, and verifying that the input tree belongs to the domain even when deleting parts of it. These features are often missing in tree-to-tree transducers, and for instance, make *STWs* incomparable with the class of top-down tree-to-tree transducers. The class of *STWs* has several interesting properties, in particular we proposed in 2011 a normal for *STWs*.

In [4], we present a Myhill-Nerode characterization of the corresponding class of sequential tree-to-word transformations. Next, we investigate what learning of *STWs* means, identify fundamental obstacles, and propose a learning model with abstain. Finally, we present a polynomial learning algorithm.

## MINT Project-Team

# 6. New Results

## 6.1. Human limits in small unidirectional mouse movements

**Participants:** Jonathan Aceituno [correspondant], Géry Casiez, Nicolas Roussel.

Computer mouse sensors keep increasing in resolution. The smallest displacement they can detect gets smaller, but little is known on our ability to control such small movements. Small target acquisition has been previously tackled, but the findings do not apply to the problem of finding the useful resolution of a user with a mouse, which corresponds to the smallest displacement (s)he can reliably produce with that device. In [16], we detail this definition and provide an associated experimental protocol to measure the useful resolution. We then report on the results of a study suggesting that high-end mice are not likely to be used to their full potential. We further comment on the different strategies used by participants to achieve best performance, and derive implications for user interfaces.

## 6.2. Small, Medium, or Large? Estimating the User-Perceived Scale of Stroke Gestures

In [27], we show that a large consensus exists among users in the way they articulate stroke gestures at various scales (i.e., small, medium, and large) and formulate a simple rule that estimates the user-intended scale of input gestures with 87% accuracy. Our estimator can enhance current gestural interfaces by leveraging scale as a natural parameter for gesture input, reflective of user perception (i.e., no training required). Gesture scale can simplify gesture set design, improve gesture- to-function mappings, and reduce the need for users to learn and for recognizers to discriminate unnecessary symbols.

## 6.3. Métamorphe : a shape changing keyboard

Métamorphe is a keyboard with mobile keys [21]. Whether keys are pressed or released, they can be at their usual height, or raised. This mechanism allows both to provide haptic feedback to ease eyes-free interaction, and to access the side of the keys. The sides of the keys can be pushed, like the top of the keys. Therefore each key can be mapped to several actions. For instance this could be useful for command selection.

## 6.4. Designing Intuitive Multi-touch 3D Navigation Techniques

**Participants:** Géry Casiez, Damien Marchal [correspondant], Nicolas Roussel, Clement Moerman.

Multi-touch displays have become commonplace over recent years. Numerous applications take advantage of this to support interactions that build on users' knowledge and correspond to daily practices within the real world. 3D applications are also becoming more common on these platforms, but the multi-touch techniques for 3D operations often lag behind 2D ones in terms of intuitiveness and ease of use. Intuitive navigation techniques are particularly needed to make multi-touch 3D applications more useful, and systematic approaches are direly needed to inform their design: existing techniques are still too often designed in ad-hoc ways. In [25], we propose a methodology based on cognitive principles to address this problem. The methodology combines standard user-centered design practices with optical flow analysis to determine the mappings between navigation controls and multi-touch input. It was used to design the navigation technique of a specific application for our industrial partner Idées3Com. The resulting technique proved to be more efficient and preferred by users when compared to existing ones, which provides a first validation of the approach.



*Figure 1. a) Métamorphe concept: the user presses the control key, keys corresponding to hotkeys rise b) key mounted on a solenoid, with force sensors on the sides c) press on the top of the key d) press on the right of the key e) press on the left of the key.*

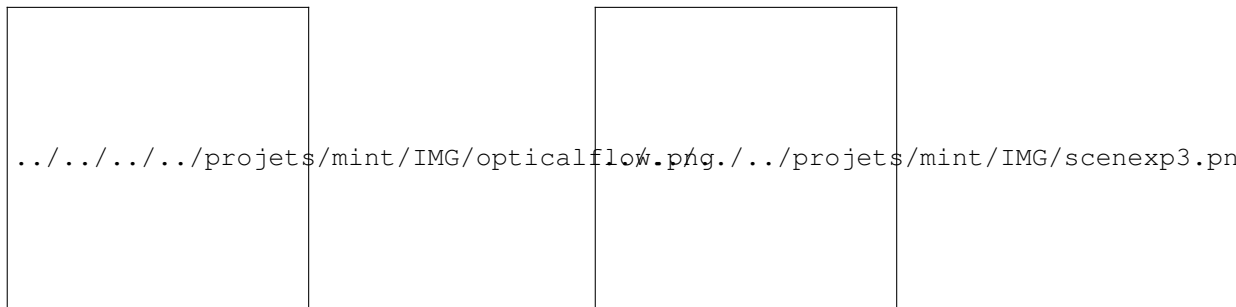


Figure 2. (Left) The optical flow for camera movements are used to design the shape of the interaction gestures. (Right) The evaluation scenario used to compare several state of the art navigation techniques.

## 6.5. Mockup Builder: 3D modeling on and above the surface

Mockup Builder [11] is a semi-immersive environment for conceptual design which allows virtual mockups to be created using gestures. Our goal is to provide familiar ways for people to conceive, create and manipulate three-dimensional shapes. To this end, we developed on-and-above-the-surface interaction techniques based on asymmetric bimanual interaction for creating and editing 3D models in a stereoscopic environment. Our approach combines both hand and finger tracking in the space on and above a multi-touch surface. This combination brings forth an alternative design environment where users can seamlessly switch between interacting on the surface or above it to leverage the benefit of both interaction spaces. A formal user evaluation conducted with experienced users shows very promising avenues for further work towards providing an alternative to current user interfaces for modeling.

## 6.6. Towards Many Gestures to One Command: A User Study for Tabletops

**Participants:** Yosra Rekik, Laurent Grisoni [correspondant], Nicolas Roussel.

This work has been accepted as a long paper at Interact 2013. Multi-touch gestures are often thought by application designers for a one-to-one mapping between gestures and commands, which does not take into account the high variability of user gestures for actions in the physical world; it can also be a limitation that leads to very simplistic interaction choices. Our motivation is to make a step toward many-to-one mappings between user gestures and commands, by understanding user gestures variability for multi-touch systems; for doing so, we set up a user study in which we target symbolic gestures on tabletops. From a first phase study we provide qualitative analysis of user gesture variability; we derive this analysis into a taxonomy of user gestures, that is discussed and compared to other existing taxonomies. We introduce the notion of atomic movement; such elementary atomic movements may be combined throughout time (either sequentially or in parallel), to structure user gesture. A second phase study is then performed with specific class of gesture-drawn symbols; from this phase, and according to the provided taxonomy, we evaluate user gesture variability with a fine grain quantitative analysis. Our findings indicate that users equally use one or two hands, also that more than half of gestures are achieved using parallel or sequential combination of atomic movements. We also show how user gestures distribute over different movement categories, and correlate to the number of fingers and hands engaged in interaction. Finally, we discuss implications of this work to interaction design, practical consequences on gesture recognition, and potential applications.

## 6.7. Sub-space gestures: elements of design for mid-air interaction with distant displays

**Participants:** Hanae Rateau, Laurent Grisoni [correspondant], Bruno de Araujo.

(Research report, accepted to publication in a modified version to IUI 2014). Multi-touch gestures are often thought by application designers for a one-to-one mapping between gestures and commands, which does not take into account the high variability of user gestures for actions in the physical world; it can also be a limitation that leads to very simplistic interaction choices. Our motivation is to make a step toward many-to-one mappings between user gestures and commands, by understanding user gestures variability for multi-touch systems; for doing so, we set up a user study in which we target symbolic gestures on tabletops. From a first phase study we provide qualitative analysis of user gesture variability; we derive this analysis into a taxonomy of user gestures, that is discussed and compared to other existing taxonomies. We introduce the notion of atomic movement; such elementary atomic movements may be combined throughout time (either sequentially or in parallel), to structure user gesture. A second phase study is then performed with specific class of gesture-drawn symbols; from this phase, and according to the provided taxonomy, we evaluate user gesture variability with a fine grain quantitative analysis. Our findings indicate that users equally use one or two hands, also that more than half of gestures are achieved using parallel or sequential combination of atomic movements. We also show how user gestures distribute over different movement categories, and correlate to the number of fingers and hands engaged in interaction. Finally, we discuss implications of this work to interaction design, practical consequences on gesture recognition, and potential applications.

## 6.8. Merging two tactile stimulation principles: Electrovibration and Squeeze film effect

**Participants:** Michel Amberg, Frédéric Giraud, Clément Nadal, Betty Semail [correspondant].

Electrovibration and squeeze film effect can modify the perception a user has of a flat surface, with opposite action. In fact, electrovibration increases the friction of the finger on the surface, while the squeeze film reduces it. These two stimulation principles are compatible, and in this work [23], we wanted to merge them in a tactile stimulator, in order to enhance the control of the lateral force. Our approach was to identify the effect of each tactile stimulation, and we proposed its modelling: the dynamic of the mechanical response of the fingerpulp has to be taken into account between the programmed stimulus and the resulting lateral force. We have shown also that the two techniques may be used simultaneously accounting to a few precautions. From the first experimental trials, the conclusion here is that the squeeze film effect is able to reduce tangential forces generated by the electrostatic forces, by going on acting on the friction coefficient.

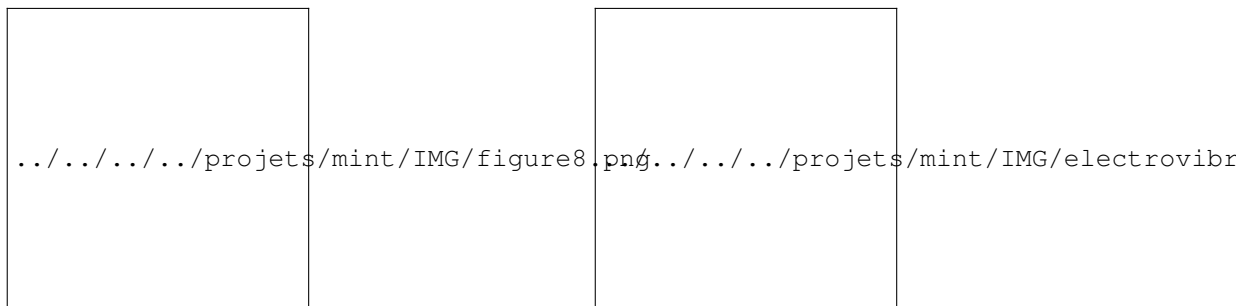


Figure 3. (Left) The experimental test bench to measure the forces produced during the stimulation. (Right) The tactile stimulator merging the two stimulation principles.