# Activity Report 2016

# Section New Results

<span style="color:red">**BONSAI Project-Team**</span>

# 7. New Results

## 7.1. Approximate pattern matching

The problem of measuring the similarity between two strings arises in many areas of sequence analysis. A common metric for it is the *Levenshtein distance*. This distance is defined as the smallest number of substitutions, insertions, and deletions of symbols required to transform one of the words into the other. We have investigated the basic problem of the size of the neighborhood of a given pattern $P$: count how many strings are within a bounded distance of a fixed reference string. There has been no efficient algorithm for calculating it so far. We have proposed a dynamic programming algorithm that scales linearly with the size of the pattern $P$. For that, we have introduced a new variant of the universal Levenshtein automaton, that is interesting by itself and that can have many other applications in text algorithms [31].

We have also addressed the related problem of approximate pattern matching: Given a text $T$ and a pattern $P$, find all locations in $T$ that differ by at most $k$ errors (in the sense of the Levenshtein distance) from $P$. We have proposed a new kind of seeds (the 01*0 seeds) that combines exact parts and parts with a fixed number of errors, and that are specifically well-suited for short DNA motifs with high error-rate. We have demonstrated the applicability of those seeds on two main case studies : pattern matching on a genomic scale with a Burrows-Wheeler transform, and multi-pattern matching with indexation of the set of patterns [30].

## 7.2. Parallel algorithm for de Bruijn graph compaction

Constructing a *de Bruijn graph* is an important step in the analysis of NGS data. This data structure is used in several applications, such as *de novo* assembly, variant detection, and transcriptome quantification. However, the representation of this graph often consumes prohibitive amounts of memory for large datasets. An operation, called compaction, enables to represent the graph more efficiently. However, so far, there was no algorithm for compacting the graph quickly and in low memory.

Along with colleagues at Inria Rennes and at Penn State University, we introduced a parallel algorithm and an implementation, BCALM 2, for constructing directly a compacted de Bruijn graph given a set of reads. Our results show that this algorithm enables to construct the graph for very large datasets, such as the spruce and pine genomes, in reasonable time and memory on a single machine. This represents a performance improvement of two orders of magnitude compared to previously available methods. BCALM 2 is open-source and was published at ISMB 2016 [20].

## 7.3. Range minimum query

The *range minimum query* problem consists in finding the minimum value inside any queried range of a preprocessed integer sequence. Several methods exist to compute the minimum in constant time, using almost the theoretical minimal amount of space. Those methods consist in splitting the problem in several subproblems and precomputing the solutions for them.

With Alice Héliou (AMIB Inria team, Saclay), Martine Léonard and Laurent Mouchard (LITIS, Rouen), we designed a new method, which is worse in terms of time complexity [24]. Our solution relies on a totally different concept as previous ones: We only store the values that are local minima. This approach is therefore simple and can, on specific inputs, require much less memory than the general theoretical minimal bound. Moreover the simplicity of the method can be easily adapted to allow updates in the original integer sequence.

## 7.4. Coding isoform structures

Our researches on gene isoform structures started in 2014 with the CG-Alcode Associated Team and in collaboration with Anne Bergeron from the LACIM (Montréal, Canada). We aimed at defining better definitions of isoform orthology at the coding level, which are based on the preservation of all the exon junctions in two orthologous isoforms. This estimation is achieved at the gene level, where sequence homology is detected for both exons and their flanking intronic splice sites [19]. The approach largely outperforms competing programs in terms of precision and recall. Using the successive releases of the ten years old CCDS database, we show that the discovery rate of orthologous isoforms between human and mouse is growing continuously and that it displays no sign of completion.

## 7.5. Nonribosomal peptides

We were invited to contribute in a volume of "Methods in Molecular Biology" by authoring a chapter focusing on NRPS biosynthesis. This chapter [32] was about the use of the Norine platform (developed by the team) and other bioinformatics tools for the analysis of nonribosomal peptide synthetases and their products. We invited our collaborator from Denmark, Tilmann Weber, to complete this chapter with the introduction of his tool, antiSMASH.

We annotated 48 genomes of *Burkholderia* species using our annotation protocol, that starts from a genome sequence and goes to the predicted nonribosomal peptides. We have predicted 161 gene clusters producing nonribosomal peptides, leading to the synthesis of not only already known peptides, but also new ones [22] with potential applications in biocontrol.

A new version of the Norine interface is now available. The form to query the annotations is now flexible and dynamic. The user can build his own query to search for annotations in several fields combined by boolean operators. Moreover, the database structure has been modified to allow, among others, a hierarchical representation of the NRPS taxonomy. Finally, the MyNorine tool has been enhanced and updated to take into account these changes and the description page of the peptides has been reorganized.

## 7.6. High-throughput V(D)J repertoire analysis

Researches on high-throughput V(D)J repertoire analysis started in the group in 2012. We have developed Vidjil, a web platform dedicated to the analysis of lymphocyte populations. Starting from DNA sequences, uploaded by the user, Vidjil identifies and quantifies lymphocyte populations and provides an interactive visualization [21].

In 2016, with our colleagues at Lille hospital, we published two articles in haematological journals to detail our method for the diagnosis [23] and for the follow-up [28] of the acute lymphoblastic leukemia using high-throughput sequencing. Our results also show what those new techniques, together with bioinformatics software, bring in a routine practice. Being a full platform with metadata storage, Vidjil is used on a regular basis by about 20 laboratories around the world. In France, the majority of diagnosis samples from acute lymphoblastic leukemia patients are now analyzed using Vidjil.

## 7.7. Assembly of the giraffe genome and the gorilla Y-chromosome

We collaborated with two labs from the Pennsylvania State Institute (Cavener Lab and Makova Lab) for practical analysis of DNA sequencing data. The first collaboration led to the publication of the giraffe genome in Nature Communication [18]. In this article our contribution was to provide the first draft-quality whole-genome sequences of the giraffe and the okapi. The second collaboration was about assembling the Y-chromosome of the gorilla using a novel sequencing strategy as well as novel computational tools. This work was published in Genome Research [29].

<span style="color:red">**DEFROST Team**</span>

# 7. New Results

## 7.1. Cochlear Implants

Publication at MICCAI 2016 (Medical Image Computing and Computer Assisted Intervention conference): **Numerical Simulation of Cochlear-Implant Surgery: Towards Patient-Specific Planning**, *Olivier Goury, Yann Nguyen, Renato Torres, Jeremie Dequidt, Christian Duriez*. **Abstract.** During Cochlear Implant Surgery, the right placement of the implant and the minimization of the surgical trauma to the inner ear are an important issue with recurrent fails. In this study, we reproduced, using simulation, the mechanical insertion of the implant during the surgery. This simulation allows to have a better understanding of the failing cases: excessive contact force, buckling of the implant inside and outside the cochlea. Moreover, using a patient-specific geometric model of the cochlea in the simulation, we show that the insertion angle is a clinical parameter that has an influence on the forces endured by both the cochlea walls and the basilar membrane, and hence to post-operative trauma. The paper presents the mechanical models used for the implant, for the basilar membrane and the boundary conditions (contact, friction, insertion etc...) and discuss the obtained results in the perspective of using the simulation for planning and robotization of the implant insertion.

<span style="color:red">https://hal.archives-ouvertes.fr/hal-01370185</span>
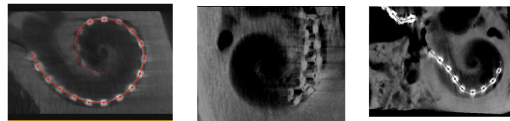


*Figure 3. Three outcomes of implant insertion (from left to right): successful insertion; failed insertion (Folding tip); incomplete insertion*
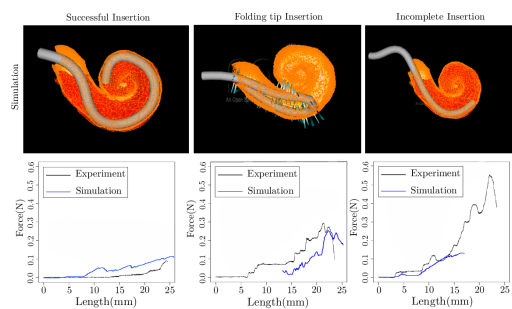


*Figure 4. Reproduction of real insertion cases with the simulation*

## 7.2. Physics based model of soft-robots

Book chapter in Soft Robotics: Trends, Applications and Challenges, Springer, 2016 **Soft Robot Modeling, Simulation and Control in Real-Time**, *Christian Duriez and Thor Bieze*. https://hal.inria.fr/hal-01410293. We were asked to write a chapter in this book on Soft Robotics. Our chapter presents new real-time and physics-based modeling methods dedicated to deformable soft robots. In our approach, continuum mechanics provides the partial derivative equations that govern the deformations, and Finite Element Method (FEM) is used to compute numerical solutions adapted to the robot. A formulation based on Lagrange Multipliers is presented to model the behavior of the actuators as well as the contact with the environment. Direct and inverse kinematic models are also obtained for real-time control. Some experiments and numerical results are presented.

## 7.3. Kinematic Modeling and control of soft robots

Publication at IROS 2016 : **Kinematic Modeling and Observer Based Control of Soft Robot using Real-Time Finite Element Method**, *Zhongkai Zhang, Jeremie Dequidt, Alexandre Kruszewski, Frederick Largilliere, Christian Duriez.*. **Abstract.**This paper aims at providing a novel approach to modeling and controlling soft robots. Based on real-time Finite Element Method (FEM), we obtain a globally defined discrete-time kinematic model in the workspace of soft robots. From the kinematic equations, we deduce the soft-robot Jacobian matrix and discuss the conditions to avoid singular configurations. Then, we propose a novel observer based control methodology where the observer is built by Finite Element Model in this paper to deal with the control problem of soft robots. A closed-loop controller for position control of soft robot is designed based on the discrete-time model with feedback signal being extracted by means of visual servoing. Finally, experimental results on a parallel soft robot show the efficiency and performance of our proposed controller. https://hal.inria.fr/hal-01370347

## 7.4. Stiffness rendering

Publication at IROS 2016 : **Stiffness rendering on soft tangible devices controlled through inverse FEM simulation**, *Frederick Largilliere, Eulalie Coevoet, Mario Sanz-Lopez, Laurent Grisoni, Christian Duriez*. **Abstract.**Haptic rendering of soft bodies is essential in medical simulations of procedures such as surgery or palpation. The most commonly used approach is to recreate the sense of touch using a specific design and control of a robotic arm. In this paper, we propose a new approach, based on soft-robotics technology. We create a tangible deformable device that allows users to " touch " soft tissues and perceive mechanical material properties, in a realistic manner. The device is able to dynamically provide user touch with different stiffness perceptions, thanks to actuators placed at the boundaries. We introduce a control algorithm, based on inverse Finite Element Analysis, which controls the actuators in order to recreate a desired stiffness that corresponds to the contact with soft tissues in the virtual environment. The approach uses antagonistic actuation principle to create a wide range of stiffness. We validate our algorithm and demonstrate the method using prototypes based on simple mechanisms. https://hal.inria.fr/hal-01386787

## 7.5. Framework for soft robot simulation

Publication at SIMPAR 2016 : **Framework for online simulation of soft robots with optimization-based inverse model**, *C. Duriez, E. Coevoet, F. Largilliere, T. Morales-Bieze, Z. Zhang, M. Sanz-Lopez, B. Carrez, D. Marchal, O. Goury, J. Dequidt*. **Abstract.**Soft robotics is an emerging field of robotics which requires computer-aided tools to simulate soft robots and provide models for their control. Until now, no unified software framework covering the different aspects exists. In this paper, we present such a framework from its theoretical foundations up to its implementation on top of Sofa, an open-source framework for deformable online simulation. The framework relies on continuum mechanics for modeling the robotic parts and boundary conditions like actuators or contacts using a unified representation based on Lagrange multipliers. It enables the digital robot to be simulated in its environment using a direct model. The model can also be inverted online using an optimization-based method which allows to control the physical robots in the task space. To demonstrate the effectiveness of the approach, we present various soft robots scenarios including ones where the robot is interacting with its environment. https://hal.inria.fr/hal-01425349

## 7.6. Closed-loop control

Closed-loop control based on dynamic models of soft robots. Model-order reduction provides a system of achievable size to apply traditional control science techniques. During the internship of Maxime Thieffry, we obtain the first results in that direction that will be extended during a PhD thesis.

<p style="text-align:center"><span style="color:red">**DOLPHIN Project-Team**</span></p>

# 7. New Results

## 7.1. Optimization under uncertainty

Participants: El-Ghazali Talbi, Raca Todosijevic, Oumayma Bahri (external collaborators: Nahla BenAmor - Univ. Tunis, Tunisia, J. Puente, C. R. Vela, I. Gonzalez-Rodriguez - Univ. Oviedo Spain)

At the problem level, the sources of uncertainty are due to many factors such as the environment parameters of the model, the decision variables and the objective functions. Examples of such uncertainties can be the demand and travel times in vehicle routing problems, the execution time in scheduling problems, the wind or solar production in energy power systems, the price of resources in manufacturing, and the mechanical properties of a structure. Then, we need precise and efficient modeling and resolution approaches which are robust and non-sensitive to those uncertainties. The appeal of optimization under uncertainty is that its performance results remain relatively unchanged when exposed to uncertain data.

We have considered the fuzzy job shop, a job shop scheduling problem with uncertain processing times modelled as triangular fuzzy numbers. While the usual approaches to solving this problem involve adapting existing metaheuristics to the fuzzy setting, we have proposed instead to follow the framework of simheuristics from stochastic optimisation. More precisely, we integrate the simulation of possible realisations of the fuzzy problem with a genetic algorithm that solves the deterministic job shop. We test the resulting method, simGA, on a testbed of 23 benchmark instances and obtain results that suggest that this is a promising approach to solving problems with uncertainty by means of metaheuristics [38].

## 7.2. Indicator-based Multiobjective Optimization

Participants: Bilel Derbel, Arnaud Liefooghe (external collaborators: Matthieu Basseur, Adrien Goëffon, Univ. Angers, France)

A large spectrum of quality indicators has been proposed so far to assess the performance of discrete Pareto set approximations in multiobjective optimization. Such indicators assign, to any solution set, a real-value reflecting a given aspect of approximation quality. This is an important issue in multiobjective optimization, not only to compare the performance and assets of different approximate algorithms, but also to improve their internal selection mechanisms. In [37], we adopt a statistical analysis to experimentally investigate by how much a selection of state-of-the-art quality indicators agree with each other for a wide range of Pareto set approximations from well-known two- and three-objective continuous benchmark functions. More particularly, we measure the correlation between the ranking of low-, medium-, and high-quality limited-size approximation sets with respect to inverted generational distance, additive epsilon, multiplicative epsilon, R2, R3, as well as hypervolume indicator values. Since no pair of indicators obtains the same ranking of approximation sets, we confirm that they emphasize different facets of approximation quality. More importantly, our statistical analysis allows the degree of compliance between these indicators to be quantified.

Subset selection constitutes an important stage of any evolutionary multiobjective optimization algorithm when truncating the current approximation set for the next iteration. This appears to be particularly challenging when the number of solutions to be removed is large, and when the approximation set contains many mutually non-dominating solutions. In particular, indicator-based strategies have been intensively used in recent years for that purpose. However, most solutions for the indicator-based subset selection problem are based on a very simple greedy backward elimination strategy. We experiment additional heuristics that include a greedy forward selection and a greedy sequential insertion policies, a first-improvement hill-climbing local search, as well as combinations of those. We evaluate the effectiveness and the efficiency of such heuristics in order to maximize the enclosed hypervolume indicator of candidate subsets during a hypothetical evolutionary process, or as a post-processing phase. Our experimental analysis, conducted on randomly generated as well

as structured two-, three- and four-objective mutually non-dominated sets, allows us to appreciate the benefit of these approaches in terms of quality, and to highlight some practical limitations and open challenges in terms of computational resources.

## 7.3. Decomposition-based Multiobjective Optimization

Participants: Bilel Derbel, Arnaud Liefooghe (external collaborators: Hernan Aguirre and Kiyoshi Tanaka, Shinshu Univ., Japan; Qingfu Zhang, City Univ., Hong Kong)

It is generally believed that local search (LS) should be used as a basic tool in multi-objective evolutionary computation for combinatorial optimization. However, not much effort has been made to investigate how to efficiently use LS in multi-objective evolutionary computation algorithms. In [28], we study some issues in the use of cooperative scalarizing local search approaches for decomposition-based multiobjective combinatorial optimization. We propose and study multiple move strategies in the MOEA/D framework. By extensive experiments on a new set of bi-objective traveling salesman problems with tunable correlated objectives, we analyze these policies with different MOEA/D parameters. Our empirical study has shed some insights about the impact of the Ls move strategy on the anytime performance of the algorithm.

## 7.4. Learning and Adaptation for Landscape-aware Algorithm Design

Participants: Bilel Derbel, Arnaud Liefooghe (external collaborators: Hernan Aguirre, Fabio Daolio, Miyako Sagawa and Kiyoshi Tanaka, Shinshu Univ., Japan; Cyril Fonlupt, Christopher Jankee and Sébastien Verel, Univ. Littoral, France)

In [13], we attempt to understand and to contrast the impact of problem features on the performance of randomized search heuristics for black-box multi-objective combinatorial optimization problems. At first, we measure the performance of two conventional dominance-based approaches with unbounded archive on a benchmark of enumerable binary optimization problems with tunable ruggedness, objective space dimension, and objective correlation ($\rho$MNK-landscapes). Precisely, we investigate the expected runtime required by a global evolutionary optimization algorithm with an ergodic variation operator (GSEMO) and by a neighborhood-based local search heuristic (PLS), to identify a $(1 + \varepsilon)$-approximation of the Pareto set. Then, we define a number of problem features characterizing the fitness landscape, and we study their intercorrelation and their association with algorithm runtime on the benchmark instances. At last, with a mixed-effects multi-linear regression we assess the individual and joint effect of problem features on the performance of both algorithms, within and across the instance classes defined by benchmark parameters. Our analysis reveals further insights into the importance of ruggedness and multi-modality to characterize instance hardness for this family of multi-objective optimization problems and algorithms.

Designing portfolio adaptive selection strategies is a promising approach to gain in generality when tackling a given optimization problem. However, we still lack much understanding of what makes a strategy effective, even if different benchmarks have been already designed for these issues. In [35], we propose a new model based on fitness cloud allowing us to provide theoretical and empirical insights on when an on-line adaptive strategy can be beneficial to the search. In particular, we investigate the relative performance and behavior of two representative and commonly used selection strategies with respect to static (off-line) and purely random approaches, in a simple, yet sound realistic, setting of the proposed model.

In evolutionary multi-objective optimization, variation operators are crucially important to produce improving solutions, hence leading the search towards the most promising regions of the solution space. In [39], we propose to use a machine learning modeling technique, namely random forest, in order to estimate, at each iteration in the course of the search process, the importance of decision variables with respect to convergence to the Pareto front. Accordingly, we are able to propose an adaptive mechanism guiding the recombination step with the aim of stressing the convergence of the so-obtained offspring. By conducting an experimental analysis using some of the WFG and DTLZ benchmark test problems, we are able to elicit the behavior of the proposed approach, and to demonstrate the benefits of incorporating machine learning techniques in order to design new efficient adaptive variation mechanisms.

## 7.5. Feature Selection using Tabu Search with Learning Memory: Learning Tabu Search

Participants: C. Dhaenens, L. Jourdan, M-E. Kessaci

Feature selection in classification can be modeled as a combinatorial optimization problem. One of the main particularities of this problem is the large amount of time that may be needed to evaluate the quality of a subset of features. We propose to solve this problem with a tabu search algorithm integrating a learning mechanism. To do so, we adapt to the feature selection problem, a learning tabu search algorithm originally designed for a railway network problem in which the evaluation of a solution is time-consuming. Experiments conducted show the benefit of using a learning mechanism to solve hard instances of the literature [hal-01370396v1].

## 7.6. MO-ParamILS: A Multi-objective Automatic Algorithm Configuration Framework

Participants: C. Dhaenens, L. Jourdan, M-E. Kessaci

Automated algorithm configuration procedures play an increasingly important role in the development and application of algorithms for a wide range of computationally challenging problems. Until very recently, these configuration procedures were limited to optimising a single performance objective, such as the running time or solution quality achieved by the algorithm being configured. However, in many applications there is more than one performance objective of interest. This gives rise to the multi-objective automatic algorithm configuration problem, which involves finding a Pareto set of configurations of a given target algorithm that characterises trade-offs between multiple performance objectives. In this work, we introduced MO-ParamILS, a multiobjective extension of the state-of-the-art single-objective algorithm configuration framework ParamILS, and demonstrated that it produces good results on several challenging bi-objective algorithm configuration scenarios compared to a base-line obtained from using a state-of-the-art single-objective algorithm configurator. [hal-01370392].

## 7.7. Parallel optimization methods revisited for multi-core and many-core (co)processors

Participants: J. Gmys and N. Melab

This contribution is a joint work with M. Mezmaz, E. Alekseeva and D. Tuyttens from University of Mons (UMONS) and T. C. Pessoa and F. H. De Carvalho Junior from Universidade Federal Do Cearà (UFC), Brazil. On the road to exascale, coprocessors are increasingly becoming key building blocks of High Performance Computing platforms. In addition to their energy efficiency, these many-core devices boost the performance of multi-core processors. During 2016, we first have revisited the design and implementation of parallel Branch-and-Bound (B&B) algorithms using the work stealing paradigm on GPU accelerators [16][40], multi-GPU systems [17], multi-core processors [15] and MIC (Xeon Phi) coprocessors [20]. The challenge is to take into account the high irregular nature of the B&B algorithm and the hardware characteristics of GPU, Xeon Phi and multi-core (co)processors. Several work stealing strategies have been investigated while addressing several issues: host-device data transfer, thread divergence and data placement on the hierarchy of memories of the GPU and vectorization on Xeon Phi. The proposed approaches have been extensively experimented considering permutation-based optimization problems (e.g. FSP). The results reported in the cited papers demonstrate the efficiency of the many-core approaches compared to their multi-core counterpart. An extension of the proposed approaches to large hybrid clusters, including multi-core and many-core (co)processors is already started in [27].

The second part of the contribution consists in proposing a new hyper-heuristic (generalized GRASP) together with its parallelization for multi-core processors [11]. A cost function based on a bounding operator (used in B&B) is integrated to GRASP for the first time. Multi-core computing is used to investigate 315 GRASP configurations. In order to improve the performance of the local search procedure used in GRASP, we have proposed in [33] an original vectorization of the cost function of the makespan of FSP on Xeon Phi coprocessors. The reported results show that speed-ups up to 4.5 can be achieved compared to a non-vectorized apprpoach.

<span style="color:red">**DREAMPAL Project-Team**</span>

# 6. New Results

## 6.1. A Language-Independent Proof System for Full Program Equivalence

Two programs are mutually equivalent if, for the same input, either they both diverge or they both terminate with the same result. Mutual equivalence is an adequate notion of equivalence for programs written in deterministic languages. It is useful in many contexts, such as capturing the correctness of program transformations within the same language, or capturing the correctness of compilers between two different languages. In [11] we introduce a language-independent proof system for mutual equivalence, which is para-metric in the operational semantics of two languages and in a state-similarity relation. The proof system is sound: if it terminates then it establishes the mutual equivalence of the programs given to it as input. We illustrate it on two programs in two different languages (an imperative one and a functional one), that both compute the Collatz sequence. The Collatz sequence is an interesting case study since it is not known wether the sequence terminates or not; nevertheless, our proof system shows that the two programs are mutually equivalent (even if we cannot establish termination or divergence of either one).

## 6.2. A Generic Framework for Symbolic Execution: a Coinductive Approach

In [12] we propose a language-independent symbolic execution framework. The approach is parameterised by a language definition, which consists of a signature for the lan-guage's syntax and execution infrastructure, a model interpreting the signature, and rewrite rules for the language's operational semantics. Then, symbolic execution amounts to computing symbolic paths using a derivative operation. We prove that the symbolic execution thus defined has the properties naturally expected from it, meaning that the feasible symbolic executions of a program and the concrete executions of the same program mutually simulate each other. We also show how a coinduction-based extension of symbolic execution can be used for the deductive verification of programs. We show how the proposed symbolic-execution approach, and the coinductive verification technique based on it, can be seamlessly implemented in language definition frameworks based on rewriting such as the K framework. A prototype implementation of our approach has been developed in K. We illustrate it on the symbolic analysis and deductive verification of nontrivial programs.

## 6.3. Circuit Merging versus Dynamic Partial Reconfiguration -The HoMade Implementation

One goal of reconfiguration is to save power and occupied resources. In [13] we compare two different kinds of reconfiguration available on field-programmable gate arrays (FPGA) and we discuss their pros and cons. The first method that we study is circuit merging. This type of reconfiguration methods consists in sharing common resources between different circuits. The second method that we explore is dynamic partial reconfiguration (DPR). It is specific to some FPGA, allowing well defined reconfigurable parts to be modified during runtime. We show that DPR, when available, has good and more predictable result in terms of occupied area. There is still a huge overhead in term of time and power consumption during the reconfiguration phase. Therefore we show that circuit merging remains an interesting solution on FPGA because it is not vendor specific and the reconfiguration time is around a clock cycle. Besides, good merging algorithms exist eventhough FPGA physical synthesis flow makes it hard to predict the real performance of the merged circuit during the optimization. We establish our comparison in the context of the HoMade processo

## 6.4. Language Definitions as Rewrite Theories

K is a formal framework for defining operational semantics of programming languages. The K-Maude compiler translates K language definitions to Maude rewrite theories. The compiler enables program execution by using the Maude rewrite engine with the compiled definitions, and program analysis by using various Maude analysis tools. K supports symbolic execution in Maude by means of an automatic transformation of language definitions. The transformed definition is called the symbolic extension of the original definnition. In [14] we investigate the theoretical relationship between K language definitions and their Maude translations, between symbolic extensions of K definitions and their Maude translations, and how the relationship between K definitions and their symbolic extensions is reflected on their respective representations in Maude. In particular, the results show how analysis performed with Maude tools can be formally lifted up to the original language definitions.

## 6.5. SCAC-Net: Reconfigurable Interconnection Network in SCAC Massively parallel SoC

Parallel communication plays a critical role in massively parallel systems, especially in distributed memory systems executing parallel programs on shared data. Therefore, integrating an interconnection network in these systems becomes essential to ensure data inter-nodes exchange. Choosing the most effective communication structure must meet certain criteria: speed, size and power consumption. Indeed, the communication phase should be as fast as possible to avoid compromising parallel computing, using small and low power consumption modules to facilitate the interconnection network extensibility in a scalable system. To meet these criteria and based on a module reuse methodology, we chose to integrate a reconfigurable SCAC-Net interconnection network to communicate data in SCAC Massively parallel SoC. In [15] we present the detailed hardware implementation and discusses the performance evaluation of the proposed reconfigurable SCAC-Net network.

## 6.6. Proving Reachability-Logic Formulas Incrementally

Reachability Logic (RL) is a formalism for defining the operational semantics of programming languages and for specifying program properties. As a program logic it can be seen as a language-independent alternative to Hoare Logics. Several verification techniques have been proposed for RL, all of which have a circular nature: the RL formula under proof can circularly be used as a hypothesis in the proof of another RL formula, or even in its own proof. This feature is essential for dealing with possibly unbounded repetitive behaviour (e.g., program loops). The downside of such approaches is that the verification of a set of RL formulas is monolithic, i.e., either all formulas in the set are proved valid, or nothing can be inferred about any of the formula's validity or invalidity. In [16] we propose a new, incremental method for proving a large class of RL formulas. The proposed method takes as input a given RL formula under proof (corresponding to a given program fragment), together with a (possibly empty) set of other valid RL formulas (e.g., already proved using our method), which specify sub-programs of the program fragment under verification. It then checks certain conditions are shown to be equivalent to the validity of the RL formula under proof. A newly proved formula can then be incrementally used in the proof of other RL formulas, corresponding to larger program fragments. The process is repeated until the whole program is proved. We illustrate our approach by verifying the nontrivial Knuth-Morris-Pratt string-matching program.

<span style="color:red">**FUN Project-Team**</span>

# 7. New Results

## 7.1. Routing

**Participants:**  Nathalie Mitton, Mouna Masmoudi.

Geographic routing is an attractive routing strategy in wireless sensor networks. It works well in dense networks, but it may suffer from the void problem. For this purpose, a recovery step is required to guarantee packet delivery. Face routing has widely been used as a recovery strategy since proved to guarantee delivery. However, it relies on a planar graph not always achievable in realistic wireless networks and may generate long paths. In [23], [12], we propose GRACO, a new geographic routing algorithm that combines a greedy forwarding and a recovery strategy based on swarm intelligence. During recovery, ant packets search for alternative paths and drop pheromone trails to guide next packets within the network. GRACO avoids holes and produces near optimal paths. Simulation results demonstrate that GRACO leads to a significant improvement of routing performance and scalability when compared to the literature algorithms.

GRACO has first been designed in the general case. We then studied its applicability to the Virtual Power Plants and their specific data packets with different priorities [23], [12]. Indeed, the Smart Grid (SG) incorporates communication networks to the conventional electricity system in order to intelligently integrate distributed energy resources (DERs) and allow for demand side management. The move to Smart grid in developing countries has to cope with great disparities of ICT infrastructures even within the same city. Besides, individual DERs are often too small to be allowed access to energy market, likewise power utilities are unable to effectively control and manage small DERs. We propose the use of affordable and scalable wireless communication technology to aggregate geographically sparse DERs into a single virtual power plant. The enrollment of prosumers in the VPP is conditional to financial performance of the plant. Thus, the VPPs are dynamic and are expected to scale up as more and more prosumers are attracted by their financial benefits. the communication network has to follow this progression and therefore to be scalable and rapidly deploy-able. We present a routing algorithm for data communication within the VPP to support centralized, decentralized or fully distributed control of the VPP's DERs.

Based on this study, we adapted GRACO so it can fit the specific cases of Smart Grid [23], [12] and more specifically to the Neighbor Area Networks (NAN) of Smart Grids, or distribution segment of the power system in the smart grid (SG). The deployment of ICT to support conventional grid will solve legacy problems that used to prevent implementation of smart services such as smart metering, demand side management or the integration of Distributed Energy Resources (DERs) within the smart grid. We demonstrate the effectiveness of GRACO in terms of scalability, peer-to-peer routing, end-to-end delay and delivery rate.

In another context, we made the observation that typical betweenness centrality metrics neglect the potential contribution of nodes that are near but not exactly on shortest paths. The idea of [35] is to give more value to these nodes. We propose a weighted betweenness centrality, a novel metric that assigns weights to nodes based on the stretch of the paths they intermediate against the shortest paths. We compare the proposed metric with the traditional and the distance-scaled betweenness metrics using four different network datasets. Results show that the weighted betweenness centrality pinpoints and promotes nodes that are underestimated by typical metrics, which can help to avoid network disconnections and better exploit multipath protocols.

## 7.2. Cloud and IoT

**Participants:**  Valeria Loscri, Nathalie Mitton, Riccardo Petrolo.

Innovative and effective solutions to the fragmentation issues in the Internet of Things (IoT) landscape have been designed and proof of concept have been implemented to show the feasibility and effectiveness of the Cloud of Things (CoT) paradigm. In other words, we have focused on the convergence of Web semantic technologies and the Cloud computing concept as key enabler of an horizontal integration of various IoT applications and platforms [21]. The heterogeneity has to be considered not only in terms of applications and platforms, but another "type of heterogeneity" that deserves to be considered and analyzed is based on different devices and their interoperability.

A feasible solution to make different and heterogeneous devices to "interoperate" is based on the exploitation of a gateway. In particular, we have considered a Gateway-as-a-Service (Gaas) in [36], where we have shown that it is an efficient and lightweight device, which can be shared between several final users. Through the container virtualization technologies, we have been able to show how several platform requirements can be met, in a context where constrained devices have been considered. This study has demonstrated the Gateway-as-a-Service (GaaS) effectiveness and its exploitability in several IoT contexts, such as smart home, buildings, farms, agriculture environments, etc.

A different and complementary, to the previous solutions, perspective of IoT paradigm is represented by the management of the huge amount of data that have to be treated in the different IoT based applications. In [45], an infer algorithm has been proposed and more specifically an Bayesian Inference Approach (BIA) with the amin objective to avoid the transmission of high spatio-temporal correlated data.

## 7.3. Resource management in FUN

**Participants:** Cristina Cano Bastidas, Valeria Loscri, Simon Duquennoy.

A standard solution for reliable low-power mesh networks was defined in IEEE802.15.4e-2012, through the new MAC layer TSCH. TSCH (Time-Slotted Channel Hopping) provides a globally synchronized network that enables scheduling and channel hopping. Our review paper [28] details the TSCH technology as well as the 6LoWPAN and 6TiSCH protocols. It gathers authors from all major open-source IoT OSes: Contiki, OpenWSN, RIOT and TinyOS. The paper presents architectural considerations when it comes to implementing portable TSCH stacks, and presents preliminary evaluation results.

TSCH networks require global synchronization. The more precise the synchronization, the more energy-efficient the network. We address the challenge of reaching micro-second time synchronization over multiple hops in TSCH networks [31], at low power. The key idea is to use two crystal oscillators, one at low-frequency for low-power timekeeping, one at high-frequency for intra-slot precision. Along with adaptive drift compensation, this method is proven effective through an experimental assessment.

Beaconing is usually employed to allow network discovery and to maintain synchronisation in mesh networking protocols, such as those defined in the IEEE 802.15.4e and IEEE 802.11s standards. Thus, avoiding persistent or consecutive collisions of beacons is crucial in order to ensure correct network operation. Beacons are also used in receiver-initiated medium access protocols to advertise that nodes are awake. Consequently, effective beacon scheduling can enable duty-cycle operation and reduce energy consumption. We propose [56] a completely decentralised and low-complexity solution based on learning techniques to schedule beacon transmissions in mesh networks. We show the algorithm converges to beacon collision-free operation almost surely in finite time and evaluate converge times in different mesh network scenarios.

In [54] we focus on new methods, architectures, and applications for the management of Cyber Physical Objects (CPOs) in the context of the Internet of Things (IoT). The book covers a wide range of topics related to CPOs, such as resource management, hardware platforms, communication and control, and control and estimation over networks. It also discusses decentralized, distributed, and cooperative optimization as well as effective discovery, management, and querying of CPOs. Other chapters outline the applications of control, real-time aspects, and software for CPOs and introduce readers to agent-oriented CPOs, communication support for CPOs, real-world deployment of CPOs, and CPOs in Complex Systems. There is a focus on the importance of application of IoT technologies for Smart Cities.

Finally, we address software security and in particular the challenge of formally verifying the source code of IoT OSes. This is the topic of the yet-to-be-started H2020 VESSEDIA project. Our preliminary study [32] demonstrated the feasibility of applying Frama-C to a memory allocation module of the Contiki OS.

## 7.4. Smart Cities

**Participants:** Nathalie Mitton, Valeria Loscri, Riccardo Petrolo.

Smart City represents one of the most promising, prominent and challenging Internet of Things (IoT) applications, but recent ICT trends suggest more and more that cities could also benefit from Cloud computing. The convergence of IoT paradigm and Cloud computing technology, can play a fundamental role for developing of highly level and organized cities form an ICT point of view, but it is of paramount importance to deal a critical analysis to identify the issues and challenges deriving from this synergy.

A novel perspective that we have considered as key factor for the realization of Future Internet is the role of the interconnected objects as active entities in the context of the networked systems [52]. With this perspective in mind, we have proposed CACHACA [43], a ranking mechanism for Sensor Networks that facilitate the discovery of services provided by each network element. Discovery functionality has been also considered in the context of VITAL project, since effective and accurate mechanisms to discover Inter-Connected Objects (ICOs) and new services represents a sine qua non condition to have effective exploration of data-sources that are appropriate for a specific business context as defined by an end-user [42] [11].

On the other hand, a Smart City is a kind of ecosystem characterized with different IoT solutions that have to cooperate and coexist and is in continuous expansion. In order to face with the integration and interoperability challenges of this ecosystem, we have considered VITAL-OS architecture that can monitor, visualize, and control all the operations of a city [44].

## 7.5. RFID

**Participants:** Nathalie Mitton, Abdoul Aziz Mbacke.

One of the devices under consideration by the FUN team is RFID. One of the main issues to widely deploy RFID reader is reader-to-reader collision. Indeed, when the electromagnetic fields of the readers overlap, a collision occurs on the tag laying in the overlapping section and cannot be read. Numerous protocols have been proposed to attempt to reduce them, but, remaining reading errors still heavily impact the performances and fairness of dense RFID deployments. In [33], [18] we introduce a new Distributed Efficient & Fair Anticollision for RFID (DEFAR) protocol. It reduces both monochannel and multichannel collisions as well as interference by a factor of almost 90% in comparison with the best state of the art protocols. The fairness of the medium access among the readers is improved to a 99% level. Such improvements are achieved applying a TDMA-based "server-less" approach and assigning different priorities to readers depending on their behavior over precedent rounds. A distributed reservation phase is organized between readers with at least one winning reader afterwards. Then, multiple reading phases occur within a single frame in order to obtain fast coverage and high throughput. The use of different reader priorities based on reading behaviors of previous frames also contributes to improve both fairness and efficiency. Simulation results show the robustness of the proposed solution in terms of different metrics such collision avoidance, fairness and coverage and in comparison with a centralized literature solution.

In order to ensure collision-free reading, a scheduling scheme is needed to read tags in the shortest possible time. We study in [37] this scheduling problem in a stationary setting and the reader minimization problem in a mobile setting. We show that the optimal schedule construction problem is NP-complete and provide an approximation algorithm that we evaluate our techniques through simulation.

## 7.6. Interferences and failures management

**Participants:** Nathalie Mitton, Viktor Toldov, Valeria Loscri, Simon Duquennoy.

In the recent years, the Machine-to-Machine (M2M) paradigm together with the integration of wireless sensors networks with the generic infrastructure via $6LoWPAN$ require the implementation of ad hoc communication protocols at the Medium Access Control layer, that do not depend on pre-existing infrastructure. Channel hopping concept has more and more gained consensus as a viable and effective solution for wireless MAC layer coordination with time-synchronized channel hopping (TSCH). In [24] we propose a decentralized multichannel MAC coordination framework (DT-SCS) leveraging the concept of *pulse-coupled oscillators* at the MAC layer. In DT-SCS, nodes randomly join a channel and are automatically spread across the available channels. The nodes then achieve PCO-based coordination via the periodic transmission of beacon packets at the MAC layer. As such, for channels with an equal number of nodes, DT-SCS converges to synchronized beacon packet transmission at the MAC layer in a completely uncoordinated manner. In order to combat the well-know phenomenon of Cross-Technology Interference (CTI) a cross-layer mechanism, CrossZig, has been implemented in [39], based on the exploitation of information at the physical layer in order to detect the presence of CTI in a corrupted packet.

A different perspective of the interference management has been considered in [47] and [41], where a novel solution to allow to secondary users the access of allocated spectrum has been proposed. The study has been based on the major consideration that a big bottleneck in cognitive radio systems is based on finding the best available channel as fast as possible.

A totally different approach to face the enormous quantity of data generated by IoT devices, is to try to reduce the sending of useless data, based on the adoption of effective predictive approaches.

In [50] we have considered the concept of high spatio-temporal correlated data and we have proposed a Belief Propagation (BP) algorithm to derive methods to drastically reduce the number of transmitted messages, by keeping an high accuracy in terms of global information.

Together with interference management approaches it is also important to figure out tools to support network operator for mitigation of the impact of failures on their infrastructures. The need of advanced Network Planning and Management Tool (NPMT) has been considered in [30].

## 7.7. Vehicular Networks

**Participants:** Nathalie Mitton, Valeria Loscri.

[27] studies the information delivery delay analysis for roadside unit deployment in a vehicular ad hoc network (VANET) with intermittent connectivity. A mathematical model is developed to describe the relationship between the average delay for delivering road condition information and the distance between two neighbor RSUs deployed along a road. The derived mathematical model considers a straight highway scenario where two RSUs are deployed at a distance without any direct connection and vehicles are sparsely distributed on the road with road condition information randomly generated between the two neighbor RSUs. Moreover, the model takes into account the vehicle speed, the vehicle density, the likelihood of an incident, and the distance between two RSUs. The effectiveness of the derived mathematical model is verified through simulation results. Given the information delivery delay constraint of a time-critical application, this model can be used to estimate the maximum distance allowed between two neighbor RSUs, which can provide a reference for the deployment of RSUs in such scenarios.

But Vehicular Networks can also convey social networks. In [53], we survey recent literature on Vehicular Social Networks that are a particular class of vehicular ad hoc networks, characterized by social aspects and features. Starting from this pillar, we investigate perspectives of next generation vehicles under the assumption of social networking for vehicular applications (i.e., safety and entertainment applications). This paper plays a role as a starting point about socially-inspired vehicles, and main related applications, as well as communication techniques. Vehicular communications can be considered as the "first social network for automobiles", since each driver can share data with other neighbors. As an instance, heavy traffic is a common occurrence in some areas on the roads (e.g., at intersections, taxi loading/unloading areas, and so on); as a consequence, roads become a popular social place for vehicles to connect to each other. Human factors are then involved in vehicular ad hoc networks, not only due to the safety related applications, but also for entertainment

purpose. Social characteristics and human behavior largely impact on vehicular ad hoc networks, and this arises to the vehicular social networks, which are formed when vehicles (individuals) "socialize" and share common interests. This survey describes the main features of vehicular social networks, from novel emerging technologies to social aspects used for mobile applications, as well as main issues and challenges. Vehicular social networks are described as decentralized opportunistic communication networks formed among vehicles. They exploit mobility aspects, and basics of traditional social networks, in order to create novel approaches of message exchange through the detection of dynamic social structures. An overview of the main state-of-the-art on safety and entertainment applications relying on social networking solutions is also provided.

Cognitive Radio (CR) together with vehicular networks have been considered with an integrated and synergic perspective in [55], since CR technology is foreseen as a very effective tool to improve the communication efficiency in the context of vehicular networked systems.

## 7.8. Self-deployment and coverage

**Participants:** Nathalie Mitton, Tahiry Razafindralambo.

Controlled mobility in wireless sensor networks can provide many services. One of the most challenging one is coverage. Coverage can be needed either for monitoring control of specific area or point of interest or for deploying a communication network. This latter case is required for instance in post-disaster situations. In post-disaster scenarios, for example, after earthquakes or floods, the traditional communication infrastructure may be unavailable or seriously disrupted and overloaded. Therefore, rapidly deployable network solutions are needed to restore connectivity and provide assistance to users and first responders in the incident area. This work surveys the solutions proposed to address the deployment of a network without any a priori knowledge about the communication environment for critical communications. The design of such a network should also allow for quick, flexible, scalable, and resilient deployment with minimal human intervention. We survey this kind of approaches in [20].

In [13], we present a decentralized deployment algorithm for wireless mobile sensor networks focused on deployment Efficiency, connectivity Maintenance and network Reparation (EMR). We assume that a group of mobile sensors is placed in the area of interest to be covered, without any prior knowledge of the environment. The goal of the algorithm is to maximize the covered area and cope with sudden sensor failures. By relying on the locally available information regarding the environment and neighborhood, and without the need for any kind of synchronization in the network, each sensor iteratively chooses the next-step movement location so as to form a hexagonal lattice grid. Relying on the graph of wireless mobile sensors, we are able to provide the properties regarding the quality of coverage, the connectivity of the graph and the termination of the algorithm. We run extensive simulations to provide compactness properties of the deployment and evaluate the robustness against sensor failures. We show through the analysis and the simulations that EMR algorithm is robust to node failures and can restore the lattice grid. We also show that even after a failure, EMR algorithm call still provide a compact deployment in a reasonable time.

Routing a fleet of robots in a known surface is a complex problem. It consists in the determination of the exact trajectory each robot has to follow to collect information. The objective pursued in [38] is to maximize the exploration of the given surface. To ensure the robots can execute the mission in a collaborative manner, connectivity constraints are considered. These constraints guarantee that robots can communicate among each other and share the collected information. Moreover, the trajectories of the robots need to respect autonomy constraints.

## 7.9. Controlled Mobility for additional services

**Participants:** Nathalie Mitton, Valeria Loscri, Jean Cristanel Razafimandimby Anjalalaina.

Wireless sensor networks (WSNs) have been of very high interest for the research community since years, but most of the time, the mobility of nodes have been considered as an obstacle to overcome. In the contrary, in have tried to adopt another perspective and see it as an asset to exploit to provide additional services.

In [19], we leverage on the ability of mobile nodes to replace or recharge static sensors. Two main approaches can be identified that target this objective: either "recharging" or "replacing" the sensor nodes that are running out of energy. Of particular interest are solutions where mobile robots are used to execute the above mentioned tasks to automatically and autonomously maintain the WSN, thus reducing human intervention. Recently, the progress in wireless power transfer techniques has boosted research activities in the direction of battery recharging, with high expectations for its application to WSNs. Similarly, also sensor replacement techniques have been widely studied as a means to provide service continuity in the network. Objective of [19] is to investigate the limitations and the advantages of these two research directions. Key decision points must be identified for effectively supporting WSN self-maintenance: (i) which sensor nodes have to be recharged/replaced; (ii) in which order the mobile robot is serving (i.e., recharging/replacing) the nodes and by following which path; (iii) how much energy is delivered to a sensor when recharged. The influence that a set of parameters, relative to both the sensors and the mobile robot, on the decisions will be considered. Centralized and distributed solutions are compared in terms of effectiveness in prolonging the network lifetime and in allowing network self-sustainability. The performance evaluation in a variety of scenarios and network settings offers the opportunity to draw conclusions and to discuss the boundaries for one technique being preferable to the other.

Mobility can also help for collecting data in wireless sensor networks [29]. The sensor data collection problem using data mules have been studied fairly extensively in the literature. However, in most of these studies, while the mule is mobile, all sensors are stationary. The objective of most of these studies is to minimize the time needed by the mule to collect data from all the sensors and return to the data collection point, from where it embarked on its data collection journey. The problem studied in this paper has two major differences with the earlier studies. First, in this study we assume that both the mule as well as the sensors are mobile. Second, we do not attempt to minimize the data collection time. Instead we minimize the number of mules that will be needed to collect data from all the sensors, subject to the constraint that the data collection process has to be completed within some pre-specified time. We show that the mule minimization problem is NP-Complete and provide a solution by first transforming it to a generalized version of the minimum flow problem in a network and then solving it optimally using Integer Linear Programming. Finally, we evaluate our algorithms through extensive simulation and present the results.

Internet of Robotic Things (IoRT) is a new concept introduced for the first time by ABI Research. Unlike the Internet of Things (IoT), IoRT provides an active sensorization and is considered as the new evolution of IoT.

This new concept will bring new opportunities and challenges, while providing new business ideas for IoT and robotics' entrepreneurs.

In [46], we focus particularly on two issues: (i) connectivity maintenance among multiple IoRT robots, and (ii) their collective coverage.

We propose (i) IoRT-based, and (ii) a neural network control scheme to efficiently maintain the global connectivity among multiple mobile robots to a desired quality-of-service (QoS) level. The proposed approaches will try to find a trade-off between collective coverage and communication quality.

The IoT-based approach is based on the computation of the algebraic connectivity and the use of virtual force algorithm.

The neural network controller, in turn, is completely distributed and mimics perfectly the IoT-based approach. Results show that our approaches are efficient, in terms of convergence, connectivity, and energy consumption.

## 7.10. New and other communication paradigms

**Participants:** Nathalie Mitton, Valeria Loscri.

Interconnection and self-organized systems are normally populated with heterogeneous and different devices. The differences range from computational capabilities, storage size, etc. Instead of considering the heterogeneity as a limitation, it is possible to "turn it" as a primitive control of the system, in order to realize more robust and more resilient communication systems.

Based on those considerations, we have studied and analyzed the specific features of devices belonging to the category of micro-nano nodes that are however, required to interact with up-sized devices.

In order to improve the understanding of the behavior of micro/nano-sized devices, we have considered fundamental the analysis in specific applications and environment, where this kind of devices can be largely exploited, such as on/in-body networks applications.

Indeed, we retain that bio-medical applications can be advantaged by an effective and efficient communication and cooperation of devices deployed both on top of the body and inside it. Even if the research community recognizes a great importance to the study of interaction between the Human Immune System (HIS) and nano devices, this branch of research is in its infancy due to the major issue to model the HIS. A theoretical derivation of HIS and its interaction with a nanoparticulate system have been proposed in [15]. Some experimental results have been derived in [16], where specific parameters, e.g. temperature variations, Ph, etc. have been considered to establish the biocpmpatibility of TiO2 particles with human tissues.

A step ahead in this direction has consisted in the consideration of alternative particles as potential information carriers always in the context of biological environments. In [40] we have studied *phonons* as information carriers, we have derived a channel modeling and evaluated the theoretical capacity. The main reasons for taking into consideration this type of nanoparticles are twofold. Firstly, phonons represent something that is naturally generated in a biological context with the application of a tolerable electromagnetic field and secondly they represent a straightforward way to implement nanomachines, since their native size.

## 7.11. Modelling and experimentations of interferences and other PHY effects

**Participants:** Nathalie Mitton, Valeria Loscri.

In the era of Internet of Things (IoT), the development of Wireless Sensor Networks (WSN) arises different challenges. Two of the main issues are electromagnetic interference and the lifetime of WSN nodes. In [48], we show and evaluate experimentally the relation between interference and energy consumption, which impacts the network lifetime. We present a platform based on commercially available low-cost hardware in order to evaluate the impact of electromagnetic interference in 2.4 GHz ISM band on energy consumption of WSN. The energy measurements are obtained separately from each electronic component in the node. Interference and energy measurements are conducted in an anechoic chamber and in an office-type lab environment. X-MAC protocol is chosen to manage the Radio Duty Cycle of the nodes and its energy performance is evaluated. The energy consumption transmitter nodes is analyzed particularly in this work. Moreover, this energy consumption has been quantified and differentiated according to the number of (re-)transmissions carried out by the transmitter as well as the number of ACK packets sent by the receiver for a single packet. Finally, we use a model of real battery to calculate the lifetime of the node for operation within different interference level zones. This study lays the basis for further design rules of communication protocols and development of WSNs.

In [49], we propose a WSN architecture for wild animal monitoring. The key requirements of the system are long range transmissions and low power consumption. Indeed, the animals could be spread over vast areas. Kruger National Park in South Africa (19485 km2) is the potential zone of implementation of the network. On the other hand, size and weight limitations of wearable devices must be respected, which limits the size and capacity of battery. Moreover, battery replacement is a difficult and expensive process. So, low energy consumption is essential to extend the network lifetime. Some animal tracking projects [3] use GSM to transmit collected data to insure the coverage over a large area. However, high energy consumption of GSM and lack of coverage of the deployment area do not meet the essential requirements of the application. LoRa technology provides both long range transmissions and low power operation. This technology could be an appropriate solution for PREDNET project. The contribution of this work is multiple: 1) we defined communication parameters of LoRa radio for PREDNET WSN; 2) we performed radio propagation simulation for chosen parameters to estimate the coverage area for both urban and wilderness (rural) scenarios; 3) we confirmed the propagation simulations with range tests; 4) we measured experimentally the Packet Error Rate (PER) of transmissions.

Terahertz frequency band is an emerging research area related to nano-scale communications. In this frequency range, specific features can provide the possibility to overcome the issues related to the spectrum scarcity and capacity limitation.

Apart high molecular absorption, and very high reflection loss that represent main phenomena in THz band, we can derive the characteristics of the channel affected by chirality effects occurring in the propagation medium, specifically , in the case where a Giant Optical Activity is present. This effect is typical of the so-called chiral-metamaterials in (4-10) THz band, and is of stimulating interest particularly for millimeter wireless communications.

In [51], [25], we analyze the behavior of specific parameters of a chiral-metamaterial, like the relative electrical permitivity, magnetic permeability and chirality coefficients, and from that we derive the channel behavior both for Line-of-Sight and No Line-of-Sight propagations. We notice the presence of spectral windows, due to peaks of resonance of chiral parameter.

Finally, performances of the chirality-affected channel have been assessed in terms of (i) channel capacity, (ii) propagation delay, and (iii) coherence band-width, for different distances.

Thanks to the exploitation of frequencies in the interval ranging from 0.06 to 10 THz, it is envisioned the possibility to overcome the issues related to the spectrum scarcity and capacity limitation. On the other hand, the design of new channel models, able to capture the inherent features of the phenomenons related with this specific field is of paramount importance. Very high molecular absorption, and very high reflection loss are peculiarities phenomenons that need to be included in these models. In [26], we present a full-wave propagation model of the electromagnetic field that propagates in the THz band both for Line-of-Sight and Non-Line-of-Sight propagation models. In the full-wave model, we also introduce the chirality effects occurring in the propagation medium, i.e., a chiral metamaterial.

<p style="text-align:center;color:red"><strong>INOCS Team</strong></p>

# 6. New Results

## 6.1. Large scale complex structure optimization

**New decomposition methods for the time-dependent combined network design and routing problem:** A significant amount of work has been focussed on the design of telecommunication networks. The performance of different Integer Programming models for various situations has been computationally assessed. One of the settings that has been thoroughly analyzed is a variant where routing decisions (for time-dependent traffic demand), and network design, are combined in a single optimization model. Solving this model with a state-of-the-art solver on representative network topologies, shows that this model quickly becomes intractable. With an extended formulation, both the number of continuous flow variables and the number of fixed charge capacity constraints are multiplied by a factor $|V|$ (where $V$ represents the set of nodes) leading to large model. However, the linear relaxation of this extended formulation yields much better lower bounds. Nevertheless, even if the extended model provides stronger lower bounds than the aggregated formulation, it suffers from its huge size: solving the linear relaxation of the problem quickly becomes intractable when the network size increases, making the linear relaxation expensive to solve. This observation motivates the analysis of decomposition methods [30].

**Convex piecewise linear unsplittable multicommodity flow problems** We studied the multi-commodity flow problem with unsplittable flows, and piecewise-linear costs on the arcs. They show that this problem is NP-hard when there is more than one commodity. We propose a new MILP models for this problem, that was compared to two formulations commonly used in the literature. The computational experiments reveal that the new model is able to obtain very strong lower bounds, and is very efficient to solve the considered problem [40].

**Tree Reconstruction Problems:** We studied the problem of reconstructing a tree network by knowing only its set of terminal nodes and their pairwise distances, so that the reconstructed network has its total edge weight minimized. This problem has applications in several areas, namely the inference of phylogenetic trees and the inference of routing networks topology. Phylogenetic trees allow the understanding of the evolutionary history of species and can assist in the development of vaccines and the study of biodiversity. The knowledge of the routing network topology is the basis for network tomography algorithms and it is a key strategy to the development of more sophisticated and ambitious traffic control protocols and dynamic routing algorithms [31].

**Comparison of formulations and solution methods for the discrete ordered p-median problem:** We presented several new formulations for the Discrete Ordered Median Problem (DOMP) based on its similarity with some scheduling problems. Some of the new formulations present a considerably smaller number of constraints to define the problem with respect to some previously known formulations. Furthermore, the lower bounds provided by their linear relaxations improve the ones obtained with previous formulations in the literature even when strengthening is not applied. We also present a polyhedral study of the assignment polytope of our tightest formulation showing its proximity to the convex hull of the integer solutions of the problem. Several resolution approaches, among which we mention a branch and cut algorithm, are compared. Extensive computational results on two families of instances, namely randomly generated and from Beasley's OR-library, show the power of our methods for solving DOMP [34].

**New models and algorithms for integrated vehicle routing problems**

We address a real-life inventory routing problem, which consists in designing routes and managing the inventories of the customers simultaneously. The problem was introduced during the 2016 ROADEF/EURO challenge. The proposed problem is original and complex for several reasons : the logistic ratio optimization objective, the hourly time-granularity for inventory constraints, the driver/trailer allocation management. Clearly, this problem is an optimization problem with complexe structure, for which we proposed a branch-cut-and-price based method : a cut and-column generation procedure was developed, along with a heuristic pricing algorithm to generate new columns and a heuristic fixing procedure to generate integer solutions. The solution method allowed the team including INOCS members to qualify to the final phase of the ROADEF/EURO challenge 2016 [41].

**Column generation approach for pure parsimony haplotyping:**  The knowledge of nucleotides chains that compose the double DNA chain of an individual has a relevant role in detecting diseases and studying populations. However, determining experimentally the single nucleotides chains that, paired, form a certain portion of the DNA is expensive and time-consuming. Mathematical programming approaches have been proposed instead, e.g. formulating the Haplotype Inference by Pure Parsimony problem (HIPP). Abstractly, we are given a set of genotypes (strings over a ternary alphabet 0, 1, 2) and we want to determine the smallest set of haplotypes (binary strings over the set 0, 1) so that each genotype can be 'generated' by some pair of haplotypes, meaning that they are compatible with the genotype and can fully explain its structure. In order to deal with larger instances, we proposed a new model involving an exponential number of variables to be solved via column generation, where variables are dynamically introduced into the model by iteratively solving a pricing problem. We compared different ways of solving the pricing problem, based on integer programming, smart enumeration and local search heuristic. The efficiency of the approach is improved by stabilization and by a heuristic to provide a good initial solution. Results show that, with respect to the linear relaxations of both the polynomial and exponential-size models, our approach yields a tighter formulation and outperforms in both efficiency and effectiveness the previous model for instances with a large number of genotypes [39].

## 6.2. Bilevel Programming

**Bilevel approaches for energy management problems:** We have proposed the first bilevel pricing models to explore the relationship between energy suppliers and customers who are connected to a smart grid. Due to their definition, bilevel models enable to integrate customer response into the optimization process of supplier who aims to maximize revenue or minimize capacity requirements. In our setting, the energy provider acts as a leader (upper level) that takes into account a smart grid (lower level) that minimizes the sum of users' disutilities. The latter bases its decisions on the hourly prices set by the leader, as well as the schedule preferences set by the users for each task. The pricing problems, we model, belong to the category of single leader single follower problems. Considering both the monopolistic and competitive environment we present two bilevel bilinear bilinear problems with continuous variables. Heuristics solutions methods are defined to solve large size instances of the models. They are based on the interactions between prices, schedules and peaks. Numerical results on randomly generated instances illustrate numerically the validity of the approach, which achieves an 'optimal trade-off between three objectives: revenue, user cost, and peak demand. Moreover, they put into highlight the ability of the heuristics to produce high quality results compared to the solution of MIP reformulations of the models[36].

**New formulations for solving Stackelberg games:**  We analyzed general Stackelberg games (SGs) and Stackelberg security games (SSGs). SGs are hierarchical adversarial games where players select actions or strategies to optimize their payoffs in a sequential manner. SSGs are a type of SGs that arise in security applications, where the strategies of the player that acts first consist in protecting subsets of targets and the strategies of the followers consist in attacking one of the targets. We review existing mixed integer optimization formulations in both the general and the security setting and present new formulations for the the second one. We compare the SG formulations and the SSG formulations both from a theoretical and a computational point of view. We indentify which formulations provide tighter linear relaxations and show that the strongest formulation for the security version is ideal in the case of one single attacker. Our computational experiments show that the new formulations can be solved in shorter times [46].

# 6.3. Robust/Stochastic programming

**Decomposition method for stochastic staff management problems :** We addressed an integrated shift scheduling and load assignment optimization problem for attended home delivery, which is a last-mile delivery service requiring the presence of the customer for the delivery. We were mainly interested in generating a daily master plan for each courier. We proposed a tactical problem integrating a shift scheduling problem and a load assignment problem under demand uncertainty, which was modeled as a two-stage stochastic programming model. This model integrates two types of decisions. First-stage decisions are related to the design of a schedule that includes the periods of the day in which each courier must work and the o-d pairs to visit at each time period. Second-stage decisions (recourse actions) consist of the allocation of a number of packages to be delivered at each time period, for each o-d pair, by each courier, such that the demand (number of packages to deliver) for each scenario is satisfied. Recourse is the ability to take corrective actions after a random event has taken place. The objective is to minimize the sum of the daily staffing cost plus the expected daily recourse cost. To solve this problem, we proposed and implemented a multi-cut integer L-shaped algorithm, where the second stage decomposes by time periods and by demand scenarios. To strengthen the first stage model, some valid inequalities are added, and some of the existing constraints are lifted. Results on real-world based instances from a delivery company demonstrate that our approach provides robust tactical solutions that easily accommodate to fluctuations in customer orders, preventing additional costs related to the underutilization of couriers and the use of external couriers to satisfy all delivery requests [37], [43].

<span style="color:red">**LINKS Project-Team**</span>

# 7. New Results

## 7.1. Querying Heterogeneous Linked Data

### 7.1.1. Provenance

The computation of the provenance of a query answer is a classical problem in database theory. It consists in aggregating the impact of tuples of a database to a query answer. This allows to give an explanation of the query answers, that can help to judge their reliability. The computation of the provenance of a query answer is thus an aggregation problem as studied by the ANR project *Aggreg* .

P. Bourhis [20] showed at **PODS** — the top conference on database theory — that the lineage of MSO queries on treelike database instances is tractable, but not on other instances. This work was in cooperation with Telecom ParisTech and ENS Paris. As a first application, he can show that MSO query evaluation on probabilistic databases is tractable for tree like database instances, but not otherwise.

P. Bourhis applied in cooperation with Tel Aviv, provenance problems to recommendation systems. This allows to explain the end result by summarising with similar data without changing significantly results obtained in general by aggregation on the data. The corresponding tool was demonstrated at **EDBT** [32].

### 7.1.2. Certain Query Answering and Access Control

The problem of certain query answering consists in finding which are the certain answer of a query in a database with incomplete data, and a set of constraints representing available the knowledge on the incomplete data.

P. Bourhis [24] presented at **LICS** — the top conference in logic in computer science — a general framework for querying databases with visible and invisible relations. This work was done in cooperation with Oxford, Santa Cruez, and Bordeaux. His framework is motivated by the problem of access control for relational databases, i.e. of data leakage in relational views, but generalizes at the same time the problem of certain query answering. Invisible relations are subject to the open world assumption possibly under constraints as usual in certain query answering, while visible relations are subject to the closed world assumption. Bourhis then show that it is decidable, whether a conjunctive has an answer in this framework, when given the visible relation, the constraints, and the query as inputs. He also studies the complexity of this problem. It turns out the complexity increases from polynomial to doubly exponential, compared to certain query answering, since adding visible relations subject to the closed world assumption.

P. Bourhis studied at **IJCAI** [19] certain query answering with some transitive closure constraints, which allow to define a constraints with recursion. This work was done in collaboration with Oxford and Telecom ParisTech.

The problem of ontological query containment consists in establishing whether the certain answers of two queries subject to an ontology are included in each other. P. Bourhis [26] studied at **KR** this problem for several closely related formalisms: monadic disjunctive Datalog (MDDLog), MMSNP (a logical generalization of constraint satisfaction problems) and ontology-mediated queries (OMQs). This work was done in cooperation with Bremen.

### 7.1.3. Recursive Queries

At **LICS** [21] again, P. Bourhis showed in collaboration with Oxford how to lift a major restriction on decidable fixpoint logics that can define recursive queries (such as C2RPQs), specifically on guarded logic. This allows to improve significantly expressiveness of decidable fixpoint logics.

A. Lemay contributed at **TKDE** [14] the *gMark* benchmark, a tool to generate large size graph database and an associated set of queries. This work was done in cooperation with Eindhoven and previous members of Links that are now in Lyon and Clérmont-Ferrant. The tool was also demonstrated at **VLDB** [13]. Its main interest is a great flexibility (the generation of the graph can be done from a simple schema, but can also incorporate elaborate a parameters), an ability to generate recursive queries, and the possibility to generate large sets of queries of a desired selectivity. This benchmark allowed for instance to highlight difficulties for the existing query engines to deal with recursive queries of high selectivity.

### 7.1.4. Data Integration

P. Bourhis and S. Staworko in cooperation with Bordeaux and Oxford presented at **TODS** [17] their work on bounded repairability for regular tree languages, which is a study on whether a tree document (typically XML) can be repaired to fit a given target tree language within a bounded amount of tree editing operations. The article studies the complexity of different classes of tree languages such as non-recursive DTDs, recursive DTDs, or languages by arbitrary bottom-up tree automaton.

J.M. Lozano started his PhD project under the supervision of I. Boneva and S. Staworko. His topic subscribes the ANR project *Datacert* on data integration and certification.

### 7.1.5. Schema Validation

A. Boiret, V. Hugot and J. Niehren studied schemas for JSON documents in **Information and Computation** [15]. This work was done in collaboration with Paris 7. A JSON document is an unordered data trees, so schemas for such documents are best seen as automata for unordered data trees. The paper generalizes several previous formalisms for automata on unordered trees in a uniform framework. Whether the equivalence of two schemas can be tested in P-time is studied for various instances of the framework.

This work subscribes to the ANR project *Colis* where unranked data trees are used as models of linux file systems. In this context, N. Bacquey started his postdoc on the verification of linux installation scripts.

## 7.2. Managing Dynamic Linked Data

### 7.2.1. Complex Event Processing

Complex event processing can be seen as the problem is to answer queries on data graphs, for graphs that arrive on streams. These queries may contain aggregates, so this work subscribes to the ANR project *Aggreg*.

In his PhD thesis, T. Sebastian [12] developed with his supervisor J. Niehren streaming algorithms covering all of XPath 3.0 queries on XML streams. For this, they proposed a higher-order query language $\lambda$XP, showed how to give a formal semantics of all of XPath 3.0 by compilation to $\lambda$XP, and then how to evaluate $\lambda$XP queries on XML streams. These algorithms were implemented in the QuiXPath tool.

At **SOFSEM**, they proposed a new technique to speed up the evaluation of navigational XPath queries on XML streams based on document projection. The idea is to skip those parts of the stream that are irrelevant for the query. This speeds up the evaluation of navigation XPath queries by a factor of 4 in usual Xpath benchmarks.

M. Sakho started his PhD project on hyperstreaming query answering algorithms for graphs under the supervision of J. Niehren and I. Boneva. Part of this work will be continued with out visitor D. Vrgoc from Santiago di Chili.

### 7.2.2. Data Centric Workflows

Data-centric workflows are complex programs that can query and update a database. The usage of data-centric workflows for crowd sourcing is the topic of the ANR Project *HeadWork*.

In collaboration with ENS Cachan and San Diego, P. Bourhis presented at **ICDT** [18] techniques on collaborative access control in a distributed query and data exchange language (Webdamlog). The goal of this work was to provide a semantic to data exchange rules defined by Webdamlog. It also allowed to prove that it is possible to formally verify whether there are data leakages.

P. Bourhis with Tel Aviv defined at **ICDE** [25] a notion of provenance for data-centric workflows, and proved that it can be used to explain the provenance of fact in the final instance of an execution. This provenance is used to answer three main questions: *why* does a specific tuple appear in the answer of a query, *what if* the initial database is changed (Revision problem), and *how to* change the query to obtain a missing tuple.

## 7.3. Linking Data Graphs

### 7.3.1. *Learning Transformations*

We consider the problem to learn queries and query-based transformations on semi-structured data from examples.

A. Boiret obtained his PhD for his work on the "Normalization and Learning of Transducers on Trees and Words" under the supervision of J. Niehren and A. Lemay. In this year, he showed how to learn top-down tree transformations with regular schema restrictions [31], [33], [34]. At **LATA** [22], he deepened a result of a previous PhD student of Links on learning sequential tree-to-word transducers (with output concatenation), by showing who to find normal forms for less restrictive linear tree-to-word transducers. At **DLT** [23], he could show in cooperation with Munich, that the equivalence problem of this class of transducers is in polynomial time, even though their normal forms may be of exponential size.

In the context of learning RDF graph transformations, S. Staworko presented a cooperation with Edinburg at **VLDB** [27]. Using bisimulation technique, he aims at aligning datas of two RDF Graphs that takes into account blank velues, changes in ontology and small differences in data values and in the structure of the graph. the alignement of graphs is an important first step for the inference of transformations.

### 7.3.2. *Learning Join Queries*

S. Staworko published in **TODS** an article [16] on learning join queries from user examples in collaboration with Universities of Lyon and Clermont-Ferrand that present techniques that allow the automatic construction of a join query through interaction with a user that simply labels sets of tuples to indicate whether the tuple is in the target query or not.

<span style="color:red; text-align:center">**MAGNET Project-Team**</span>

# 7. New Results

## 7.1. Decentralized and Private Learning

In [13], we address the problem of decentralized minimization of pairwise functions of the data points, where these points are distributed over the nodes of a graph defining the communication topology of the network. This general problem finds applications in ranking, distance metric learning and graph inference, among others. We propose new gossip algorithms based on dual averaging which aims at solving such problems both in synchronous and asynchronous settings. The proposed framework is flexible enough to deal with constrained and regularized variants of the optimization problem. Our theoretical analysis reveals that the proposed algorithms preserve the convergence rate of centralized dual averaging up to an additive bias term. We present numerical simulations on Area Under the ROC Curve (AUC) maximization and metric learning problems which illustrate the practical interest of our approach.

In [19], we consider a set of learning agents in a collaborative peer-to-peer network, where each agent learns a *personalized model* according to its own learning objective. The question addressed in this paper is: how can agents improve upon their locally trained model by communicating with other agents that have similar objectives? We introduce and analyze two asynchronous gossip algorithms running in a fully decentralized manner. Our first approach, inspired from label propagation, aims to smooth pre-trained local models over the network while accounting for the confidence that each agent has in its initial model. In our second approach, agents jointly learn and propagate their model by making iterative updates based on both their local dataset and the behavior of their neighbors. Our algorithm for solving this challenging optimization problem relies on the Alternating Direction Method for Multipliers (ADMM).

In [20], we propose a decentralized protocol for a large set of users to privately compute averages over their joint data, which can later be used to learn more complex models. Our protocol can find a solution of arbitrary accuracy, does not rely on a trusted third party and preserves the privacy of users throughout the execution in both the honest-but-curious and malicious adversary models. Furthermore, we design a verification procedure which offers protection against malicious users joining the service with the goal of manipulating the outcome of the algorithm.

## 7.2. Natural Language Processing

In [12], we introduce a simple semi-supervised approach to improve implicit discourse relation identification. This approach harnesses large amounts of automatically extracted discourse connectives along with their arguments to construct new distributional word representations. Specifically, we represent words in the space of discourse connectives as a way to directly encode their rhetorical function. Experiments on the Penn Discourse Treebank demonstrate the effectiveness of these task-tailored representations in predicting implicit discourse relations. Our results indeed show that, despite their simplicity, these connective-based representations outperform various off-the-shelf word embeddings, and achieve state-of-the-art performance on this problem.

Along the PhD thesis of THIBAULT LIÉTARD, we are working on learning a similarity between text entities for the task of coreference resolution. Unlike indirect classification criteria often used in the literature, the similarity function naturally operates on pairs of mentions and several relevant objectives can be considered. For instance, we can learn the parameters of the similarity function such that the similarity of a given mention to its closest antecedent coreferent mention is larger than to any closer non-coreferential antecedent candidate. The resulting similarity scores can then be plugged into a greedy clustering procedure, or used to build a weighted graph of mentions to be clustered by spectral algorithms. For the representations of (pairs of) mentions on which the similarity function is learned, we consider both traditional linguistic features as well as external information about the general context of occurrence of the mentions using word embeddings.

Along the PhD thesis of MATHIEU DEHOUCK, we study the problem of cross-lingual dependency parsing, aiming at leveraging training data from different source languages to learn a parser in a target language. Specifically, this approach first constructs word vector representations that exploit structural (i.e., dependency-based) contexts but only considering the morpho-syntactic information associated with each word and its contexts. These delexicalized word embeddings, which can be trained on any set of languages and capture features shared accross languages are then used in combination with standard language-specific features to train a lexicalized parser in the target language. We evaluate our approach through experiments on a set of eight different languages that are part the Universal Dependencies Project. Our main results show that using such embeddings (monolingual or multilingual) achieves significant improvements over monolingual baselines. The work is submitted.

## 7.3. Edge Prediction in Networks

In [18] we address the problem of classifying the links of signed social networks given their full structural topology. In the problem of edge sign prediction, we are given a directed graph (representing a social network), and our task is to predict the binary labels of the edges (i.e., the positive or negative nature of the social relationships). Many successful heuristics for this problem are based on the troll-trust features, estimating at each node the fraction of outgoing and incoming positive/negative edges. We show that these heuristics can be understood, and rigorously analyzed, as approximators to the Bayes optimal classifier for a simple probabilistic model of the edge labels. We then show that the maximum likelihood estimator for this model approximately corresponds to the predictions of a label propagation algorithm run on a transformed version of the original social graph. Extensive experiments on a number of real-world datasets show that this algorithm is competitive against state-of-the-art classifiers in terms of both accuracy and scalability. Finally, we show that troll-trust features can also be used to derive online learning algorithms which have theoretical guarantees even when edges are adversarially labeled.

In [16], we address the problem of predicting connections between a set of data points. We focus on the *graph reconstruction* problem, where the prediction rule is obtained by minimizing the average error over all $n(n-1)/2$ possible pairs of the $n$ nodes of a training graph. Our first contribution is to derive learning rates of order $O(\log n/n)$ for this problem, significantly improving upon the slow rates of order $O(1/\sqrt{n})$ established in the seminal work of [27]. Strikingly, these fast rates are universal, in contrast to similar results known for other statistical learning problems (e.g., classification, density level set estimation, ranking, clustering) which require strong assumptions on the distribution of the data. Motivated by applications to large graphs, our second contribution deals with the computational complexity of graph reconstruction. Specifically, we investigate to which extent the learning rates can be preserved when replacing the empirical reconstruction risk by a computationally cheaper Monte-Carlo version, obtained by sampling with replacement $B \ll n^2$ pairs of nodes. Finally, we illustrate our theoretical results by numerical experiments on synthetic and real graphs.

## 7.4. Mining Geotagged Social Data

Data generated on location-based social networks provide rich information on the whereabouts of urban dwellers. Specifically, such data reveal who spends time where, when, and on what type of activity (e.g., shopping at a mall, or dining at a restaurant). That information can, in turn, be used to describe city regions in terms of activity that takes place therein. For example, the data might reveal that citizens visit one region mainly for shopping in the morning, while another for dining in the evening. Furthermore, once such a description is available, one can ask more elaborate questions. For example, one might ask what features distinguish one region from another – some regions might be different in terms of the type of venues they host and others in terms of the visitors they attract. As another example, one might ask which regions are similar across cities. In [11], we present a method to answer such questions using publicly shared Foursquare data. Our analysis makes use of a probabilistic model, the features of which include the exact location of activity, the users who participate in the activity, as well as the time of the day and day of week the activity takes place. Compared to previous approaches to similar tasks, our probabilistic modeling approach allows us to make

minimal assumptions about the data – which relieves us from having to set arbitrary parameters in our analysis (e.g., regarding the granularity of discovered regions or the importance of different features). We demonstrate how the model learned with our method can be used to identify the most likely and distinctive features of a geographical area, quantify the importance features used in the model, and discover similar regions across different cities. Finally, we perform an empirical comparison with previous work and discuss insights obtained through our findings. Our results were also presented through an interactive demo at the 25th World Wide Web Conference [21].

## 7.5. Learning from Non-iid Data

In [14] we deal with the generalization ability of classifiers trained from non-iid evolutionary-related data in which all training and testing examples correspond to leaves of a phylogenetic tree. For the realizable case, we prove PAC-type upper and lower bounds based on symmetries and matchings in such trees.

In [9], we studied learning problems where the performance criterion consists of an average over tuples (e.g., pairs or triplets) of observations rather than over individual observations, as in many learning problems involving networked data (e.g., link prediction), but also in metric learning and ranking. In this setting, the empirical risk to be optimized takes the form of a U-statistic, and its terms are highly dependent and thus violate the classic i.i.d. assumption. From a computational perspective, the calculation of such statistics is highly expensive even for a moderate sample size $n$, as it requires averaging $O(n^d)$ terms. We show that, strikingly, such empirical risks can be replaced by drastically computationally simpler Monte-Carlo estimates based on $O(n)$ terms only, usually referred to as incomplete U-statistics, without damaging the $O(1/\sqrt{n})$ learning rate of Empirical Risk Minimization (ERM) procedures. For this purpose, we establish uniform deviation results describing the error made when approximating a U-process by its incomplete version under appropriate complexity assumptions. Extensions to model selection, fast rate situations and various sampling techniques are also considered , as well as an application to stochastic gradient descent for ERM. Finally, numerical examples are displayed in order to provide strong empirical evidence that the approach we promote largely surpasses more naive subsampling techniques.

## 7.6. Adaptive Graph Construction

The efficiency of graph-based semi-supervised algorithms depends on the graph of instances on which they are applied. The instances are often in a vectorial form before a graph linking them is built. The construction of the graph relies on a metric over the vectorial space that help define the weight of the connection between entities. The classic choice for this metric is usually a distance measure or a similarity measure based on the euclidean norm. We claim that in some cases the euclidean norm on the initial vectorial space might not be the more appropriate to solve the task efficiently. In the work [17], we propose an algorithm that aims at learning the most appropriate vectorial representation for building a graph on which the task at hand is solved efficiently. In addition to experimental results showing the interest of such an approach, we define initial conditions under which the graph-based classification is ensured to perform optimally.

<p style="text-align:center"><span style="color:red">**MEPHYSTO Project-Team**</span></p>

# 7. New Results

## 7.1. Macroscopic behaviors of large interacting particle systems

### 7.1.1. *Stochastic acceleration and approach to equilibrium*

S. De Bièvre, Carlos Mejia-Monasterio (Madrid) and Paul E. Parris (Missouri) [57] studied thermal equilibration in a two-component Lorentz gas, in which the obstacles are modeled by rotating disks. They show that a mechanism of dynamical friction leads to a fluctuation-dissipation relation that is responsible for driving the system to equilibrium.

Stephan De Bièvre, Jeremy Faupin (Metz) and Schuble (Metz) [59] studied a related model quantum mechanically. Here a quantum particle moves through a field of quantized bose fields, modeling membranes that exchange energy and momentum with the particle. They establish a number of spectral properties of this model, that will be essential to study the time-asymptotic behavior of the system.

S. De Bièvre and collaborators analyse in [20] a multi-particle, kinetic version of a Hamiltonian model describing the interaction of a gas of particles with a vibrating medium. They prove existence results for weak solutions, and identify an asymptotic regime where the model, quite surprisingly, approaches the attractive Vlasov—Poisson system.

### 7.1.2. *Towards the weak KPZ universality conjecture*

One may start by considering the microscopic system in equilibrium (its measure is parametrized by the thermodynamical quantities under investigation). By removing the mean to the empirical measure and by scaling it properly, one would like to show that the random process, obtained by this rescaling, converges, as the size of the system is taken to infinity, to another random process which is a solution of some generalized stochastic PDE. Thanks to the remarkable recent result of M. Jara and P. Gonçalves [66], one has now all in hands to establish the latter result for a particular stochastic PDE known as the stochastic Burgers equation, and its companion, the Kardar-Parisi-Zhang (KPZ) equation. Indeed, in the latter paper, the authors introduce a new tool, called the second order Boltzmann-Gibbs principle, which permits to replace certain additive functionals of the dynamics by similar functionals given in terms of the density of the particles.

In [28], M. Simon in collaboration with T. Franco and P. Gonçalves, investigate the case of a microscopic dynamics with local defects, which is much harder. More precisely, the microscopic particle system is locally perturbed, and depending on the type of perturbation, the macroscopic laws can hold different boundary conditions. Since the ideas of [66] do not apply to the model considered there, they propose a new way to estimate the error in the replacement performed in the Boltzmann-Gibbs principle.

In the same spirit, M. Simon in collaboration with O. Blondel and P. Gonçalves investigate in [7] the class of kinetically constrained lattice gases that have been introduced and intensively studied in the literature in the past few years. In these models, particles are subject to restrictive constraints that make both approaches of [66] and [28] not work, so that new mathematical tools are needed. The main technical difficulty is that their model exhibits configurations that do not evolve under the dynamics and are locally non-ergodic. Their proof does not impose any knowledge on the spectral gap for the microscopic models. Instead, it relies on the fact that, under the equilibrium measure, the probability to find a blocked configuration in a finite box is exponentially small in the size of the box.

With these two recent results, M. Simon and coauthors contribute towards the *weak KPZ universality conjecture*, which states that a large class of one-dimensional weakly asymptotic conservative systems should converge to the KPZ equation.

### 7.1.3. Diffusion and fractional diffusion of energy

The rigorous derivation of the heat equation from deterministic systems of Newtonian particles is one of the most fundamental questions in mathematical physics. The main issue is that the existence of conservation laws and the high number of degrees of freedom impose very poor ergodic properties to the associated dynamical systems. A possible way out of this lack of ergodicity is to introduce stochastic models, in such a way that in one hand ergodicity issues are solved by the stochastic dynamics and in the other hand the qualitative behaviour of the system is not modified by the randomness. In these models, one starts with a chain of oscillators with a Hamiltonian dynamics, and one adds a stochastic component in such a way that the fundamental conservation laws (energy, momentum and *stretch* in this case) are maintained, and the corresponding Gibbs measures become ergodic.

It was already proved in [51] that these stochastic chains model correctly the behaviour of the conductivity. In particular, it is prove that Fourier law holds in dimension $d \geq 3$ if energy and momentum are conserved, and in any dimension if only energy is conserved. Once the conductivity has been successfully understood, one investigates the existence of the *hydrodynamic limit*, which fully describes the macroscopic evolution of the *empirical profiles* associated to the conserved quantity. In [41], M. Simon in collaboration with T. Komorowski and S. Olla consider the unpinned harmonic chain where the velocities of particles can randomly change sign. The only conserved quantities of the dynamics are the energy and the elongation. Using a diffusive space-time scaling, the profile of elongation evolves independently of the energy and follows a linear diffusive equation. The energy profile evolves following a non-linear diffusive equation involving the elongation. The presence of non-linearity makes the macroscopic limit non-trivial, and its mathematical proof requires very sophisticated arguments.

In [52] and [69] it has been previously shown that in the case of one-dimensional harmonic oscillators with noise that preserves the momentum, the scaling limit of the energy fluctuations is ruled by the *fractional* heat equation

$$\partial_t u = -(-\Delta)^{3/4} u.$$

This equation does not only predict the superdiffusivity of energy in momentum-conserving models, but it also predicts the speed at which it diverges. This result opens a way to a myriad of open problems. The main goal is to observe anomalous fractional superdiffusion type limit in the context of low dimensional asymmetric systems with several conserved quantities. In two recent papers by M. Simon in collaboration with C. Bernardin, P. Gonçalves, M. Jara, M. Sasada [53] & [32], they confirmed rigorously recent Spohn's predictions on the Lévy form of the energy fluctuations for a harmonic chain perturbed by an energy-volume conservative noise. In [32] they also showed the existence of a crossover between a normal diffusion regime and a fractional superdiffusion regime by tuning a parameter of a supplementary stochastic noise conserving the energy but not the volume.

## 7.2. Qualitative results in homogenization

### 7.2.1. Isotropy and loss of ellipticity in periodic homogenization

Since the seminal contribution of Geymonat, Müller, and Triantafyllidis, it is known that strong ellipticity is not necessarily conserved by homogenization in linear elasticity. This phenomenon is typically related to microscopic buckling of the composite material. In [24] G. Francfort and A. Gloria study the interplay between isotropy and strong ellipticity in the framework of periodic homogenization in linear elasticity. Mixtures of two isotropic phases may indeed lead to loss of strong ellipticity when arranged in a laminate manner. They show that if a matrix/inclusion type mixture of isotropic phases produces macroscopic isotropy, then strong ellipticity cannot be lost.

### 7.2.2. *From polymer physics to nonlinear elasticity*

OIn [23], M. Duerinckx and A. Gloria succeeded in relaxing one of the two unphysical assumptions made in [1] on the growth of the energy of polymer chains. In particular, [23] deals with the case when the energy of the polymer chain is allowed to blow up at finite deformation.

### 7.2.3. *The Clausius-Mossotti formula*

In the mid-nineteenth century, Clausis, Mossotti and Maxwell essentially gave a first order Taylor expansion for (what is now understood as) the homogenized coefficients associated with a constant background medium perturbed by diluted spherical inclusions. Such an approach was recently used and extended by the team MATHERIALS to reduce the variance in numerical approximations of the homogenized coefficients, cf. [46], [45], [72]. In [22], M. Duerinckx and A. Gloria gave the first rigorous proof of the Clausius-Mossotti formula and provided the theoretical background to analyze the methods introduced in [72].

## 7.3. Quantitative results in stochastic homogenization

### 7.3.1. *Quantitative results for almost periodic coefficients*

In [6], S. Armstrong, A. Gloria and T. Kuusi (Aalto University) obtained the first improvement over the thirty year-old result by Kozlov [70] on almost periodic homogenization. In particular they introduced a class of almost periodic coefficients which are not quasi-periodic (and thus strictly contains the Kozlov class) and for which almost periodic correctors exist. Their approach combines the regularity theory developed by S. Armstrong and C. Smart in [49] and adapted to the almost periodic setting by S. Armstrong and Z. Shen [48], a new quantification of almost-periodicity, and a sensitivity calculus in the spirit of [3].

### 7.3.2. *Optimal stochastic integrability in stochastic homogenization*

In [40] A. Gloria and F. Otto consider uniformly elliptic coefficient fields that are randomly distributed according to a stationary ensemble of a finite range of dependence. They show that the gradient and flux $(\nabla\phi, a(\nabla\phi + e))$ of the corrector $\phi$, when spatially averaged over a scale $R \gg 1$ decay like the CLT scaling $R^{-d/2}$. They establish this optimal rate on the level of *sub-Gaussian* bounds in terms of the stochastic integrability, and also establish a suboptimal rate on the level of optimal Gaussian bounds in terms of the stochastic integrability. The proof unravels and exploits the self-averaging property of the associated semigroup, which provides a natural and convenient disintegration of scales, and culminates in a propagator estimate with strong stochastic integrability. As an application, they characterize the fluctuations of the homogenization commutator, and prove sharp bounds on the spatial growth of the corrector, a quantitative two-scale expansion, and several other estimates of interest in homogenization.

### 7.3.3. *A theory of fluctuations in stochastic homogenization*

In [39], M. Duerinckx, A. Gloria, and F. Otto establish a path-wise theory of fluctuations in stochastic homogenization of linear elliptic equations in divergence form. More precisely they consider the model problem of a discrete equation with independent and identically distributed conductances (as considered in [27]). They identify a single quantity, which they call the homogenization commutator, that drives the fluctuations in stochastic homogenization in the following sense. On the one hand, this tensor-valued stationary random field satisfies a functional central limit theorem, and (when suitably rescaled) converges to a Gaussian white noise. On the other hand, the fluctuations of the gradient of the corrector, the fluctuations of the flux of the corrector, and the fluctuations of any solution of the PDE with random coefficients and localized right-hand side are characterized at leading order by the fluctuations of this homogenization commutator in a path-wise sense. As a consequence, when properly rescaled, the solution satisfies a functional central limit theorem, the gradient of the corrector converges to the Helmholtz projection of a Gaussian white noise, and the flux of the corrector converges to the Leray projection of the same white noise. Compared to previous contributions, our approach, based on the homogenization commutator, unravels the complete structure of fluctuations. It holds in any dimension $d \geq 2$, yields the first path-wise results, quantifies the limit theorems in Wasserstein distance, and only relies on arguments that extend to the continuum setting and to the case of systems.

## 7.4. Numerical methods for evolution equations

In [36] G. Dujardin analyzes an exponential integrator applied to the nonlinear Schrödinger equation with white noise dispersion. This models appears in optic fibers. Together with his co-author, he proves that this explicit scheme applied to the sctochastic PDE is of mean-square order 1. He uses it to illustrate a conjecture on the well-posedness of the equation in some regimes of the nonlinearity. Comparisons with several other schemes of the litterature are proposed. A last, another new (implicit) exponential integrators is proposed, which preserves the $L^2$-norm of the solution and is compared with the explicit one introduced beforehand.

## 7.5. Schrödinger equations

### 7.5.1. *Nonlinear optical fibers*

S. Rota Nodari, G. Dujardin, S. De Bièvre and collaborators continued their previous work on periodically modulated optical fibers with the experimental physicists of PhLAM [19]. They show that the nonlinear stage of modulational instability induced by parametric driving in the *defocusing* nonlinear Schrödinger equation can be accurately described by combining mode truncation and averaging methods, valid in the strong driving regime. The resulting integrable oscillator reveals a complex hidden heteroclinic structure of the instability. A remarkable consequence, validated by the numerical integration of the original model, is the existence of breather solutions separating different Fermi-Pasta-Ulam recurrent regimes.

In [42] S. de Bièvre and G. Dujardin analyze the formation of the Kuznetsov-Ma soliton of the 1D Schrödinger equation in the presence of periodic modulation satisfying an integrability condition. They show that this particular soliton has several compression points, the number, position and shape of which are controlled by the amplitude and the frequency of the modulation. They analyze the interplay between the frenquency of the soliton and the frequency of the modulation. Moreover, they show that one can suppress any component of the output spectrum of the soliton by a suitable choice of the amplitude and frequency of the modulation.

These works are part of the activities developped in the LabEx CEMPI.

### 7.5.2. *Nonlinear Schrödinger equations*

In [54], D. Bonheure, J.-B. Casteras and R. Nascimento obtained new results on the existence and qualitative properties of waveguides for a mixed-diffusion NLS. In particular, they proved the first existence results for waveguides with fixed mass and provided several qualitative descriptions of these.

S. De Bièvre and S. Rota Nodari continued their work on orbital stability of relative equilibria of Hamiltonian dynamical systems on Banach spaces, with a second paper [37], dealing with the situation where multi-dimensional invariance groups are present in the systems considered. They present a generalization of the Vakhitov-Kolokolov slope condition to this higher dimensional setting, and show how it allows to prove the local coercivity of the Lyapunov function, which in turn implies orbital stability. The method is applied to study the orbital stability of the plane waves of a system of two coupled nonlinear Schrödinger equations. They provide a comparison of their approach to the classical one by Grillakis-Shatah-Strauss.

## 7.6. Miscellaneous results

In [21] Mitia Duerinckx establishes the global well-posedness of a family of equations, which are obtained in certain regimes — in a joint work in preparation with Sylvia Serfaty — as the mean-field evolution of the supercurrent density in a (2D section of a) type-II superconductor with pinning and with imposed electric current. General vortex-sheet initial data are also considered, and the uniqueness and regularity properties of the solution are investigated.

In [33], [8], [11], [12], D. Bonheure, J.-B. Casteras and collaborators made bifurcation analysis and constructed multi-layer solutions of the Lin-Ni-Takagi and Keller-Segel equations, which come from the Keller-Segel system of chemotaxis in specific cases. A remarkable feature of the results is that the layers do not accumulate to the boundary of the domain but satisfy an optimal partition problem contrary to the previous type of solutions constructed for these models.

In [16], [17], [35], J.-B. Casteras and collaborators study different problems related to the existence of constant mean curvature hypersurfaces with prescribed asymptotic boundary on Cartan-Hadamard manifold. In particular, they obtained the first existence results for minimal graphs with prescribed asymptotic Dirichlet data under a pointwise pinching condition for sectionals curvatures.

S. De Bièvre and co-workers present in [67] a general approach to calculating the entanglement of formation for superpositions of two-mode coherent states, placed equidistantly on a circle in phase space. In the particular case of rotationally-invariant circular states the value of their entanglement is shown to be given by analytical expressions. They analyse the dependence of the entanglement on the radius of the circle and number of components in the superposition.

A. Benoit continues his analysis of hyperbolic equations in corner spaces. He addresses in [30] the rigorous construction of geometric optics expansions for weakly well-posed hyperbolic corner problems. He studies in [31] the semi-group stability for finite difference discretizations of hyperbolic systems of equations in corner domains, extending previous results of Coulombel & Gloria and Coulombel in the case of the halfspace.

<span style="color:red">**MINT Project-Team**</span>

# 6. New Results

## 6.1. ControllAR: Appropriation of visual Feedback on Control Surfaces

Florent Berthaut, Alex Jones

Despite the development of touchscreens, many expert systems for working with digital multimedia content, such as in music composition and performance, video editing or visual performance, still rely on control surfaces. This can be due to the accuracy and appropriateness of their sensors, the haptic feedback that they offer, and most importantly the way they can be adapted to the specific subset of gestures and tasks that users need to perform. On the other hand, visual feedback on controllers remains limited and/or fixed, preventing similar personalizing. In this paper, we propose ControllAR, a novel system that facilitates the appropriation of rich visual feedback on control surfaces through remixing of graphical user interfaces and augmented reality display. We then use our system to study current and potential appropriation of visual feedback in the case of digital musical instruments and derive guidelines for designers and developers.



*Figure 1. ControllAR: (left) ControllAR is used to augment a control surface with the remixed graphical user interface of music software, (right) Visual feedback designed by electronic musicians during our study belong to three categories: mappings feedback, processes feedback and content feedback.*

## 6.2. Talaria: Continuous Drag & Drop on a Wall Display

Hanaë Rateau, Yosra Rekik, Laurent Grisoni, Joaquim Jorge

We present an interaction technique combining tactile actions and Midair pointing to access out-of-reach content on large displays without the need to walk across the display. Users can start through a Touch gesture on the display surface and finish Midair by pointing to push content away or inversely to retrieve a content. The technique takes advantage of wellknown semantics of pointing in human-to-human interaction. These, coupled with the semantics of proximal relations and deictic proxemics make the proposed technique very powerful as it leverages on well-understood human-human interaction modalities. Experimental results show this technique to outperform direct tactile interaction on dragging tasks. From our experience we derive four guidelines for interaction with large-scale displays.

## 6.3. Multi fngers interaction on a surface haptic display

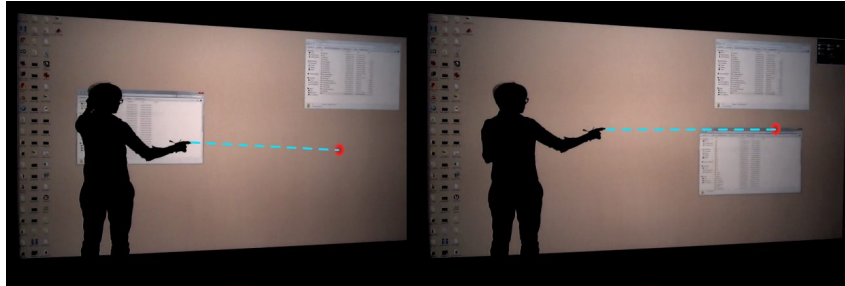Sofiane Ghenna, Christophe Giraud-Audine, Michel Amberg, Frédéric Giraud, Betty Lemaire-Semail

*Figure 2. Talaria*

In this study, we develop and implement a method for superimposing two vibration modes in order to produce different tactile stimuli on two fingers located in different positions. The tactile stimulation is based on the squeeze film effect which decreases the friction between a fingertip and a vibrating plate.

Experimental test have been conducted on a 1D tactile device. They show that it is possible to continuously control the friction on two fingers moving independently. Then, we developed the design of a 2D device based on the same principle, which gives rise to the design of a two fingers tactile display. Evaluations were conducted using a modal analysis with experimental validation.
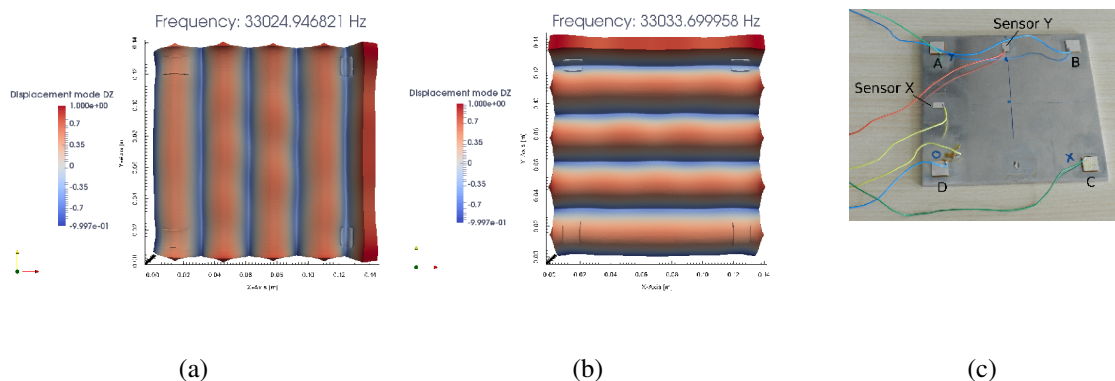


(a)                                      (b)                                      (c)

*Figure 3. Vibration modes and mode shapes using FEM with the position of the actuators (in white in a and b), and the prototype (c).*

## 6.4. Finding the Minimum Perceivable Size of a Tactile Element on an Ultrasonic Based Haptic Tablet

Farzan Kalantari, Laurent Grisoni, Frédéric Giraud, Yosra Rekik

Tactile devices with ultrasonic vibrations (based on squeeze film effect) using piezoelectric actuators are one of the existing haptic feedback technologies. In this study we have performed two psychophysical experiments on an ultrasonic haptic tablet, in order to find the minimum size of a tactile element on which all the users are able to perfectly identify different types of textures. Our results show that the spatial resolution of the tactile element on haptic touchscreen actually varies, depending on the number and types of tactile feedback information. A first experiment exhibits three different tactile textures, chosen as being easily recognized by

users. We use these textures in a second experiment, and evaluate minimal spatial area on which the chosen set of textures can be recognized. Among other, we find the minimal size depends on the texture nature.

## 6.5. BOEUF: A Unified Framework for Modeling and Designing Digital Orchestras

Florent Berthaut, Luke Dahl, Patricia Plénacoste

Orchestras of Digital Musical Instruments (DMIs) enable new musical collaboration possibilities, extending those of acoustic and electric orchestras. However the creation and development of these orchestras remain constrained. In fact, each new musical collaboration system or orchestra piece relies on a fixed number of musicians, a fixed set of instruments (often only one), and a fixed subset of possible modes of collaboration. In this paper, we describe a unified framework that enables the design of Digital Orchestras with potentially different DMIs and an expand-able set of collaboration modes. It relies on research done on analysis and classification of traditional and digital orchestras, on research in Collaborative Virtual Environments, and on interviews of musicians and composers. The BOEUF framework consists of a classification of modes of collaboration and a set of components for modelling digital orchestras. Integrating this framework into DMIs will enable advanced musical collaboration modes to be used in any digital orchestra, including spontaneous jam sessions.

Current work on this project consists in the implementation of BOEUF in the PureData programming language and in the study of its impact on musical collaboration during short improvised jam sessions.

<div align="center">

## Mjolnir Team

</div>

# 7. New Results

## 7.1. Introduction

The following sections summarize our main results of the year. For a complete list, see the list of publications at the end of this report.

## 7.2. Understanding and modeling users

**Participants:** Géry Casiez, Christian Frisson, Alix Goguey, Stéphane Huot, Sylvain Malacria, Mathieu Nancel, Thibault Raffaillac, Nicolas Roussel.

### 7.2.1. Touch interaction with finger identification: which finger(s) for what?

The development of robust methods to identify which finger is causing each touch point, called "finger identification," will open up a new input space where interaction designers can associate system actions to different fingers [11]. However, relatively little is known about the performance of specific fingers as single touch points or when used together in a "chord". We presented empirical results for accuracy, throughput, and subjective preference gathered in five experiments with 48 participants exploring all 10 fingers and 7 two-finger chords. Based on these results, we developed design guidelines for reasonable target sizes for specific fingers and two-finger chords, and a relative ranking of the suitability of fingers and two-finger chords for common multi-touch tasks. Our work contributes new knowledge regarding specific finger and chord performance and can inform the design of future interaction techniques and interfaces utilizing finger identification [28].

### 7.2.2. Training and use of brain-computer interfaces

Brain-Computer Interfaces (BCIs) are much less reliable than other input devices, with error rates ranging from 5% up to 60%. To assess the subjective frustration, motivation, and fatigue of users when confronted to different levels of error rate, we conducted an BCI experiment in which it was artificially controlled. Our results show that a prolonged use of BCI significantly increases the perceived fatigue, and induces a drop in motivation [38]. We also found that user frustration increases with the error rate of the system but this increase does not seem critical for small differences of error rate. For future BCIs, we thus advise to favor user comfort over accuracy when the potential gain of accuracy remains small.

We have also investigated if the stimulation used for training an SSVEP-based BCI have to be similar to the one used *in fine* for interaction. We recorded 6-channels EEG data from 12 subjects in various conditions of distance between targets, and of difference in color between targets. Our analysis revealed that the stimulation configuration used for training which leads to the best classification accuracy is not always the one which is closest to the end use configuration [15]. We found that the distance between targets during training is of little influence if the end use targets are close to each other, but that training at far distance can lead to a better accuracy for far distance end use. Additionally, an interaction effect is observed between training and testing color: while training with monochrome targets leads to good performance only when the test context involves monochrome targets as well, a classifier trained on colored targets can be efficient both for colored and monochrome targets. In a nutshell, in the context of SSVEP-based BCI, training using distant targets of different colors seems to lead to the best and more robust performance in all end use contexts.

### 7.2.3. Evaluation metrics for touch latency compensation

Touch systems have a delay between user input and corresponding visual feedback, called input "latency" (or "lag"). Visual latency is more noticeable during continuous input actions like dragging, so methods to display feedback based on the most likely path for the next few input points have been described in research papers and patents. Designing these "next-point prediction" methods is challenging, and there have been no standard metrics to compare different approaches. We introduced metrics to quantify the probability of 7 spatial error "side-effects" caused by next-point prediction methods [35]. Types of side-effects were derived using a thematic analysis of comments gathered in a 12 participants study covering drawing, dragging, and panning tasks using 5 state-of-the-art next-point predictors. Using experiment logs of actual and predicted input points, we developed quantitative metrics that correlate positively with the frequency of perceived side-effects. These metrics enable practitioners to compare next-point predictors using only input logs.

### 7.2.4. Application use in the real world

Interface designers, HCI researchers or usability experts often need to collect information regarding usage of interactive systems and applications in order to interpret quantitative and behavioral aspects from users – such as our study on the use of trackpads described before – or to provide user interface guidelines. Unfortunately, most existing applications are closed to such probing methods: source code or scripting support are not always available to collect and analyze users' behaviors in real world scenarios.

InspectorWidget [26] is an open-source cross-platform application we designed to track and analyze users' behaviors in interactive software. The key benefits of this application are: 1) it works with closed applications that do not provide source code nor scripting capabilities; 2) it covers the whole pipeline of software analysis from logging input events to visual statistics through browsing and programmable annotation; 3) it allows post-recording logging; and 4) it does not require programming skills. To achieve this, InspectorWidget combines low-level event logging (e.g. mouse and keyboard events) and high-level screen capturing and interpretation features (e.g. interface widgets detection) through computer vision techniques.

### 7.2.5. Trackpad use in the real world

*Trackpads* (or *touchpads*) allow to control an on-screen cursor with finger movements on their surface. Recent models also support force sensing and multi-touch interactions, which make it possible to scroll a document by moving two fingers or to switch between virtual desktops with four fingers, for example. But despite their widespread use, little is known about how users interact with them, and which gestures they are most familiar with. To better understand this, we conducted a three-steps field study with Apple Macbook's multi-touch trackpads.

The first step of our study consisted in collecting low-level interaction data such as contact points with the trackpad and the multi-touch gestures performed while interacting. We developed a dedicated interaction logging application that we deployed on the workstation of 11 users for a duration of 14 days, and collected a total of over 82 millions contact points and almost 220 000 gestures. We then investigated finger chords (i.e., fingers used) and hand usage when interacting with a trackpad. In that purpose, we designed a dedicated mirror stand that can be easily positioned in front of the laptop's embedded web camera to divert its capturing field (Figure 1 , left). This mirror stand is combined with a background application taking photos when a multi-finger gesture is performed. We deployed this setup on the computer of 9 users for a duration of 14 days. Finally, we deployed a system preference collection application to gather the trackpad system preferences (such as transfer function and gestures associated) of 80 users. Our main findings are that touch contacts on the trackpad are performed on a limited sub-surface and are relatively slow (Figure 1 , right); that the consistency of user finger chords varies depending on the frequency of a gesture and the number of fingers involved; and that users tend to rely on the default system preferences of the trackpad [34].

## 7.3. Interactive visualization and animations

**Participants:** Amira Chalbi-Neffati, Fanny Chevalier, Nicolas Roussel.

*Figure 1. Left: Mirror positioned in front of a built-in camera to divert its field of view and analyze finger-chords usages; Right: frequency distribution of all touch events of our participants, overlaid on a trackpad.*

### 7.3.1. Social network analysis

The egocentric analysis of dynamic networks focuses on discovering the temporal patterns of a subnetwork around a specific central actor, i.e. an ego-network. These types of analyses are useful in many application domains, such as social science and business intelligence, providing insights about how the central actor interacts with the outside world. *EgoLines* is an interactive visualization we designed to support the egocentric analysis of dynamic networks . Using a "subway map" metaphor, a user can trace an individual actor over the evolution of the ego-network (Figure 2 ). The design of EgoLines is grounded in a set of key analytical questions pertinent to egocentric analysis, derived from interviews with three domain experts and general network analysis tasks. The results of controlled experiments with end-users and domain experts show its effectiveness in egocentric analysis tasks. Egolines can be tested at http://fannychevalier.net/egolines.html
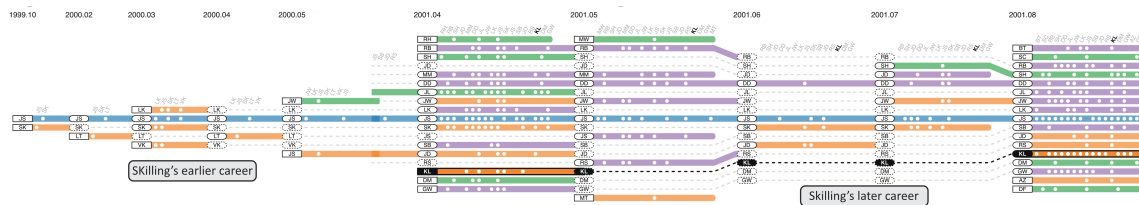


*Figure 2. Egolines used to explore the dynamic network of email communications among employees at the Enron company.*

### 7.3.2. Cross-sectional cohort phenotype

Cross-sectional phenotype studies are used by genetics researchers to better understand how phenotypes vary across patients with genetic diseases, both within and between cohorts. Analyses within cohorts identify patterns between phenotypes and patients, e.g. co-occurrence, and isolate special cases, e.g. potential outliers). Comparing the variation of phenotypes between two cohorts can help distinguish how different factors affect disease manifestation, e.g. causal genes, age of onset.). *PhenoStacks* is a novel visual analytics tool we designed to support the exploration of phenotype variation within and between cross-sectional patient cohorts . By leveraging the semantic hierarchy of the Human Phenotype Ontology, phenotypes are presented in context, can be grouped and clustered, and are summarized via overviews to support the exploration of

phenotype distributions (Figure 3 ). The HPO is rarely used for visualization and was only recently first employed in PhenoBlocks [49]. In this prior work, we used the HPO to visualize phenotypes in clinical diagnosis settings, supporting the pairwise comparison of patient phenotypes using explicit encoding. In this new work, we turn our focus to genetics researchers conducting cross-sectional cohort studies, where the distribution of phenotypes is compared across many patients. The design of PhenoStacks was motivated by formative interviews with genetics researchers. The results of a deployment evaluation with four expert genetics researchers suggest that PhenoStacks can help identify phenotype patterns, investigate data quality issues, and inform data collection design. PhenoStacks is available from http://phenostacks.org/
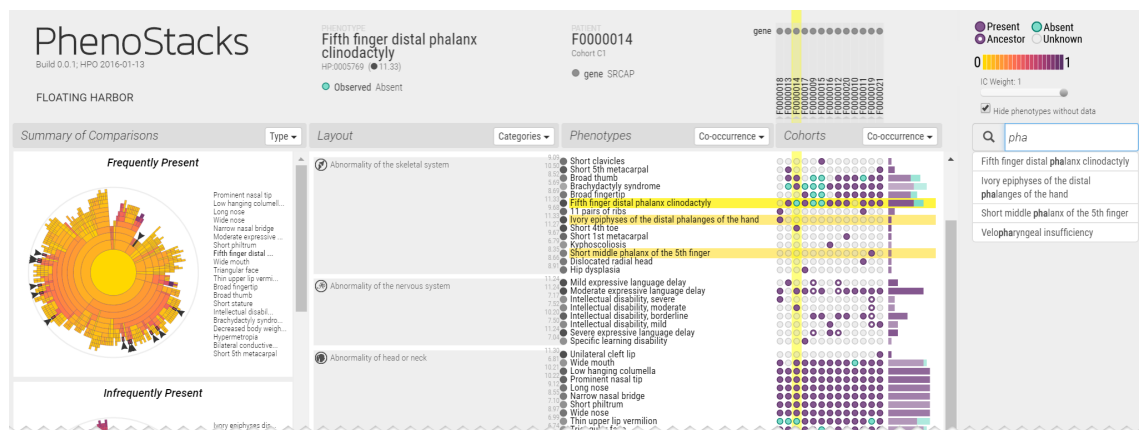


*Figure 3. Exploration of phenotypic variation in cross-sectional cohorts of patients with a rare genetic disease using PhenoStacks.*

### 7.3.3. *Human routine behavior*

Human routines are blueprints of behavior, which allow people to accomplish purposeful repetitive tasks at many levels, ranging from the structure of their day to how they drive through an intersection. People express their routines through actions that they perform in the particular situations that triggered those actions. An ability to model routines and understand the situations in which they are likely to occur could allow technology to help people improve their bad habits, inexpert behavior, and other suboptimal routines. However, existing routine models do not capture the causal relationships between situations and actions that describe routines. Byproducts of an existing activity prediction algorithm can be used to model those causal relationships in routines [23]. We applied this algorithm on two example datasets, and showed that the modeled routines are meaningful — that they are predictive of people's actions and that the modeled causal relationships provide insights about the routines that match findings from previous research. Our approach offers a generalizable solution to model and reason about routines. We show that the extracted routine patterns are at least as predictive of behaviors in the two behavior logs as the baseline we establish with existing algorithms.

To make the routine behavior models created using our approach accessible to participants and allow them to investigate the extracted routine patterns, we developed a simple visualization tool. To maintain a level of familiarity, we base our visual encoding of routine behavior elements on a traditional visual representation of an MDP as a graph (Figure 4 ). Our MDP graph contains nodes representing states (as circles) and actions (as squares), directed edges from state nodes to action nodes (indicating possible actions people can perform in those states), and directed edges from actions to states (indicating state transitions for any given state and action combination).
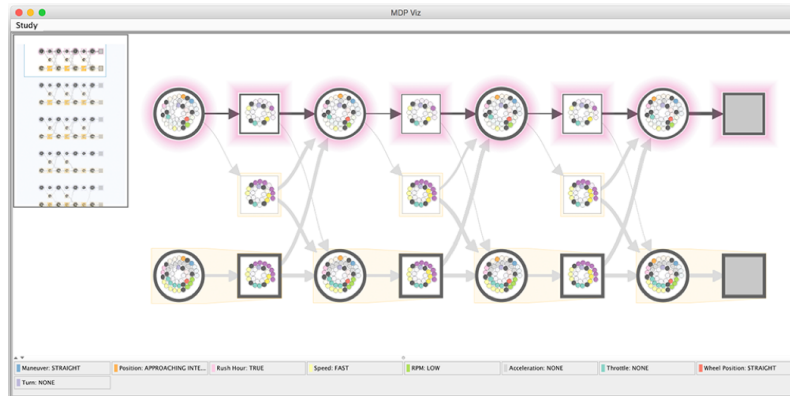
*Figure 4. Our visual analytics tool showing the main routine and one likely variation of non-aggressive drivers extracted using our approach.*

### 7.3.4. *Meta-analysis of data based on user-authored annotations*

User-authored annotations of data can support analysts in the activity of hypothesis generation and sense-making, where it is not only critical to document key observations, but also to communicate insights between analysts. *Annotation Graphs* are dynamic graph visualizations that enable meta-analysis of data based on user-authored annotations . The annotation graph topology encodes annotation semantics, which describe the content of and relations between data selections, comments, and tags. We present a mixed-initiative approach to graph layout that integrates an analyst's manual manipulations with an automatic method based on similarity inferred from the annotation semantics. Annotation graphs are implemented within a system, C8, that supports authoring annotations during exploratory analysis of a dataset (Figure 5 ). In this work, we develop and evaluate the system through an iterative user-centered design process with three experts, situated in the domain of analyzing HCI experiment data. The results suggest that annotation graphs are effective as a method of visually extending user-authored annotations to data meta-analysis for discovery and organization of ideas.

### 7.3.5. *Fundamentals of animated transitions*

Animations are increasingly used in interactive systems in order to enhance the usability and aesthetics of user interfaces. While animations are proven to be useful in many cases, we still find defective ones causing many problems, such as distracting users from their main task or making data exploration slower. The fact that such animations still exist proves that animations are not yet very well understood as a cognitive aid, and that we have not yet definitely decided what makes a well designed one. Our work on this topic aims at better understanding the different aspects of animations for user interfaces and exploring new methods and guidelines for designing them.

From bouncing icons that catch attention, to transitions helping with orientation, to tutorials, animations can serve numerous purposes. In , we revisit Baecker and Small's pioneering work *Animation at the Interface*, 25 years later. We review academic publications and commercial systems, and interviewed 20 professionals of various backgrounds. Our insights led to an expanded set of roles played by animation in interfaces today for keeping in context, teaching, improving user experience, data encoding and visual discourse. We illustrate each role with examples from practice and research, discussed evaluation methods and point to opportunities for future research. This expanded description of roles aims at inspiring the HCI research community to find novel uses of animation, guide them towards evaluation and spark further research.

We have also studied different aspects of animations for visual analysis tasks. We have worked on the design of a new model for animated transitions, explored certain aspects of visual grouping for these transitions, and
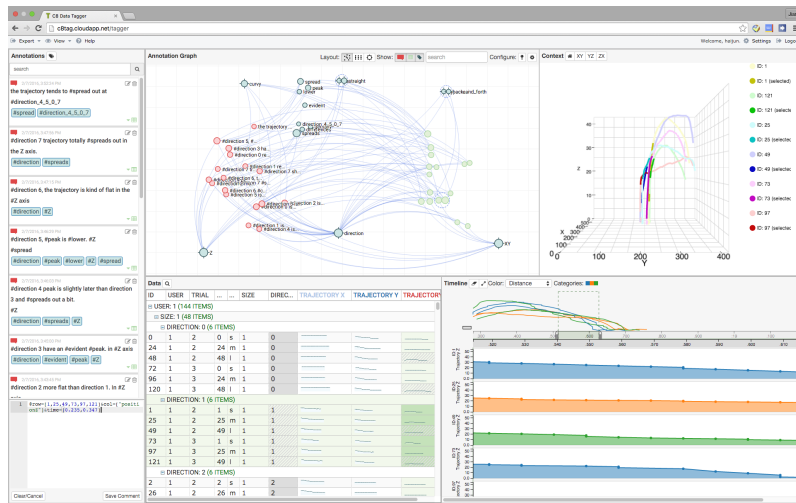
*Figure 5. C8 used to analyze the results of an HCI user study that records participants pointing at a target on a tabletop display with different experimental conditions.*

studied the impact of their temporal structure on data interpretation. These works, while still in progress, have been presented at the IHM doctoral consortium [39].

## 7.4. Interaction techniques

**Participants:** Géry Casiez, Fanny Chevalier, Stéphane Huot, Sylvain Malacria, Justin Mathew, Thomas Pietrzak, Nicolas Roussel.

### 7.4.1. Interaction in 3D environments

In virtual environments, interacting directly with our hands and fingers greatly contributes to the sense of immersion, especially when force feedback is provided for simulating the touch of virtual objects. Yet, common haptic interfaces are unfit for multi-finger manipulation and only costly and cumbersome grounded exoskeletons do provide all the efforts expected from object manipulation. To make multi-finger haptic interaction more accessible, we propose to combine two affordable haptic interfaces into a bimanual setup named DesktopGlove [18]. With this approach, each hand is in charge of different components of object manipulation: one commands the global motion of a virtual hand while the other controls its fingers for grasping. In addition, each hand is subjected to forces that relate to its own degrees of freedom so that users perceive a variety of haptic effects through both of them. Our results show that (1) users are able to integrate the separated degrees of freedom of DesktopGlove to efficiently control a virtual hand in a posing task, (2) DesktopGlove shows overall better performance than a traditional data glove and is preferred by users, and (3) users considered the separated haptic feedback realistic and accurate for manipulating objects in virtual environments.

We also investigated how head movements can serve to change the viewpoint in 3D applications, especially when the viewpoint needs to be changed quickly and temporarily to disambiguate the view. We studied how to use yaw and roll head movements to perform orbital camera control, i.e., to rotate the camera around a specific point in the scene [33]. We reported on four user studies. Study 1 evaluated the useful resolution of head movements and study 2 informed about visual and physical comfort. Study 3 compared two interaction techniques, designed by taking into account the results of the two previous studies. Results show that head roll is more efficient than head yaw for orbital camera control when interacting with a screen. Finally, Study

4 compared head roll with a standard technique relying on the mouse and the keyboard. Moreover, users were allowed to use both techniques at their convenience in a second stage. Results show that users prefer and are faster (14.5%) with the head control technique.

### 7.4.2. Storyboard sketching for stereo 3D films and Virtual Reality stories

The resurgence of stereoscopic and Virtual Reality (VR) media has motivated filmmakers to evolve new stereo- and VR-cinematic vocabularies, as many principles for stereo 3D film and VR story are unique. Concepts like plane separation, parallax position, and depth budgets in stereo, and presence, active experience, blocking and stitching in VR are missing from early planning due to the 2D nature of existing storyboards. Motivated to foresee difficulties exclusive to stereoscopy and VR, but also to exploit the unique possibilities of these medium, the 3D and VR cinematography communities encourages filmmakers to start thinking in stereo/VR as early as possible. Yet, there are very few early stage tools to support the ideation and discussion of a stereoscopic film or a VR story. Traditional solutions for early visual development and design, in current practices, are either strictly 2D or require 3D modeling skills, producing content that is consumed passively by the creative team.

To fill the gap in the filmmakers' toolkit, we proposed *Storeoboard* [31], a system for stereo-cinematic conceptualization, via storyboard sketching directly in stereo (Figure 6 ); and a novel multi-device system supporting the planning of virtual reality stories. Our tools are the first of their kind, allowing filmmakers to explore, experiment and conceptualize ideas in stereo or VR early in the film pipeline, develop new stereo- and VR-cinematic constructs and foresee potential difficulties. Our solutions are the design outcome of interviews and field work with directors, stereographers, storyboard artists and VR professionals. Our core contributions are thus: 1) a principled approach to the design and development of the first stereoscopic storyboard system that allows the director and artists to explore both the stereoscopic space and concepts in real-time, addressing key HCI challenges tied to sketching in stereoscopy; and 2) a principled survey of the state of the art in cinematic VR planning to design the first multi-device system that supports a storyboard workflow for VR film. We evaluated our tools with focus group and individual user studies with storyboard artists and industry professionals. In [31], we also report on feedback from the director of a live action, feature film on which Storeoboard was deployed. Results suggest that our approaches provide the speed and functionality needed for early stage planning, and the artifacts to properly discuss steroscopic and VR films.



*Figure 6. Storeoboard augments sketch-based storyboards with stereoscopic 3D planes for a fluid and flexible authoring of stereoscopic storyboards.*

### 7.4.3. Tactile displays and vibrotactile feedback

Tactile displays have predominantly been used for information transfer using patterns or as assistive feedback for interactions. With recent advances in hardware for conveying increasingly rich tactile information that mirrors visual information, and the increasing viability of wearables that remain in constant contact with the skin, there is a compelling argument for exploring tactile interactions as rich as visual displays. As Direct

Manipulation underlies much of the advances in visual interactions, we introduced *Direct Manipulation-enabled Tactile display* [29]. We defined the concepts of a tactile screen, tactile pixel, tactile pointer, and tactile target which enable tactile pointing, selection and drag & drop. We built a proof of concept tactile display and studied its precision limits. We further developed a performance model for DMTs based on a tactile target acquisition study, and studied user performance in a real-world DMT menu application. The results show that users are able to use the application with relative ease and speed.

We have also explored vibrotactile feedback with wearable devices such as smartwatches and activity trackers, which are becoming prevalent. These devices provide continuous information about health and fitness, and offer personalized progress monitoring, often through multimodal feedback with embedded visual, audio, and vibrotactile displays. Vibrations are particularly useful when providing discreet feedback, without users having to look at a display or anyone else noticing, thus preserving the flow of the primary activity. Yet, current use of vibrations is limited to basic patterns, since representing more complex information with a single actuator is challenging. Moreover, it is unclear how much the user's current physical activity may interfere with their understanding of the vibrations. We addressed both issues through the design and evaluation of ActiVibe, a set of vibrotactile icons designed to represent progress through the values 1 to 10 [24]. We demonstrate a recognition rate of over 96% in a laboratory setting using a commercial smartwatch. ActiVibe was also evaluated in situ with 22 participants for a 28-day period. We show that the recognition rate is 88.7% in the wild and give a list of factors that affect the recognition, as well as provide design guidelines for communicating progress via vibrations.

### 7.4.4. *Force-based autoscroll*

Autoscroll, also known as edge-scrolling, is a common interaction technique in graphical interfaces that allows users to scroll a viewport while in dragging mode: once in dragging mode, the user moves the pointer near the viewport's edge to trigger an "automatic" scrolling. In spite of their wide use, existing autoscroll methods suffer from several limitations  [45]. First, most autoscroll methods over-rely on the size of the control area, that is, the larger it is, the faster scrolling rate can be. Therefore, the level of control depends on the available distance between the viewport and the edge of the display, which can be limited. This is for example the case with small displays or when the view is maximized. Second, depending on the task, the users' intention can be ambiguous (e.g. dragging and dropping a file is ambiguous as the user's target may be located within the initial viewport or in a different one on the same display). To reduce this ambiguity, the size of the control area is drastically smaller for drag-and-drop operations which consequently also affects scrolling rate control as the user has a limited input area to control the scrolling speed.

We explored how force-sensing input, which is now available on commercial devices such as the Apple Magic Trackpad 2 or iPhone 6S, can be used to overcome the limitations of autoscroll. Indeed, force-sensing is an interesting candidate because: 1) users are often already applying a (relatively soft) force on the input device when using autoscroll and 2) varying force on the input device does not require to move the pointer, thus making it possible to offer control to the user while using a small and consistent control area regardless of the task and the device. We designed and proposed ForceEdge, a novel interaction technique mapping the force applied on a trackpad to the autoscrolling rate [19]. We implemented a software interface that can be used to design different transfer functions that map the force to autoscrolling rate and test these mappings for text selection and drag-and-drop tasks. Our pilot studies showed encouraging results and future work will focus on conducting more robust evaluations, as well as testing ForceEdge on mobile devices.

### 7.4.5. *Combined Brain and gaze inputs for target selection*

Gaze-based interfaces and Brain-Computer Interfaces (BCIs) allow for hands-free human–computer interaction. We investigated the combination of gaze and BCIs and proposed a novel selection technique for 2D target acquisition based on input fusion. This new approach combines the probabilistic models for each input, in order to better estimate the intent of the user. We evaluated its performance against the existing gaze and brain–computer interaction techniques. Twelve participants took part in our study, in which they had to search and select 2D targets with each of the evaluated techniques. Our fusion-based hybrid interaction technique

was found to be more reliable than the previous gaze and BCI hybrid interaction techniques for 10 participants over 12, while being 29% faster on average. However, similarly to what has been observed in hybrid gaze-and-speech interaction, gaze-only interaction technique still provides the best performance. Our results should encourage the use of input fusion, as opposed to sequential interaction, in order to design better hybrid interfaces [14].

### 7.4.6. Actuated desktop devices

Desktop workstation remains the most common setup for office work tasks such as text editing, CAD, data analysis or programming. While several studies investigated how users interact with their devices (e.g. pressing keyboard keys, moving the cursor, etc.), it is not clear how they arrange their devices on the desk and whether we can leverage existing users' behaviors.

We designed the LivingDesktop [22], an augmented desktop with devices capable of moving autonomously. The LivingDesktop can control the position and orientation of the mouse, keyboard and monitors, offering different degrees of control for both the system (autonomous, semi-autonomous) and the user (manual, semi-manual) as well as different perceptive qualities (visual, haptic) thanks to a large range of device motions. We implemented a proof-of-concept of the LivingDesktop combining rail, robotic base and magnetism to control the position and orientation of the devices. This new setup presents several interesting features: (1) it improves ergonomics by continuously adjusting the position of its devices to help users adopting ergonomic postures and avoiding static postures for extended periods; (2) it facilitates collaborative works between local (e.g. located in the office) and remote co-workers; (3) it leverages context by reacting to the position of the user in the office, the presence of physical objects (e.g. tablets, food) or users' current activity in order to maintain a high level of comfort; (4) it reinforces physicality within the desktop workstation to increase immersion.

We conducted a scenario evaluation of the LivingDesktop. Our results showed the perceived usefulness of collaborative and ergonomics applications, as well as how it inspired our participants to elaborate novel application scenario, including social communication or accessibility.

### 7.4.7. Latency compensation

Human-computer interactions are greatly affected by the latency between the human input and the system visual response and the compensation of this latency is an important problem for the HCI community. We have developed a simple forecasting algorithm for latency compensation in indirect interaction using a mouse, based on numerical differentiation. Several differentiators were compared, including a novel algebraic version, and an optimized procedure was developed for tuning the parameters of the algorithm. The efficiency was demonstrated on real data, measured with a 1ms sampling time. These results are developed in [37] and patent has been filed on a subsequent technique for latency compensation [42].

<span style="color:red">**MODAL Project-Team**</span>

# 7. New Results

## 7.1. An oracle inequality for Quasi-Bayesian Non-Negative Matrix Factorisation

**Participant:** Benjamin Guedj.

We have extended the quasi-Bayesian perspective to the popular setting of non-negative matrix factorisation. This is a pivotal problem in machine learning (image segmentation, recommendation systems, audio source separation, ...) and we were able to propose an original estimator of the unobserved matrix. An oracle inequality is derived, along with several possible implementations. This work is now submitted to an international journal [38].

Joint work with Pierre Alquier.

## 7.2. PAC-Bayesian Online Clustering

**Participants:** Benjamin Guedj, Le Li.

We have extended the PAC-Bayesian framework to online learning. Our algorithm (called PACBO) performs online clustering of random sequences, and is supported by strong theoretical (regret bounds) and algorithmic (ergodicity of an MCMC implementation) results. This work is now submitted to an international journal [46].

Joint work with Sébastien Loustau.

## 7.3. Simpler PAC-Bayesian Bounds for Hostile Data

**Participant:** Benjamin Guedj.

We have introduced an original and much simpler way of deriving PAC-Bayesian bounds, through the use of $f$-divergences (therefore generalizing earlier works on Renyi's divergence and Kullback-Leibler divergence). This work is now submitted to an international conference [39].

Joint work with Pierre Alquier.

## 7.4. Clustering categorical functional data: Application to medical discharge letters

**Participants:** Cristian Preda, Cristina Preda, Vincent Vandewalle.

Categorical functional data represented by paths of a stochastic jump process are considered for clustering. For paths of the same length, the extension of the multiple correspondence analysis allows the use of well-known methods for clustering finite dimensional data. When the paths are of different lengths, the analysis is more complex. In this case, for Markov models we have proposed an EM algorithm to estimate a mixture of Markov processes. This work has been presented in a workshop [48].

## 7.5. Simultaneous dimension reduction and multi-objective clustering

**Participant:** Vincent Vandewalle.

In model based clustering of quantitative data it is often supposed that only one clustering variable explains the heterogeneity of all the others variables. However, when variables come from different sources, it is often unrealistic to suppose that the heterogeneity of the data can only be explained by one variable. If such an assumption is made, this could lead to a high number of clusters which could be difficult to interpret. A model based multi-objective clustering is proposed, it assumes the existence of several latent clustering variables, each one explaining the heterogeneity of the data on some clustering projection. In order to estimate the parameters of the model an EM algorithm is proposed, it mainly relies on a reinterpretation of the standard factorial discriminant analysis in a probabilistic way. The obtained results are projections of the data on some principal clustering components allowing some synthetic interpretation of the principal clusters raised by the data. This work has been presented in a conference [49].

## 7.6. Spatial Prediction of solar energy

**Participant:** Sophie Dabo.

Sophie Dabo-Niang's new result concern a work on spatial prediction of solar Energy in collaboration with some physicians and is now published [15].

This paper introduces a new approach for the forecasting of solar radiation series at a located station for very short time scale. We built a multivariate model in using few stations (3 stations). The proposed model is a spatio temporal vector autoregressive VAR model specifically designed for the analysis of spatially sparse spatio-temporal data. This model differs from classic linear models in using spatial and temporal parameters where the available predictors are the lagged values at each station. A spatial structure of stations is defined by the sequential introduction of predictors in the model. Moreover, an iterative strategy in the process of our model will select the necessary stations removing the uninteresting predictors and also selecting the optimal p-order. We studied the performance of this model. The metric error, the relative root mean squared error (rRMSE), is presented at different short time scales. Moreover, we compared the results of our model to simple and well known persistence model and those found in literature.

## 7.7. Multiple change-point detection

**Participants:** Alain Celisse, Guillemette Marot.

This is a joint work with Morgane Pierre-Jean and Guillem Rigaill (Univ. Evry).

The paper related to the work described in previous MODAL team reports (sections Kernel change point) has been pursuied and made available on Arxiv [42]. For recall, this work focuses on the problem of detecting abrupt changes arising in the full distribution of the observations (not only in the mean or variance). It provides greatly improved algorithms in terms of computational complexity (both in time and space). The computational and statistical performances of these new algorithms have been assessed through empirical experiments, which are detailed in the preprint.

## 7.8. Differential gene expression analysis

**Participants:** Alain Celisse, Guillemette Marot.

The use of empirical Bayesian techniques implemented in the R package metaMA has enabled to better understand Waldenstrom's macroglobulinemia. The new findings in Biology have been published in [18].

## 7.9. New concentration inequalities for the leave-$p$-out CV estimator

**Participant:** Alain Celisse.

New concentration inequalities have been established for the leave-$p$-out cross-validation estimator applied to assess the performance the $k$-nearest neighbour binary classifier. Joint work with Tristan Mary-Huard.

## 7.10. A new notion of stability for learning algorithms

**Participants:** Alain Celisse, Benjamin Guedj.

We introduced a new notion of stability for learning algorithms, which bridges the gap between the earlier uniform and hypothesis stability notions. It allows us to derive new PAC exponential concentration inequalities that apply to the Ridge regression algorithm as a first step. The first version of this work is presented in the preprint [41] and is now an active line of research.

## 7.11. Model for conditionally correlated categorical data

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

It is a model-based clustering proposal (called CMM for Conditional Modes Model) 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 work is now published in the international journal Advances in Data Analysis and Classification (Marbac et al, 2016). Furthermore, an R package (CoModes) is available on Rforge.

## 7.12. 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 work has been now accepted to an international journal [21]. Furthermore, an R package (MixCluster) is available on Rforge.

## 7.13. Degeneracy in multivariate Gaussian mixtures (complete data case)

**Participant:** Christophe Biernacki.

In the case of Gaussian mixtures, unbounded likelihood is an important theoretical and practical problem. Using the weak information that the latent sample size of each component has to be greater than the space dimension, a simple non-asymptotic stochastic lower bound on variances is derived. It is proved also that maximizing the likelihood under this data-driven constraint leads to consistent estimates. This work has been presented as an invited talk to the international workshop [28] and a paper for an international journal is been prepared.

This is a joined work with Gwënaelle Castellan of University of Lille.

## 7.14. Degeneracy in multivariate Gaussian mixtures (missing data case)

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

In the case of multivariate Gaussian mixtures, unbounded likelihood is an important theoretical and practical problem. However, in the case of missing data situations, this drawback is exacerbated for too reasons. Firstly, degeneracy frequence increases with missing data occurrence. Secondly, the EM dynamic is hardly detected since it implies linear grows of the log-likelihood, contrary to exponential grows in the complete data case, leading to computation waste and also high risk of erroneous estimates. Using the weak information that the latent sample size of each component (restricted to complete data) has to be greater than the space dimension, it is derived a simple contraint EM algorithm variant allowing to solve simultaneously both problems. This work has been presented to the international workshop [28] and a paper for an international journal is been prepared.

## 7.15. Data units selection in statistics

**Participant:** Christophe Biernacki.

Usually, the data unit definition is fixed by the practitioner but it can happen that it hesitates between several data unit options. In this context, it is highlighted that it is possible to embed data unit selection into a classical model selection principle. The problem is introduced in a regression context before to focus on the model-based clustering and co-clustering context, for data of different kinds (continuous, categorical, counting, ...). It has led to an invitation to an international workshop [29] and a preprint is being to be prepared.

It is a joint work with Alexandre Lourme from University of Bordeaux.

## 7.16. Label switching in Bayesian mixture model estimation

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

In the case of mixtures of distributions, it is well-known that the Bayesian posterior distribution is invariant to label switching, it means invariant to any renumbering of components. Consequences are important, typically leading to unuseful estimates like the posterior mean. Many attempts exist to solve this problem but it is advocated in this work that such a quest should be unfruitful since it is a direct consequence of the label non-identifiability of mixtures themselves. The present work proposes an original way to manage the label switching problem based on the Gibbs algorithm dynamic. The basic idea is to control the label switching probability along Gibbs iterations, controlled by both the sample size and the component overlap. An early version of this work has been presented as an invited talk to the international workshop [28].

## 7.17. Trade-off computation time and accuracy

**Participants:** Christophe Biernacki, Maxime Brunin, Alain Celisse.

Most estimates practically arise from algorithmic processes aiming at optimizing some standard, but usually only asymptotically relevant, criteria. Thus, the quality of the resulting estimate is a function of both the iteration number and also the involved sample size. An important question is to design accurate estimates while saving computation time, and we address it in the simplified context of linear regression here. Fixing the sample size, we focus on estimating an early stopping time of a gradient descent estimation process aiming at maximizing the likelihood. It appears that the accuracy gain of such a stopping time increases with the number of covariates, indicating potential interest of the method in real situations involving many covariates. A first version of this work has been presented to an international conference [27], and a preprint is being in progress.

## 7.18. Projection under pairwise control

**Participant:** Christophe Biernacki.

Visualization of high-dimensional and possibly complex (non continuous for instance) data onto a low-dimensional space may be difficult. Several projection methods have been already proposed for displaying such high-dimensional structures on a lower-dimensional space, but the information lost is not always easy to use. Here, a new projection paradigm is presented to describe a non-linear projection method that takes into account the projection quality of each projected point in the reduced space, this quality being directly available in the same scale as this reduced space. More specifically, this novel method allows a straightforward visualization data in $R^2$ with a simple reading of the approximation quality, and provides then a novel variant of dimensionality reduction.

This work is under revision in an international journal [37] and it has also been presented to an international conference [25].

It is a joint work with Hiba Alawieh and Nicolas Wicker, both from University of Lille.

## 7.19. Matching of descriptors evolving over time

**Participants:** Christophe Biernacki, Anne-Lise Bedenel.

In the web domain, and in particular for insurance comparison, data constantly evolve, implying that it is difficult to directly exploit them. For example, to do a classification, performing standard learning processes require data descriptor equal for both learning and test samples. Indeed, for answering to web surfer expectation, online forms whence data come from are regularly modified. So, features and data descriptors are also regularly modified. In this work, it is introduced a process to estimate and understand connections between transformed data descriptors. This estimated matching between descriptors will be a preliminary step before applying later classical learning methods. This work has been presented to a national conference [33], with international audience.

It is a joint work with Laetitia Jourdan, from University of Lille and Inria.

## 7.20. Real-time audio sources classification

**Participants:** Christophe Biernacki, Maxime Baelde.

Recent research on machine learning focuses on audio source identification in complex environments. They rely on extracting features from audio signals and use machine learning techniques to model the sound classes. However, such techniques are often not optimized for a real-time implementation and in multi-source conditions. It is proposed here a new real-time audio single-source classification method based on a dictionary of sound models (that can be extended to a multi-source setting). The sound spectrums are modeled with mixture models and form a dictionary. The classification is based on a comparison with all the elements of the dictionary by computing likelihoods and the best match is used as a result. It is found that this technique outperforms classic methods within a temporal horizon of 0.5s per decision (achieved 6% of errors on a database composed of 50 classes). Future works will focus on the multi-sources classification and reduce the computational load. This work has been accepted in 2016 to be presented in 2017 to an international conference in Signal Processing [32].

It is a joint work with Raphaël Greff, from the A-Volute company.

## 7.21. Model-Based Co-clustering for Ordinal Data

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

A model-based coclustering algorithm for ordinal data is presented. This algorithm relies on the latent block model using the BOS model (Biernacki and Jacques, 2015, Stat. Comput.) for ordinal data and a SEM-Gibbs algorithm for inference. Numerical experiments on simulated data illustrate the efficiency of the inference strategy. This work has been presented to an international workshop [30] and also to a national conference with an international audience [35].

## 7.22. Computational and statistical trade-offs in change-point detection

**Participants:** Christophe Biernacki, Maxime Brunin, Alain Celisse.

The change-point detection problem aims to detect changes in the distribution of observations collected over the time between the instants 1,...,T in the offline context. These changes occur at some instants called change-points. Our method provides consistent estimates of the change-points obtained by the Kernel Binary Segmentation algorithm with stopping rule (KBS). Moreover, the proposed method has a lower complexity in time and in space than the Kernel Dynamic Programming (KDP). This work has been presented to a national conference with an international audience [34].

## 7.23. MixtComp software for full mixed data

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

MixtComp (Mixture Computation) is an integration software from the MODAL team for model-based clustering of mixed data. Its computing core is written in C++ and is accessed through an R interface. Its architecture allows to easily and quickly integrate new univariate models (under the conditional independence assumption) as they are published. The first phase of development was the implementation of three basic models (Gaussian, Multinomial, Poisson) with the native management of partially observed data (including intervals). It now implements models related to ordinal data (2015), rank data (2015) and functional data (2016), still with missing or partially missing data. The code is developed internally, and has been field-tested through several contracted partnerships (see the section about contracts). It is now referenced in the BIL database and the APP. It is available through a new web interface, called MASSICCC at https://massiccc.lille. inria.fr/#/ (see also the dedicated section). MixtComp has been presented to an invited talk in October 2016 at the Academy of Sciences in Tunisia [26].

## 7.24. MASSICCC platform for SaaS software availability

**Participants:**  Christophe Biernacki, Vincent Kubicki, Matthieu Marbac Lourdelle.

MASSICCC is a demonstration platform giving access through a SaaS (service as a software) concept to data analysis libraries developed at Inria. It allows to obtain results either directly through a website specific display (specific and interactive visual outputs) or through an R data object download. It started in October 2015 for two years and is common to the Modal team (Inria Lille) and the Select team (Inria Saclay). In 2016, two packages have been integrated: Mixmod and MixtComp (see the specific section about MixtComp). In 2017, it is planned to integrate the BlockCluster package. The MASSICCC platform gardually replaces the former BigStat platform available here: https://modal-research.lille.inria.fr/BigStat/. BigStat and MASSICCC have been both presented to an invited talk in October 2016 at the Academy of Sciences in Tunisia [26].

MASSICCC has led to a first short meeting in April 2016 in Lille for obtaining a feedback from company and academic users. Here is the link towards this event: Link. A second similar event is planned in February 2017 in Paris. Joint work with Jonas Renault and Josselin Demont (both at InriaTech).

The MASSICCC platform is available on https://massiccc.lille.inria.fr

## 7.25. CoModes package for correlated categorical variables

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

CoModes is an R package for model-based clustering of categorical data. In this package, the Conditional Modes Model (CMM), published in 2016 (Marbac et al, 2016), takes into account the main conditional dependencies between variables through particular modality crossings (so-called modes). CoModes performs the model selection and provides the best model according to the exact integrated likelihood criterion and the maximum likelihood estimates. It is available online on Rforge (https://r-forge.r-project.org/R/?group_id=1809).

## 7.26. MixCluster package for correlated mixed variables

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

MixCluster is an R package for model-based clustering of mixed data (continuous, binary, integer). In this package, the model, accepted for publication in 2016 [21], takes into account the main conditional dependencies between variables through Gaussian copula. Mixcluster performs the model selection and provides the best model according to Bayesian approaches. It is available online on Rforge (https://r-forge. r-project.org/R/?group_id=1939).

<center>

<span style="color:red">**NON-A Project-Team**</span>

</center>

# 7. New Results

## 7.1. Homogeneity Theory

Homogeneity is one of the tools we develop for finite-time convergence analysis. In 2016 this concept has received various improvements:

- Frequency domain approach to analysis of homogeneous nonlinear systems [85], [46]:

  Analysis of feedback sensitivity functions for implicit Lyapunov function-based control system is developed. The Gang of Four and loop transfer function are considered for practical implementation of the control via frequency domain control design. The effectiveness of this control scheme is demonstrated on an illustrative example of roll control for a vectored thrust aircraft.

- Homogeneous distributed parameter systems [72], [32]:

  A geometric homogeneity is introduced for evolution equations in a Banach space. Scalability property of solutions of homogeneous evolution equations is proven. Some qualitative characteristics of stability of trivial solution are also provided. In particular, finite-time stability of homogeneous evolution equations is studied. Classical theorems on existence and uniqueness of solutions of nonlinear evolution equations are revised. A universal homogeneous feedback control for a finite-time stabilization of linear evolution equation in a Hilbert space is designed using homogeneity concept. The design scheme is demonstrated for distributed finite-time control of heat and wave equations.

- Robustness of Homogenous Systems:
  - [93], [36] The problem of stability robustness with respect to time-varying perturbations of a given frequency spectrum is studied applying homogeneity framework. The notion of finite-time stability over time intervals of finite length, i.e. short-finite-time stability, is introduced and used for that purpose. The results are applied to demonstrate some robustness properties of the three-tank system.
  - The uniform stability notion for a class of nonlinear time-varying systems is studied in [35] using the homogeneity framework. It is assumed that the system is weighted homogeneous considering the time variable as a constant parameter, then several conditions of uniform stability for such a class of systems are formulated. The results are applied to the problem of adaptive estimation for a linear system.
  - Robustness with respect to delays is discussed in [84], [45] for homogeneous systems with negative degree. It is shown that if homogeneous system with negative degree is globally asymptotically stable at the origin in the delay-free case then the system is globally asymptotically stable with respect to a compact set containing the origin independently of delay. The possibility of applying the result for local analysis of stability for not necessary homogeneous systems is analyzed. The theoretical results are supported by numerical examples.

- Finite-time and Fixed-time Control and Estimation:
  - [61], [46] A switched supervisory algorithm is proposed, which ensures fixed-time convergence by commutation of finite-time or exponentially stable homogeneous systems of a special class, and a finite-time convergence to the origin by orchestrating among asymptotically stable systems. A particular attention is paid to the case of exponentially stable systems. Finite-time and fixed-time observation problem of linear multiple input multiple output (MIMO) control systems is studied. The nonlinear dynamic observers , which guarantee convergence of the observer states to the original system state in a finite and a fixed (defined a priori) time, are studied. Algorithms for the observers parameters tuning are also developed.

- [16] This paper focuses on the design of fixed-time consensus for multiple unicycle-type mobile agents. A distributed switched strategy, based on local information, is proposed to solve the leader-follower consensus problem for multiple nonholonomic agents in chained form. The switching times and the prescribed convergence time are explicitly given regardless of the initial conditions. Simulation results highlight the efficiency of the proposed method.

- Discretization of Homogeneous Systems:
  - [63] Sufficient conditions for the existence and convergence to zero of numeric approximations to solutions of asymptotically stable homogeneous systems are obtained for the explicit and implicit Euler integration schemes. It is shown that the explicit Euler method has certain drawbacks for the global approximation of homogeneous systems with non-zero degrees, whereas the implicit Euler scheme ensures convergence of the approximating solutions to zero.
  - [69] The known results on asymptotic stability of homogeneous differential inclusions with negative homogeneity degrees and their accuracy in the presence of noises and delays are extended to arbitrary homogeneity degrees. Discretization issues are considered, which include explicit and implicit Euler integration schemes. Computer simulation illustrates the theoretical results.

- Multi-Homogeneity and differential inclusions:
  - The notion of homogeneity in the bi-limit from is extended in [21] to local homogeneity and then to homogeneity in the multi-limit. The converse Lyapunov/Chetaev theorems on (homogeneous) system instability are obtained. The problem of oscillation detection for nonlinear systems is addressed. The sufficient conditions of oscillation existence for systems homogeneous in the multi-limit are formulated. The proposed approach estimates the number of oscillating modes and the regions of their location. Efficiency of the technique is demonstrated on several examples.
  - In [94], the notion of geometric homogeneity is extended for differential inclusions. This kind of homogeneity provides the most advanced coordinate-free framework for analysis and synthesis of nonlinear discontinuous systems. The main qualitative properties of continuous homogeneous systems are extended to the discontinuous setting: the equivalence of the global asymptotic stability and the existence of a homogeneous Lyapunov function; the link between finite-time stability and negative degree of homogeneity; the equivalence between attractivity and asymptotic stability are among the proved results.

## 7.2. Algebraic Technique For Estimation

- Time parameter estimation for a sum of sinusoidal waveform signals [39]:

  A novel algebraic method is proposed to estimate amplitudes, frequencies, and phases of a biased and noisy sum of complex exponential sinusoidal signals. The resulting parameter estimates are given by original closed formulas, constructed as integrals acting as time-varying filters of the noisy measured signal. The proposed algebraic method provides faster and more robust results, compared with usual procedures. Some computer simulations illustrate the efficiency of our method.

- Algebraic estimation via orthogonal polynomials [80]:

  Many important problems in signal processing and control engineering concern the reconstitution of a noisy biased signal. For this issue, we consider the signal written as an orthogonal polynomial series expansion and we provide an algebraic estimation of its coefficients. We specialize in Hermite polynomials. On the other hand, the dynamical system described by the noisy biased signal may be given by an ordinary differential equation associated with classical orthogonal polynomials. The signal may be recovered through the coefficients identification. As an example, we illustrate our algebraic method on the parameter estimation in the case of Hermite polynomials.

- An effective study of the algebraic parameter estimation problem [105]:

  Within the algebraic analysis approach, we first give a general formulation of the algebraic parameter estimation for signals which are defined by ordinary differential equations with polynomial coefficients such as the standard orthogonal polynomials (Chebyshev, Jacobi, Legendre, Laguerre, Hermite, ... polynomials). We then show that the algebraic parameter estimation problem for a truncated expansion of a function into an orthogonal basis of $L^2$ defined by orthogonal polynomials can be studied similarly. Then, using symbolic computation methods such as Gröbner basis techniques for (noncommutative) polynomial rings, we first show how to compute ordinary differential operators which annihilate a given polynomial and which contain only certain parameters in their coefficients. Then, we explain how to compute the intersection of the annihilator ideals of two polynomials and characterize the ordinary differential operators which annihilate a first polynomial but not a second one. These results are implemented in the NON-A package built upon the OREMODULES software.

## 7.3. Set-Theoretic Methods of Control, Observer Design and Estimation

- Interval Observers:
  - [19] New design of interval observers for continuous-time systems with discrete-time measurements is proposed. For this purpose new conditions of positivity for linear systems with sampled feedbacks are obtained. A sampled-data stabilizing control is synthesized based on provided interval estimates. Efficiency of the obtained solution is demonstrated on examples.
  - [66] The problem of interval state estimation is studied for systems described by parabolic Partial Differential Equations (PDEs). The proposed solution is based on a finite-element approximation of PDE, with posterior design of an interval observer for the obtained ordinary differential equation. The interval inclusion of the state function of PDE is obtained using the estimates on the error of discretization. The results are illustrated by numerical experiments with an academic example.
  - [18] New interval observers are designed for linear systems with time-varying delays in the case of delayed measurements. Interval observers employ positivity and stability analysis of the estimation error system, which in the case of delayed measurements should be delay-dependent. New delay-dependent conditions of positivity for linear systems with time-varying delays are introduced. Efficiency of the obtained solution is demonstrated on examples.
  - [22] Interval state observers provide an estimate on the set of admissible values of the state vector at each instant of time. Ideally, the size of the evaluated set is proportional to the model uncertainty, thus interval observers generate the state estimates with estimation error bounds, similarly to Kalman filters, but in the deterministic framework. Main tools and techniques for design of interval observers are reviewed in this tutorial for continuous-time, discrete-time and time-delayed systems.
  - [43] investigates the problem of observer design for a general class of linear singular time-delay systems, in which the time delays are involved in the state, the output and the known input (if there exists). The involvement of the delay could be multiple which however is rarely studied in the literature. Sufficient conditions are proposed which guarantees the existence of a Luenberger-like observer for the general system.
  - In [90] an interval observer is proposed for on-line estimation of differentiation errors in some class of high-order differentiators (like a high-gain differentiator, or homogeneous nonlinear differentiator, or super-twisting differentiator). The results are verified and validated on the telescopic link of a robotic arm for forestry applications in which the mentioned approaches are used to estimate the extension velocity while the interval observer gives bounds to this estimation.

– The problem of interval observer design is studied in [87] for a class of linear hybrid systems. Several observers are designed oriented on different conditions of positivity and stability for estimation error dynamics. Efficiency of the proposed approach is demonstrated by computer experiments for academic and bouncing ball systems.

– The problem of estimation of sequestered parasites Plasmodium falciparum in malaria, based on measurements of circulating parasites, is addressed in [60]. It is assumed that all (death, transition, recruitment and infection) rates in the model of a patient are uncertain (just intervals of admissible values are given) and the measurements are subject to a bounded noise, then an interval observer is designed. Stability of the observer can be verified by a solution of LMI. The efficiency of the observer is demonstrated in simulation.

- Observer design:

  – [81] presents a new approach for observer design for a class of nonlinear singular systems which can be transformed into a special normal form. The interest of the proposed form is to facilitate the observer synthesis for the studied nonlinear singular systems. Necessary and sufficient geometrical conditions are deduced in order to guarantee the existence of a diffeomorphism which transforms the studied nonlinear singular systems into the proposed normal form.

  – In [38], we investigate the estimation problem for a class of partially observable nonlinear systems. For the proposed Partial Observer Normal Form (PONF), necessary and sufficient conditions are deduced to guarantee the existence of a change of coordinates which can transform the studied system into the proposed PONF. Examples are provided to illustrate the effectiveness of the proposed results.

  – [71] deals with the problem of finite-time and fixed-time observation of linear multiple input multiple output (MIMO) control systems. The nonlinear dynamic observers , which guarantee convergence of the observer states to the original system state in a finite and a fixed (defined a priori) time, are studied. Algorithms for the observers parameters tuning are also developed. The theoretical results are illustrated by numerical examples.

  – [44] Sliding mode control design for linear systems with incomplete and noisy measurements of the output and additive/multiplicative exogenous disturbances is studied. A linear minimax observer estimating the system's state with minimal worst-case error is designed. An algorithm, generating continuous and discontinuous feedbacks, which steers the state as close as possible to a given sliding hyperplane in finite time, is presented. The optimality (sub-optimality) of the designed feedbacks is proven for the case of bounded noises and additive (multiplicative) disturbances of $L_2$-class.

  – [37] deals with the design of a robust control for linear systems with external disturbances using a homogeneous differentiator-based observer based on a implicit Lyapunov function approach. Sufficient conditions for stability of the closed-loop system in the presence of external disturbances are obtained and represented by linear matrix inequalities. The parameter tuning for both controller and observer is formulated as a semi-definite programming problem with linear matrix inequalities constraints. Simulation results illustrate the feasibility of the proposed approach and some improvements with respect to the classic linear observer approach.

  – The problem studied in [17] is one of improving the performance of a class of adaptive observer in the presence of exogenous disturbances. The $H^\infty$ gains of both a conventional and the newly proposed sliding-mode adaptive observer are evaluated, to assess the effect of disturbances on the estimation errors. It is shown that if the disturbance is "matched" in the plant equations, then including an additional sliding-mode feedback injection term, dependent on the plant output, improves the accuracy of observation.

- In [95], we consider the classical reaching problem of sliding mode control design, that is to find a control law which steers the state of a Linear Time-Invariant (LTI) system towards a given hyperplane in a finite time. Since the LTI system is subject to unknown but bounded disturbances we apply the minimax observer which provides the best possible estimate of the system's state. The reaching problem is then solved in observer's state space by constructing a feedback control law. The cases of discontinuous and continuous admissible feedbacks are studied. The theoretical results are illustrated by numerical simulations.

- Estimation and Identification:

  - The problem of output control for linear uncertain system with external perturbations is studied in [77]. It is assumed that the output available for measurements is the higher order derivative of the state only (acceleration for a second order plant), which is also corrupted by noise. Then via series of integration an identification algorithm is proposed for identification of values of all parameters and unknown initial conditions for the state vector. Finally, two control algorithms are developed, adaptive and robust, providing boundedness of trajectories for the system. Efficiency of the obtained solutions is demonstrated by numerical experiments.

  - [24] focuses on the problem of velocity and position estimation. A solution is presented for a class of oscillating systems in which position, velocity and acceleration are zero mean signals. The proposed scheme considers that the dynamic model of the system is unknown. Only noisy acceleration measurements, that may be contaminated by zero mean noise and constant bias, are considered to be available. The proposal uses the periodic nature of the signals obtaining finite-time estimations while tackling integration drift accumulation.

  - In [41], we investigate the problem of simultaneous state and parameter estimation for a class of nonlinear systems which can be transformed into an output depending normal form. A new and simple adaptive observer for such class of systems is presented. Sufficient condition for the existence of the proposed observer is derived. A concrete application is given in order to highlight the effectiveness of the proposed result.

  - In [75], the problem of time-varying parameter identification is studied. To this aim, an identification algorithm is developed in order to identify time-varying parameters in a finite-time. The convergence proofs are based on a notion of finite-time stability over finite intervals of time, i.e. Short-finite-time stability; homogeneity for time-varying systems; and Lyapunov function approach. The algorithm asks for a condition over the regressor term which is related to the classic identifiability condition corresponding to the injectivity of such a term. Simulation results illustrate the feasibility of the proposed algorithm.

## 7.4. Stability, Stabilization, Synchronization

- Input-to-state stability:

  - Supported by a novel field definition and recent control theory results, a new method to avoid local minima is proposed in [25]. It is formally shown that the system has an attracting equilibrium at the target point, repelling equilibriums in the obstacles centers and saddle points on the borders. Those unstable equilibriums are avoided capitalizing on the established Input-to-State Stability (ISS) property of this multistable system. The proposed modification of the PF method is shown to be effective by simulation for a two variables integrator and then applied to an unicycle-like wheeled mobile robots which is subject to additive input disturbances.

  - [62] Motivated by the problem of phase-locking in droop-controlled inverter-based microgrids with delays, the recently developed theory of input-to-state stability (ISS) for multistable systems is extended to the case of multistable systems with delayed dynamics. Sufficient conditions for ISS of delayed systems are presented using Lyapunov-Razumikhin

functions. It is shown that ISS multistable systems are robust with respect to delays in a feedback. The derived theory is applied to two examples. First, the ISS property is established for the model of a nonlinear pendulum and delay-dependent ro-bustness conditions are derived. Second, it is shown that, under certain assumptions, the problem of phase-locking analysis in droop-controlled inverter-based microgrids with delays can be reduced to the stability investigation of the nonlinear pendulum. For this case, corresponding delay-dependent conditions for asymptotic phase-locking are given.

– [103] A necessary and sufficient criterion to establish input-to-state stability (ISS) of nonlinear dynamical systems, the dynamics of which are periodic with respect to certain state variables and which possess multiple invariant solutions (equilibria, limit cycles, etc.), is provided. Unlike standard Lyapunov approaches, the condition is relaxed and formulated via a sign-indefinite function with sign-definite derivative, and by taking the system's periodicity explicitly into account. The new result is established by using the framework of cell structure introduced in [24] and it complements the methods developed in [3], [4] for periodic systems. The efficiency of the proposed approach is illustrated via the global analysis of a nonlinear pendulum with constant persistent input.

– In [53] we revisit the problem of stabilizing a triple integrator using a control that depends on the signs of the state variables. For a more general class of linear systems it is shown that the stabilization by sign feedback is possible, depending on some properties of the system's matrix. The conditions for the stability are established by means of linear matrix inequalities. For the triple integrator, the domain of stability is evaluated. Also, the control law is augmented by a linear feedback and the stability properties for this case, checked. The results are illustrated by numerical experiments for a chain of integrators of third order.

- Stabilization:

  – A solution to the problem of global fixed-time output stabilization of a chain of integrators is proposed in [70]. A nonlinear state feedback and a dynamic observer are designed in order to guarantee both fixed-time estimation and fixed-time control. Robustness with respect to exogenous disturbances and measurement noises is established. The performance of the obtained control and estimation algorithms are illustrated by numeric experiments.

  – In [20], the rate of convergence to the origin for a chain of integrators stabilized by homogeneous feedback is accelerated by a supervisory switching of control parameters. The proposed acceleration algorithm ensures a fixed-time convergence for otherwise exponentially or finite-time stable homogeneous closed-loop systems. Bounded disturbances are taken into account. The results are especially useful in the case of exponentially stable systems widespread in the practice. The proposed switching strategy is illustrated by computer simulation.

  – [33] The problem of robust finite-time stabilization of perturbed multi-input linear system by means of generalized relay feedback is considered. A new control design procedure, which combines convex embedding technique with Implicit Lyapunov Function (ILF) method, is developed. The sufficient conditions for both local and global finite-time stabilization are provided. The issues of practical implementation of the obtained implicit relay feedback are discussed. Our theoretical result is supported by numerical simulation for a Buck converter.

  – [100] contributes to the stability analysis for impulsive dynamical systems based on a vector Lyapunov function and its divergence operator. The new method relies on a 2D time domain representation. The result is illustrated for the exponential stability of linear impulsive systems based on LMIs. The obtained results provide some notions of minimum and maximum dwell-time. Some examples illustrate the feasibility of the proposed approach.

– The Universal Integral Control, introduced in H.K. Khalil, is revisited in [34] by employing mollifiers instead of a high-gain observer for the differentiation of the output signal. The closed loop system is a classical functional differential equation with distributed delays on which standard Lyapunov arguments are applied to study the stability. Low-pass filtering capability of mollifiers is demonstrated for a high amplitude and rapidly oscillating noise. The controller is supported by numerical simulations.

- Synchronization:

  – In [12], we study a robust synchronization problem for multistable systems evolving on manifolds within an Input-to-State Stability (ISS) framework. Based on a recent generalization of the classical ISS theory to multistable systems, a robust synchronization protocol is designed with respect to a compact invariant set of the unperturbed system. The invariant set is assumed to admit a decomposition without cycles, that is, with neither homoclinic nor heteroclinic orbits. Numerical simulation examples illustrate our theoretical results.

  – In [51], [96], motivated by a recent work of R. Brockett Brockett (2013), we study a robust synchronization problem for multistable Brockett oscillators within an Input-to-State Stability (ISS) framework. Based on a recent generalization of the classical ISS theory to multistable systems and its application to the synchronization of multistable systems, a synchronization protocol is designed with respect to compact invariant sets of the unperturbed Brockett oscillator. The invariant sets are assumed to admit a decomposition without cycles (i.e. with neither homoclinic nor heteroclinic orbits). Contrarily to the local analysis of Brockett (2013), the conditions obtained in our work are global and applicable for family of non-identical oscillators. Numerical simulation examples illustrate our theoretical results.

## 7.5. Non-Linear, Sampled-Data And Time-Delay Systems

- Time-delay systems:

  – The problem of delay estimation for a class of nonlinear time-delay systems is considered in [82]. The theory of non-commutative rings is used to analyze the identifiability. Sliding mode technique is utilized in order to estimate the delay showing the possibility to have a local (or global) delay estimation for periodic (or aperiodic) delay signals.

  – In [14] we give sufficient conditions guaranteeing the observability of singular linear systems with commensurable delays affected by unknown inputs appearing in both the state equation and the output equation. These conditions allow for the reconstruction of the entire state vector using past and actual values of the system output. The obtained conditions coincide with known necessary and sufficient conditions of singular linear systems without delays.

  – [67] presents a finite-time observer for linear time-delay systems. In contrast to many observers, which normally estimate the system state in an asymptotic fashion, this observer estimates the exact system state in predetermined finite time. The finite-time observer proposed is achieved by updating the observer state based on actual and pass data of the observer. Simulation results are also presented to illustrate the convergence behavior of the finite-time observer.

  – The backward observability (BO) of a part of the vector of trajectories of the system state is tackled in [57] for a general class of linear time-delay descriptor systems with unknown inputs. By following a recursive algorithm, we present easy testable sufficient conditions ensuring the BO of descriptor time-delay systems.

– Motivated by the problem of phase-locking in droop-controlled inverter-based microgrids with delays, in [23], the recently developed theory of input-to-state stability (ISS) for multistable systems is extended to the case of multistable systems with delayed dynamics. Sufficient conditions for ISS of delayed systems are presented using Lyapunov-Razumikhin functions. It is shown that ISS multistable systems are robust with respect to delays in a feedback. The derived theory is applied to two examples. First, the ISS property is established for the model of a nonlinear pendulum and delay-dependent robustness conditions are derived. Second, it is shown that, under certain assumptions, the problem of phase-locking analysis in droop-controlled inverter-based microgrids with delays can be reduced to the stability investigation of the nonlinear pendulum. For this case, corresponding delay-dependent conditions for asymptotic phase-locking are given.

– Causal and non-causal observability are discussed in [68] for nonlinear time-delay systems. By extending the Lie derivative for time-delay systems in the algebraic framework introduced by Xia et al. (2002), we present a canonical form and give sufficient condition in order to deal with causal and non-causal observations of state and unknown inputs of time-delay systems.

– [83] presents a finite-time observer for linear time-delay systems with commensurate delay. Unlike the existing observers in the literature which converges asymptotically, the proposed observer provides a finite-time estimation. This is realized by using the well-known sliding mode technique. Simulation results are also presented in order to illustrate the feasibility of the proposed method.

• Sampled-Data systems:

– [104] presents basic concepts and recent research directions about the stability of sampled-data systems with aperiodic sampling. We focus mainly on the stability problem for systems with arbitrary time-varying sampling intervals which has been addressed in several areas of research in Control Theory. Systems with aperiodic sampling can be seen as time-delay systems, hybrid systems, Input/Output interconnections, discrete-time systems with time-varying parameters, etc. The goal of the article is to provide a structural overview of the progress made on the stability analysis problem. Without being exhaustive, which would be neither possible nor useful, we try to bring together results from diverse communities and present them in a unified manner. For each of the existing approaches, the basic concepts, fundamental results, converse stability theorems (when available), and relations with the other approaches are discussed in detail. Results concerning extensions of Lyapunov and frequency domain methods for systems with aperiodic sampling are recalled, as they allow to derive constructive stability conditions. Furthermore, numerical criteria are presented while indicating the sources of conservatism, the problems that remain open and the possible directions of improvement. At last, some emerging research directions, such as the design of stabilizing sampling sequences, are briefly discussed.

– In [31] we investigate the stability analysis of nonlinear sampled-data systems, which are affine in the input. We assume that a stabilizing controller is designed using the emulation technique. We intend to provide sufficient stability conditions for the resulting sampled-data system. This allows to find an estimate of the upper bound on the asynchronous sampling intervals, for which stability is ensured. The main idea of the paper is to address the stability problem in a new framework inspired by the dissipativity theory. Furthermore, the result is shown to be constructive. Numerically tractable criteria are derived using linear matrix inequality for polytopic systems and using sum of squares technique for the class of polynomial systems.

– [76] deals with the sampled-data control problem based on state estimation for linear sampled-data systems. An impulsive system approach is proposed based on a vector Lyapunov function method. Observer-based control design conditions are expressed in terms of LMIs. Some examples illustrate the feasibility of the proposed approach.

# 7.6. Effective algebraic systems theory

- Algebraic analysis approach:
  - The purpose of [97] is to present a survey on the effective algebraic analysis approach to linear systems theory with applications to control theory and mathematical physics. In particular, we show how the combination of effective methods of computer algebra − based on Gröbner basis techniques over a class of noncommutative polynomial rings of functional operators called Ore algebras − and constructive aspects of module theory and homological algebra enables the characterization of structural properties of linear functional systems. Algorithms are given and a dedicated implementation, called ORE-ALGEBRAICANALYSIS, based on the Mathematica package HOLONOMICFUNCTIONS, is demonstrated.
  - As far as we know, there is no algebraic (polynomial) approach for the study of linear differential time-delay systems in the case of a (sufficiently regular) time-varying delay. Based on the concept of skew polynomial rings developed by Ore in the 30s, the purpose of [73] is to construct the ring of differential time-delay operators as an Ore extension and to analyze its properties. Classical algebraic properties of this ring, such as noetherianity, its homological and Krull dimensions and the existence of Gröbner bases, are characterized in terms of the time-varying delay function. In conclusion, the algebraic analysis approach to linear systems theory allows us to study linear differential time-varying delay systems (e.g. existence of autonomous elements, controllability, parametrizability, flatness, behavioral approach) through methods coming from module theory, homological algebra and constructive algebra.
  - Within the algebraic analysis approach to linear systems theory, in [98], we investigate the equivalence problem of linear functional systems, i.e., the problem of characterizing when all the solutions of two linear functional systems are in a one-to-one correspondence. To do that, we first provide a new characterization of isomorphic finitely presented modules in terms of inflations of their presentation matrices. We then prove several isomorphisms which are consequences of the unimodular completion problem. We then use these isomorphisms to complete and refine existing results concerning Serre's reduction problem. Finally, different consequences of these results are given. All the results obtained are algorithmic for rings for which Gröbner basis techniques exist and the computations can be performed by the Maple packages OREMODULES and OREMORPHISMS.
  - In [99], we study algorithmic aspects of the algebra of linear ordinary integro-differential operators with polynomial coefficients. Even though this algebra is not Noetherian and has zero divisors, Bavula recently proved that it is coherent, which allows one to develop an algebraic systems theory over this algebra. For an algorithmic approach to linear systems of integro-differential equations with boundary conditions, computing the kernel of matrices with entries in this algebra is a fundamental task. As a first step, we have to find annihilators of integro-differential operators, which, in turn, is related to the computation of polynomial solutions of such operators. For a class of linear operators including integro-differential operators, we present an algorithmic approach for computing polynomial solutions and the index. A generating set for right annihilators can be constructed in terms of such polynomial solutions. For initial value problems, an involution of the algebra of integro-differential operators then allows us to compute left annihilators, which can be interpreted as compatibility conditions of integro-differential equations with boundary conditions.
  - Recent progress in computer algebra has opened new opportunities for the parameter estimation problem in nonlinear control theory, by means of integro-differential input-output equations. In [102], we recall the origin of integro-differential equations. We present new opportunities in nonlinear control theory. Finally, we review related recent theoretical approaches on integro-differential algebras, illustrating what an integro-differential elimination method might be and what benefits the parameter estimation problem would gain from

it.

- Computational real algebraic geometric approach:

  – In [74], we present a symbolic-numeric method for solving the $H_\infty$ loop-shaping design problem for low order single-input single-output systems with parameters. Due to the system parameters, no purely numerical algorithm can indeed solve the problem. Using Gröbner basis techniques and the Rational Univariate Representation of zero-dimensional algebraic varieties, we first give a parametrization of all the solutions of the two Algebraic Riccati Equations associated with the $H_\infty$ control problem. Then, following some works on the spectral factorization problem, a certified symbolic-numeric algorithm is obtained for the computation of the positive definite solutions of these two Algebraic Riccati Equations. Finally, we present a certified symbolic-numeric algorithm which solves the $H_\infty$ loop-shaping design problem for the above class of systems.

  – In [58], the asymptotic stability of linear differential systems with commensurate delays is studied. A classical approach for checking that all the roots of the corresponding quasipolynomial have negative real parts consists in computing the set of critical zeros of the quasipolynomial, i.e., the roots (and the corresponding delays) of the quasipolynomial that lie on the imaginary axis, and then analyzing the variation of these roots with respect to the variation of the delay. Based on solving algebraic systems techniques, a certified and efficient symbolic-numeric algorithm for computing the set of critical roots of a quasipolynomial is proposed. Moreover, using recent algorithmic results developed by the computer algebra community, we present an efficient algorithm for the computation of Puiseux series at a critical zero which allows us to finely analyze the stability of the system with respect to the variation of the delay.

  – In [59], we present new computer algebra based methods for testing the structural stability of $n$-D discrete linear systems (with $n \geq 2$). More precisely, we show that the standard characterization of the structural stability of a multivariate rational transfer function (namely, the denominator of the transfer function does not have solutions in the unit polydisc of $\mathbb{C}^n$) is equivalent to fact that a certain system of polynomials does not have real solutions. We then use state-of-the-art algorithms of the computer algebra community to check this last condition, and thus the structural stability of multidimensional systems.

## 7.7. Applications

- A fault detection method for an automatic detection of spawning in oysters [13]:

  Using measurements of valve activity (i.e. the distance between the two valves) in populations of bivalves under natural environmental condition (16 oysters in the Bay of Arcachon, France, in 2007, 2013 and 2014), an algorithm for an automatic detection of the spawning period of oysters is proposed in this paper. Spawning observations are important in aquaculture and biological studies, and until now, such a detection is done through visual analysis by an expert. The algorithm is based on the fault detection approach and it works through the estimation of velocity of valve movement activity, that can be obtained by calculating the time derivative of the valve distance. A summarized description of the methods used for the derivative estimation is provided, followed by the associated signal processing and decision making algorithm to determine spawning from the velocity signal. A protection from false spawning detection is also considered by analyzing the simultaneity in spawning. Through this study, it is shown that spawning in a population of oysters living in their natural habitat (i.e. in the sea) can be automatically detected without any human expertise saving time and resources. The fault detection method presented in the paper can also be used to detect complex oscillatory behavior which is of interest to control engineering community.

- Robust synchronization of genetic oscillators [52]:

Cell division introduces discontinuities in the dynamics of genetic oscillators (circadian clocks, synthetic oscillators, etc.) causing phase drift. This paper considers the problem of phase synchronization for a population of genetic oscillators that undergoes cell division and with a common entraining input in the population. Inspired by stochastic simulation, this paper proposes analytical conditions that guarantee phase synchronization. This analytical conditions are derived based on Phase Response Curve (PRC) model of an oscillator (the first order reduced model obtained for the linearized system and inputs with sufficiently small amplitude). Cell division introduces state resetting in the model (or phase resetting in the case of phase model), placing it in the class of hybrid systems. It is shown through numerical experiments for a motivating example that without common entraining input in all oscillators, the cell division acts as a disturbance causing phase drift, while the presence of entrainment guarantees boundedness of synchronization phase errors in the population. Theoretical developments proposed in the paper are demonstrated through numerical simulations for two different genetic oscillator models (Goodwin oscillator and Van der Pol oscillator).

- Modeling pointing tasks in mouse-based human-computer interactions [54]:

  Pointing is a basic gesture performed by any user during human-computer interaction. It consists in covering a distance to select a target via the cursor in a graphical user interface (e.g. a computer mouse movement to select a menu element). In this work, a dynamic model is proposed to describe the cursor motion during the pointing task. The model design is based on experimental data for pointing with a mouse. The obtained model has switched dynamics, which corresponds well to the state of the art accepted in the human-computer interaction community. The conditions of the model stability are established. The presented model can be further used for the improvement of user performance during pointing tasks.

- Modeling and control of turbulent flows [64]:

  The model-based closed-loop control of a separated flow can be studied based on the model described by Navier-Stokes equation. However, such a model still rises difficult issues for control practice. An alternative bilinear and delayed model has been developed tested on the experiments allowing its identification. The identification technique combines least-square technique with a Mesh Adaptive Direct Search (MADS) algorithm.

- Practical design considerations for successful industrial application of model-based fault detection techniques to aircraft systems [47]:

  This paper discusses some key factors which may arise for successful application of model-based Fault Detection(FD) techniques to aircraft systems. The paper reports on the results and the lessons learned during flight V & V(Validation & Verification) activities, implementation in the A380 Flight Control Computer(FCC) and A380 flight tests at Airbus(Toulouse, France).The paper does not focus on new theoretical materials, but rather on a number of practical design considerations to provide viable technological solutions and mechanization schemes. The selected case studies are taken from past and on-going research actions between Airbus and the University of Bordeaux (France). One of the presented solutions has received final certification on new generation Airbus A350 aircraft and is flying (first commercial flight: January 15,2015)

- Finite-time obstacle avoidance for unicycle-like robot [26]:

  The problem of avoiding obstacles while navigating within an environment for a Unicycle-like Wheeled Mobile Robot (WMR) is of prime importance in robotics; the aim of this work is to solve such a problem proposing a perturbed version of the standard kinematic model able to compensate for the neglected dynamics of the robot. The disturbances are considered additive on the inputs and the solution is based on the supervisory control framework, finite-time stability and a robust multi-output regulation. The effectiveness of the solution is proved, supported by experiments and finally compared with the Dynamic Window Approach (DWA) to show how the proposed method can perform better than standard methods.

- Almost global attractivity of a synchronous generator connected to an infinite bus [56]:

The problem of deriving verifiable conditions for stability of the equilibria of a realistic model of a synchronous generator with constant field current connected to an infinite bus is studied in the paper. Necessary and sufficient conditions for existence and uniqueness of equilibrium points are provided. Furthermore, sufficient conditions for almost global attractivity are given. To carry out this analysis a new Lyapunov–like function is proposed to establish convergence of bounded trajectories, while the latter is proven using the powerful theoretical framework of cell structures pioneered by Leonov and Noldus.

<p style="text-align:center; color:red;">**RAPSODI Team**</p>

# 7. New Results

## 7.1. large-time behavior of some numerical schemes

In [19], 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.

In [13], M. Bessemoulin-Chatard and C. Chainais-Hillairet study the large-time behavior of a numerical scheme discretizing drift-diffusion systems for semiconductors. The numerical method is based on a  generalization of the classical Scharfetter-Gummel scheme which allows to consider both linear or nonlinear pressure laws.They study the convergence of approximate solutions towards an approximation of the thermal equilibrium state as time tends to infinity, and obtain a decay rate by controlling the discrete relative entropy with the entropy production. This result is proved under assumptions of existence and uniform-in-time $L^\infty$ estimates for numerical solutions, which are then discussed.

The question of uniform-in-time $L^\infty$ estimates for the scheme proposed in [13] has then be tackled by M. Bessemoulin-Chatard, C. Chainais-Hillairet and A. Jüngel. The result is obtained *via* a Moser's iteration technique adapted to the discrete setting. Related to this question, the existence of a positive lower bound for the numerical solution of a convection-diffusion equation has been studied by C. Chainais-Hillairet, B. Merlet and A. Vasseur. They apply a method due to De Giorgi in order to establish a positive lower bound for the numerical solution of a stationary convection-diffusion equation. These results are submitted for publication in the FVCA8 conference (to be held in June 2017).

In [11] B. Merlet *et al.* consider a second-order two-step time discretization of the Cahn-Hilliard equation with an analytic nonlinearity. They study the long time behavior of the discrete solution and show that if the time-step is chosen small enough, the sequence generated by the scheme converges to a steady state as time tends to infinity. Convergence rates are also provided. This parallels the behavior of the solutions of the non-discretized solutions and shows the reliability of the scheme for long time simulations. The method of proof is based on the Lojasiewicz-Simon inequality and on the study of the pseudo-energy associated with the discretization which is shown to be non-increasing.

## 7.2. Theoretical and numerical analysis of corrosion models

The Diffusion Poisson Coupled Model [1] 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). The DPCM model also takes into account the growth of the oxide host lattice and its dissolution, leading to moving boundary equations. Numerical experiments done for the simulation of this model with moving boundaries show the convergence in time towards a pseudo-steady-state. C. Chainais-Hillairet and T. O. Gallouët  prove in [18]  the existence of pseudo-stationary solutions for some simplified versions of the DPCM model. They also propose a new scheme in order to compute directly this pseudo-steady-state. Numerical experiments show the efficiency of this method.

The modeling of concrete carbonation also leads to a system of partial differential equations posed on a moving domain. C. Chainais-Hillairet, B. Merlet and A. Zurek propose and analyze a finite volume scheme for the concrete carbonation model. They prove the convergence of the sequence of approximate solutions towards a  weak solution. Numerical experiments show the order 2 in space of the scheme and illustrate the $\sqrt{t}$ law of propagation of  the size of the carbonated zone. This result is submitted for publication.

## 7.3. Modeling and numerics for porous media flows

In [16], C. Cancès and C. Guichard propose a nonlinear Control Volume Finite Elements method with upwinding in order to solve possibly nonlinear and degenerate parabolic equations. This method was designed in order to preserve at the discrete level the positivity and the nonlinear stability of the solutions. In [25], A. Ait Hammou Oulhaj, C. Cancès, and C. Chainais-Hillairet extend the approach of [16] to the more complex case of Richards equation modeling saturated/unsaturated flows in anisotropic porous media. The additional complexity comes from the fact that convective terms and elliptic degeneracy are considered in [25]. The scheme preserves at the discrete level the nonnegativity and the nonlinear stability of the solutions. Its convergence is rigorously proved, and numerical results are provided in order to illustrate the behavior of the scheme.

In [49], C. Cancès, T. O. Gallouët, and L. Monsaingeon show that the equations governing two-phase flows in porous media have a formal gradient flow structure. The goal of the longer contribution [29] is then twofold. First, it extends the variational interpretation of [49] to the case where an arbitrary number of phases are in competition to flow within a porous medium. Second, we provide rigorous foundations to our claim. More precisely, the convergence of a minimizing movement scheme *à la* Jordan, Kinderlehrer, and Otto [66] is shown in [29], providing by the way a new existence result for multiphase flows in porous media. The result relies on advances tools related to optimal transportation [75], [74].

## 7.4. Complex fluid flows: modeling, analysis and numerics

The analysis of the Kazhikhov-Smagulov model was given by Bresch at al. [48] (see also reference therein). These authors prove the global existence of weak solution without assuming small data and without any assumption on the diffusivity coefficient. Following the physical experiment given by Joseph [67], we introduce a Korteweg stress tensor in the previous model. The theory of Korteweg considers the possibility that motions can be driven by additional stresses associated with gradients of density. In process of slow diffusion on miscible incompressible fluids, for example water and glycerin, dynamical effects which mimic surface tension can arise in thin mixing layers where the gradients of density are large. In the context of the PhD thesis of Meriem Ezzoug (July 2016, University of Monastir, Tunisia), C. Calgaro and co-authors study a multiphase incompressible fluid model, called the Kazhikhov-Smagulov-Korteweg model. They prove in [14] that this model is globally well posed in a 3D bounded domain.

In [21], P.-E. Jabin and T. Rey investigate the behavior of granular gases in the limit of small Knudsen number, that is very frequent collisions. They deal with the physically relevant strongly inelastic case, in one dimension of space and velocity. The study of such limit, also known as hydrodynamic limit is to give a reduced description of the kinetic equation, using a fluid approximation. They are able to prove the convergence of the particle distribution function toward a monokinetic distribution, whose moments verify the pressureless Euler system. The proof relies on dispersive relations at the kinetic level, which leads to the so-called Oleinik property at the limit, and in particular stability of the solution to the fluid problem.

In [34], I. Lacroix-Violet and A. Vasseur present the construction of global weak solutions to the quantum Navier-Stokes equation, for any initial value with bounded energy and entropy. The construction is uniform with respect to the Planck constant. This allows to perform the semi-classical limit to the associated compressible Navier-Stokes equation. One of the difficulty of the problem is to deal with the degenerate viscosity, together with the lack of integrability on the velocity. The method is based on the construction of weak solutions that are renormalized in the velocity variable. The existence, and stability of these solutions do not need the Mellet-Vasseur inequality [71].

In [31], G. Dimarco, R. Loubère, J. Narski and T. Rey deal with the extension of the Fast Kinetic Scheme (FKS) [55], [56] originally constructed for solving the BGK equation, to the more challenging case of the Boltzmann equation. The scheme combines a robust and fast method for treating the transport part based on an innovative Lagrangian technique supplemented with fast spectral schemes to treat the collisional operator by means of an operator splitting approach. This approach along with several implementation features related to the parallelization of the algorithm permits to construct an efficient simulation tool which is numerically

tested against exact and reference solutions on classical problems arising in rarefied gas dynamic. They present results up to the 3D×3D case for unsteady flows for the Variable Hard Sphere model which may serve as benchmark. For this reason, they also provide for each problem details on the computational cost and memory consumption as well as comparisons with the BGK model or the limit model of compressible Euler equations.

## 7.5. Improving the numerical efficiency of numerical methods

In this section, we gather contributions in which a methodology was introduced in order to reduce the computational cost at fixed accuracy or to improve the accuracy for a fixed computational cost.

In [20], E. Creusé and his collaborators generalized some of their previous results on residual a posteriori error estimators for low electromagnetism [10] , [52] to the case where some voltage or current excitation is specified in the model (see e.g. such models in [63], [42]). It consequently led to consider different formulations and to overcome some specific difficulties in order to derive the reliability of the involved estimators.

It is now well accepted that well-balanced schemes are of great interest in order to compute accurate solutions to systems of PDEs (see for instance  [60]). In [36], L. Pareschi and T. Rey propose a systematic way to tune classical numerical schemes in order to make them well-balanced and asymptotic preserving. Inspired by micro-macro decomposition methods for kinetic equations, they present a class of schemes which are capable to preserve the steady state solution and achieve high order accuracy for a class of time dependent partial differential equations including nonlinear diffusion equations and kinetic equations. Extension to systems of conservation laws with source terms are also discussed, as well as Total Variation Diminishing preserving properties.

The contribution [26] by K. Brenner and C. Cancès is devoted to the improvement of the behavior of Newton's method when solving degenerate parabolic equations. Such equations are very common for instance in the context of complex porous media flows. In [26], the presentation focuses on Richards equation modeling saturated/unsaturated flows in porous media. The basic idea is the following: Newton's method is not invariant by nonlinear change of variables. The choice of the primary variable then impacts the effective resolution of the nonlinear system provided by the scheme. The idea developed in [26] is then to construct an abstract primary variable to facilitate Newton's method's convergence. This leads to an impressive reduction of the computational cost, a better accuracy in the results and an strong robustness of the method w.r.t. the nonlinearities appearing in the continuous model.

## 7.6. Variational modeling and analysis

Bose-Einstein condensates are a unique way to observe quantum effects at a (relatively) large scale. The fundamental states of such condensates are obtained as minimizers of a Gross-Pitaievskii functional. In [33], M. Goldman and B. Merlet consider the case of a two component Bose-Einstein condensate in the strong segregation regime (the energy favors spatial segregation of the two different Boson species). They identify two different regimes in the strong segregation and small healing length limit. In one of these regimes, the relevant limit is an interesting weighted isoperimetric problem which explains some of the numerical simulations of [70].

In [32], B. Merlet *et al.* consider the branched transportation problem in 2D associated with a cost per unit length of the form $1 + \alpha m$ where $m$ denotes the amount of transported mass and $\alpha > 0$ is a fixed parameter (the limit case $\alpha = 0$ corresponds to the classical Steiner problem). Motivated by the numerical approximation of this problem, they introduce a family of functionals $(\{F_\varepsilon\}_{\varepsilon>0})$ which approximate the above branched transport energy. They justify rigorously the approximation by establishing the equicoercivity and the $\Gamma$-convergence of $F_\varepsilon$ as $\varepsilon \downarrow 0$. These functionals are modeled on the Ambrosio-Tortorelli functional and are easy to optimize in practice (the algorithm amounts to perform repetitively the alternate optimization of two quadratic functionals). Numerical evidences of the efficiency of the method are presented.

## 7.7. Miscellaneous

This section gathers results from members of the team that are not directly related to the core of the scientific program of the team.

In [12], I. Violet-Lacroix and co-authors consider the derivation of continuous and fully discrete artificial boundary conditions for the linearized Korteweg-de-Vries equation. They are provided for two different numerical schemes. The boundary conditions being nonlocal with respect to time variable, they propose fast evaluations of discrete convolutions. Various numerical tests are presented to show the effectiveness of the constructed artificial boundary conditions.

A semi-discrete in time Crank-Nicolson scheme to discretize a weakly damped forced nonlinear fractional Schrödinger equation in the whole space ($\mathbb{R}$ is considered by C. Calgaro and co-authors in [28]. They prove that such semi-discrete equation provides a discrete infinite dimensional dynamical in $H^\alpha(\mathbb{R})$ that possesses a global attractor. They show also that if the external force is in a suitable weighted Lebesgue space then this global attractor has a finite fractal dimension.

In [35], F. Nabet considers a finite-volume approximation, based on a two point flux approximation, for the Cahn-Hilliard equation with dynamic boundary conditions. An error estimate for the fully-discrete scheme on a possibly smooth non-polygonal domain is proved and numerical simulations which validate the theoretical result are given.

<p style="text-align: center; color: red;">**RMOD Project-Team**</p>

# 7. New Results

## 7.1. Practical Validation of Bytecode to Bytecode JIT Compiler Dynamic Deoptimization.

Speculative inlining in just-in-time compilers enables many performance optimizations. However, it also introduces significant complexity. The compiler optimizations themselves, as well as the deoptimization mechanism are complex and error prone. To stabilize our bytecode to bytecode just-in-time compiler, we designed a new approach to validate the correctness of dynamic deoptimization. The approach consists of the symbolic execution of an optimized and an unop-timized bytecode compiled method side by side, deoptimizing the abstract stack at each deoptimization point (where dynamic deoptimization is possible) and comparing the deoptimized and unoptimized abstract stack to detect bugs. The implementation of our approach generated tests for several hundred thousands of methods, which are now available to be run automatically after each commit [13].

## 7.2. Recording and Replaying System-Specific Conventions.

In other situations, we found that developers sometimes perform sequences of code changes in a systematic way. These sequences consist of small code changes (*e.g.*, create a class, then extract a method to this class), which are applied to groups of related code entities (*e.g.*, some of the methods of a class). We propose to help this task by letting the developer record the sequence of code changes when he first applies it, and then generalize this sequence to apply it in other locations. The evaluation is based on real instances of such sequences that we found in different open source systems. We were able to replay 92% of the examples, which consisted in up to seven code entities modified up to 66 times. We are still working on the approach to allow for (semi-)automatic generalization of the recorded sequence of changes [71], [70].

## 7.3. Test Case Selection in Industry: an Analysis of Issues Related to Static Approaches

Automatic testing constitutes an important part of everyday development practice. But running all these tests may take hours. This is especially true for large systems involving, for example, the deployment of a web server or communication with a database. For this reason, tests are not launched as often as they should be and are mostly run at night. The company wishes to improve its development and testing process by giving to developers rapid feedback after a change. An interesting solution to give developers rapid feedback after a change is to reduce the number of tests to run by identifying only those exercising the piece of code changed. Two main approaches are proposed in the literature: static and dynamic. We evaluate these approaches on three industrial, closed source, cases to understand the strengths and weaknesses of each solution. We also propose a classification of problems that may arise when trying to identify the tests that cover a method.

# SEQUEL Project-Team

# 7. New Results

## 7.1. Decision-making Under Uncertainty

### 7.1.1. Reinforcement Learning

**Analysis of Classification-based Policy Iteration Algorithms**, [20]

We introduce a variant of the classification-based approach to policy iteration which uses a cost-sensitive loss function weighting each classification mistake by its actual regret, that is, the difference between the action-value of the greedy action and of the action chosen by the classifier. For this algorithm, we provide a full finite-sample analysis. Our results state a performance bound in terms of the number of policy improvement steps, the number of rollouts used in each iteration, the capacity of the considered policy space (classifier), and a capacity measure which indicates how well the policy space can approximate policies that are greedy with respect to any of its members. The analysis reveals a tradeoff between the estimation and approximation errors in this classification-based policy iteration setting. Furthermore it confirms the intuition that classification-based policy iteration algorithms could be favorably compared to value-based approaches when the policies can be approximated more easily than their corresponding value functions. We also study the consistency of the algorithm when there exists a sequence of policy spaces with increasing capacity.

**Reinforcement Learning of POMDPs using Spectral Methods**, [23]

We propose a new reinforcement learning algorithm for partially observable Markov decision processes (POMDP) based on spectral decomposition methods. While spectral methods have been previously employed for consistent learning of (passive) latent variable models such as hidden Markov models, POMDPs are more challenging since the learner interacts with the environment and possibly changes the future observations in the process. We devise a learning algorithm running through episodes, in each episode we employ spectral techniques to learn the POMDP parameters from a trajectory generated by a fixed policy. At the end of the episode, an optimization oracle returns the optimal memoryless planning policy which maximizes the expected reward based on the estimated POMDP model. We prove an order-optimal regret bound w.r.t. the optimal memoryless policy and efficient scaling with respect to the dimensionality of observation and action spaces.

**Bayesian Policy Gradient and Actor-Critic Algorithms**, [15]

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 dynamics within trajectories to be of any particular 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 proceed to supplement our Bayesian policy gradient framework with 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.

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

**Improved Learning Complexity in Combinatorial Pure Exploration Bandits**, [32]

We study the problem of combinatorial pure exploration in the stochastic multi-armed bandit problem. We first construct a new measure of complexity that provably characterizes the learning performance of the algorithms we propose for the fixed confidence and the fixed budget setting. We show that this complexity is never higher than the one in existing work and illustrate a number of configurations in which it can be significantly smaller. While in general this improvement comes at the cost of increased computational complexity, we provide a series of examples , including a planning problem, where this extra cost is not significant.

**Online learning with noisy side observations**, [43]

We propose a new partial-observability model for online learning problems where the learner, besides its own loss, also observes some noisy feedback about the other actions, depending on the underlying structure of the problem. We represent this structure by a weighted directed graph, where the edge weights are related to the quality of the feedback shared by the connected nodes. Our main contribution is an efficient algorithm that guarantees a regret of $O(\sqrt{\alpha * T})$ after T rounds, where $\alpha$ * is a novel graph property that we call the effective independence number. Our algorithm is completely parameter-free and does not require knowledge (or even estimation) of $\alpha$ *. For the special case of binary edge weights, our setting reduces to the partial-observability models of Mannor & Shamir (2011) and Alon et al. (2013) and our algorithm recovers the near-optimal regret bounds.

**Online learning with Erdös-Rényi side-observation graphs**, [42]

We consider adversarial multi-armed bandit problems where the learner is allowed to observe losses of a number of arms beside the arm that it actually chose. We study the case where all non-chosen arms reveal their loss with an unknown probability rt, independently of each other and the action of the learner. Moreover, we allow rt to change in every round t, which rules out the possibility of estimating rt by a well-concentrated sample average. We propose an algorithm which operates under the assumption that rt is large enough to warrant at least one side observation with high probability. We show that after T rounds in a bandit problem with N arms, the expected regret of our algorithm is of order O(sqrt(sum(t=1)T (1/rt) log N )), given that rt less than log T / (2N-2) for all t. All our bounds are within logarithmic factors of the best achievable performance of any algorithm that is even allowed to know exact values of rt.

**Revealing graph bandits for maximizing local influence**, [27]

We study a graph bandit setting where the objective of the learner is to detect the most influential node of a graph by requesting as little information from the graph as possible. One of the relevant applications for this setting is marketing in social networks, where the marketer aims at finding and taking advantage of the most influential customers. The existing approaches for bandit problems on graphs require either partial or complete knowledge of the graph. In this paper, we do not assume any knowledge of the graph, but we consider a setting where it can be gradually discovered in a sequential and active way. At each round, the learner chooses a node of the graph and the only information it receives is a stochastic set of the nodes that the chosen node is currently influencing. To address this setting, we propose BARE, a bandit strategy for which we prove a regret guarantee that scales with the detectable dimension, a problem dependent quantity that is often much smaller than the number of nodes.

**Algorithms for Differentially Private Multi-Armed Bandits**, [50]

We present differentially private algorithms for the stochastic Multi-Armed Bandit (MAB) problem. This is a problem for applications such as adaptive clinical trials, experiment design, and user-targeted advertising where private information is connected to individual rewards. Our major contribution is to show that there exist $(\epsilon, \delta)$ differentially private variants of Upper Confidence Bound algorithms which have optimal regret, $O(\epsilon^{-1} + \log T)$. This is a significant improvement over previous results, which only achieve poly-log regret

$O(\epsilon^{-2} \log^2 T)$, because of our use of a novel interval-based mechanism. We also substantially improve the bounds of previous family of algorithms which use a continual release mechanism. Experiments clearly validate our theoretical bounds.

**On the Complexity of Best Arm Identification in Multi-Armed Bandit Models**, [17]

The stochastic multi-armed bandit model is a simple abstraction that has proven useful in many different contexts in statistics and machine learning. Whereas the achievable limit in terms of regret minimization is now well known, our aim is to contribute to a better understanding of the performance in terms of identifying the m best arms. We introduce generic notions of complexity for the two dominant frameworks considered in the literature: fixed-budget and fixed-confidence settings. In the fixed-confidence setting, we provide the first known distribution-dependent lower bound on the complexity that involves information-theoretic quantities and holds when m is larger than 1 under general assumptions. In the specific case of two armed-bandits, we derive refined lower bounds in both the fixed-confidence and fixed-budget settings, along with matching algorithms for Gaussian and Bernoulli bandit models. These results show in particular that the complexity of the fixed-budget setting may be smaller than the complexity of the fixed-confidence setting, contradicting the familiar behavior observed when testing fully specified alternatives. In addition, we also provide improved sequential stopping rules that have guaranteed error probabilities and shorter average running times. The proofs rely on two technical results that are of independent interest : a deviation lemma for self-normalized sums (Lemma 19) and a novel change of measure inequality for bandit models (Lemma 1).

**Optimal Best Arm Identification with Fixed Confidence**, [33]

We give a complete characterization of the complexity of best-arm identification in one-parameter bandit problems. We prove a new, tight lower bound on the sample complexity. We propose the 'Track-and-Stop' strategy, which we prove to be asymptotically optimal. It consists in a new sampling rule (which tracks the optimal proportions of arm draws highlighted by the lower bound) and in a stopping rule named after Chernoff, for which we give a new analysis.

**On Explore-Then-Commit Strategies**, [35]

We study the problem of minimising regret in two-armed bandit problems with Gaussian rewards. Our objective is to use this simple setting to illustrate that strategies based on an exploration phase (up to a stopping time) followed by exploitation are necessarily suboptimal. The results hold regardless of whether or not the difference in means between the two arms is known. Besides the main message, we also refine existing deviation inequalities, which allow us to design fully sequential strategies with finite-time regret guarantees that are (a) asymptotically optimal as the horizon grows and (b) order-optimal in the minimax sense. Furthermore we provide empirical evidence that the theory also holds in practice and discuss extensions to non-gaussian and multiple-armed case.

### 7.1.3. Recommendation systems

**Scalable explore-exploit Collaborative Filtering**, [39]

Recommender Systems (RS) aim at suggesting to users one or several items in which they might have interest. These systems have to update themselves as users provide new ratings, but also as new users/items enter the system. While this adaptation makes recommendation an intrinsically sequential task, most researches about RS based on Collaborative Filtering are omitting this fact, as well as the ensuing exploration/exploitation dilemma: should the system recommend items which bring more information about the users (explore), or should it try to get an immediate feedback as high as possible (exploit)? Recently, a few approaches were proposed to solve that dilemma, but they do not meet requirements to scale up to real life applications which is a crucial point as the number of items available on RS and the number of users in these systems explode. In this paper, we present an explore-exploit Collaborative Filtering RS which is both efficient and scales well. Extensive experiments on some of the largest available real-world datasets show that the proposed approach performs accurate personalized recommendations in less than a millisecond per recommendation, which makes it a good candidate for true applications.

**Large-scale Bandit Recommender System**, [38]

The main target of Recommender Systems (RS) is to propose to users one or several items in which they might be interested. However, as users provide more feedback, the recommendation process has to take these new data into consideration. The necessity of this update phase makes recommendation an intrinsically sequential task. A few approaches were recently proposed to address this issue, but they do not meet the need to scale up to real life applications. In this paper , we present a Collaborative Filtering RS method based on Matrix Factorization and Multi-Armed Bandits. This approach aims at good recommendations with a narrow computation time. Several experiments on large datasets show that the proposed approach performs personalized recommendations in less than a millisecond per recommendation.

**Sequential Collaborative Ranking Using (No-)Click Implicit Feedback**, [40]

We study Recommender Systems in the context where they suggest a list of items to users. Several crucial issues are raised in such a setting: first, identify the relevant items to recommend; second, account for the feedback given by the user after he clicked and rated an item; third, since new feedback arrive into the system at any moment, incorporate such information to improve future recommendations. In this paper, we take these three aspects into consideration and present an approach handling click/no-click feedback information. Experiments on real-world datasets show that our approach outperforms state of the art algorithms.

**Hybrid Recommender System based on Autoencoders**, [49]

A standard model for Recommender Systems is the Matrix Completion setting: given partially known matrix of ratings given by users (rows) to items (columns), infer the unknown ratings. In the last decades, few attempts where done to handle that objective with Neural Networks, but recently an architecture based on Autoencoders proved to be a promising approach. In current paper, we enhanced that architecture (i) by using a loss function adapted to input data with missing values, and (ii) by incorporating side information. The experiments demonstrate that while side information only slightly improve the test error averaged on all users/items, it has more impact on cold users/items.

**Compromis exploration-exploitation pour système de recommandation à grande échelle**, [53]

Les systèmes de recommandation recommandent à des utilisateurs un ou des produits qui pourraient les intéresser. La recommandation se fonde sur les retours des utilisateurs par le passé, lors des précédentes recommandations. La recommandation est donc un problème séquentiel et le système de recommandation recommande (i) pour obtenir une bonne récompense, mais aussi (ii) pour mieux cerné l'utilisateur/les produits et ainsi obtenir de meilleures récompenses par la suite. Quelques approches récentes ciblent ce double objectif mais elles sont trop gourmandes en temps de calcul pour s'appliquer à certaines applications de la vie réelle. Dans cet article, nous présentons un système de recommandation fondé sur la factorisation de matrice et les bandits manchots. Plusieurs expériences sur de grandes base de données montrent que l'approche proposée fournit de bonnes recommandations en moins d'une milli-seconde par recommandation.

**Filtrage Collaboratif Hybride avec des Auto-encodeurs**, [54]

Le filtrage collaboratif (CF) exploite les retours des utilisateurs pour leur fournir des recommandations personnalisées. Lorsque ces algorithmes ont accès à des informations complémentaires, ils ont de meilleurs résultats et gèrent plus efficacement le démarrage à froid. Bien que les réseaux de neurones (NN) remportent de nombreux succès en traitement d'images, ils ont reçu beaucoup moins d'attention dans la communauté du CF. C'est d'autant plus surprenant que les NN apprennent comme les algorithme de CF une représentation latente des données. Dans cet article, nous introduisons une architecture de NN adaptée au CF (nommée CFN) qui prend en compte la parcimonie des données et les informations complémentaires. Nous montrons empiriquement sur les bases de données MovieLens et Douban que CFN bât l'état de l'art et profite des informations complémentaires. Nous fournissons une implémentation de l'algorithme sous forme d'un plugin pour Torch.

## 7.1.4. *Nonparametric statistics of time series*

**Things Bayes can't do**, [48]

The problem of forecasting conditional probabilities of the next event given the past is consideredin a general probabilistic setting. Given an arbitrary (large, uncountable) set C of predictors, we would like to construct a single predictor that performs asymptotically as well as the best predictor in C, on any data. Here we show that there are sets C for which such predictors exist, but none of them is a Bayesian predictor with a prior concentrated on C.In other words, there is a predictor with sublinear regret, but every Bayesian predictor must have a linear regret. This negative finding is in sharp contrast with previous resultsthat establish the opposite for the case when one of the predictors in C achieves asymptotically vanishing error.In such a case, if there is a predictor that achieves asymptotically vanishing error for any measure in C, then there is a Bayesian predictor that also has this property, and whose prior is concentrated on (a countable subset of) C.

### 7.1.5. *Imitation and Inverse Reinforcement Learning*

**Score-based Inverse Reinforcement Learning**, [29]

This paper reports theoretical and empirical results obtained for the score-based Inverse Reinforcement Learning (IRL) algorithm. It relies on a non-standard setting for IRL consisting of learning a reward from a set of globally scored trajec-tories. This allows using any type of policy (optimal or not) to generate trajectories without prior knowledge during data collection. This way, any existing database (like logs of systems in use) can be scored a posteriori by an expert and used to learn a reward function. Thanks to this reward function, it is shown that a near-optimal policy can be computed. Being related to least-square regression, the algorithm (called SBIRL) comes with theoretical guarantees that are proven in this paper. SBIRL is compared to standard IRL algorithms on synthetic data showing that annotations do help under conditions on the quality of the trajectories. It is also shown to be suitable for real-world applications such as the optimisation of a spoken dialogue system.

### 7.1.6. *Stochastic Games*

**Blazing the trails before beating the path: Sample-efficient Monte-Carlo planning**, [37]

You are a robot and you live in a Markov decision process (MDP) with a finite or an infinite number of transitions from state-action to next states. You got brains and so you plan before you act. Luckily, your roboparents equipped you with a generative model to do some Monte-Carlo planning. The world is waiting for you and you have no time to waste. You want your planning to be efficient. Sample-efficient. Indeed, you want to exploit the possible structure of the MDP by exploring only a subset of states reachable by following near-optimal policies. You want guarantees on sample complexity that depend on a measure of the quantity of near-optimal states. You want something, that is an extension of Monte-Carlo sampling (for estimating an expectation) to problems that alternate maximization (over actions) and expectation (over next states). But you do not want to StOP with exponential running time, you want something simple to implement and computationally efficient. You want it all and you want it now. You want TrailBlazer.

**Maximin Action Identification: A New Bandit Framework for Games**, [34]

We study an original problem of pure exploration in a strategic bandit model motivated by Monte Carlo Tree Search. It consists in identifying the best action in a game, when the player may sample random outcomes of sequentially chosen pairs of actions. We propose two strategies for the fixed-confidence setting: Maximin-LUCB, based on lower-and upper-confidence bounds; and Maximin-Racing, which operates by successively eliminating the sub-optimal actions. We discuss the sample complexity of both methods and compare their performance empirically. We sketch a lower bound analysis, and possible connections to an optimal algorithm.

## 7.2. Statistical analysis of time series

### 7.2.1. *Change Point Analysis*

**Nonparametric multiple change point estimation in highly dependent time series**, [18]

Given a heterogeneous time-series sample, the objective is to find points in time, called change points, where the probability distribution generating the data has changed. The data are assumed to have been generated by arbitrary unknown stationary ergodic distributions. No modelling, independence or mixing assumptions 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.

### 7.2.2. *Clustering Time Series, Online and Offline*

**Consistent Algorithms for Clustering Time Series**, [19]

The problem of clustering is considered for the case where every point is a time series. The time series are either given in one batch (offline setting), or they are allowed to grow with time and new time series can be added along the way (online setting). We propose a natural notion of consistency for this problem, and show that there are simple, com-putationally efficient algorithms that are asymptotically consistent under extremely weak assumptions on the distributions that generate the data. The notion of consistency is as follows. A clustering algorithm is called consistent if it places two time series into the same cluster if and only if the distribution that generates them is the same. In the considered framework the time series are allowed to be highly dependent, and the dependence can have arbitrary form. If the number of clusters is known, the only assumption we make is that the (marginal) distribution of each time series is stationary ergodic. No paramet-ric, memory or mixing assumptions are made. When the number of clusters is unknown, stronger assumptions are provably necessary, but it is still possible to devise nonparametric algorithms that are consistent under very general conditions. The theoretical findings of this work are illustrated with experiments on both synthetic and real data.

### 7.2.3. *Automata Learning*

**PAC learning of Probabilistic Automaton based on the Method of Moments**, [36]

Probabilitic Finite Automata (PFA) are gener-ative graphical models that define distributions with latent variables over finite sequences of symbols, a.k.a. stochastic languages. Traditionally , unsupervised learn-ing of PFA is performed through algorithms that iteratively improves the likelihood like the Expectation-Maximization (EM) algorithm. Recently, learning algorithms based on the so-called Method of Moments (MoM) have been proposed as a much faster alternative that comes with PAC-style guarantees. However, these algorithms do not ensure the learnt automata to model a proper distribution , limiting their applicability and preventing them to serve as an initialization to iterative algorithms. In this paper, we propose a new MoM-based algorithm with PAC-style guarantees that learns automata defining proper distributions. We assess its performances on synthetic problems from the PAutomaC challenge and real datasets extracted from Wikipedia against previous MoM-based algorithms and EM algorithm.

### 7.2.4. *Online Kernel and Graph-Based Methods*

**Analysis of Nyström method with sequential ridge leverage score sampling**, [26]

Large-scale kernel ridge regression (KRR) is limited by the need to store a large kernel matrix Kt. To avoid storing the entire matrix Kt, Nystro¨m methods subsample a subset of columns of the kernel matrix, and efficiently find an approximate KRR solution on the reconstructed Kt . The chosen subsampling distribution in turn affects the statistical and computational tradeoffs. For KRR problems, [15, 1] show that a sampling distribution proportional to the ridge leverage scores (RLSs) provides strong reconstruction guarantees for Kt. While exact RLSs are as difficult to compute as a KRR solution, we may be able to approximate them well enough. In this paper, we study KRR problems in a sequential setting and introduce the INK-ESTIMATE algorithm, that incrementally computes the RLSs estimates. INK-ESTIMATE maintains a small sketch of Kt, that at each step is used to compute an intermediate estimate of the RLSs. First, our sketch update does not require access to previously seen columns, and therefore a single pass over the kernel matrix is sufficient. Second, the algorithm requires a fixed, small space budget to run dependent only on the effective dimension of the kernel matrix. Finally, our sketch provides strong approximation guarantees on the distance $||Kt - Kt||^2$

, and on the statistical risk of the approximate KRR solution at any time, because all our guarantees hold at any intermediate step.

## 7.3. Statistical Learning and Bayesian Analysis

### 7.3.1. *Non-parametric methods for Function Approximation*

**Pliable rejection sampling**, [30]

Rejection sampling is a technique for sampling from difficult distributions. However, its use is limited due to a high rejection rate. Common adaptive rejection sampling methods either work only for very specific distributions or without performance guarantees. In this paper, we present pliable rejection sampling (PRS), a new approach to rejection sampling, where we learn the sampling proposal using a kernel estimator. Since our method builds on rejection sampling, the samples obtained are with high probability i.i.d. and distributed according to f. Moreover, PRS comes with a guarantee on the number of accepted samples.

### 7.3.2. *Non-parametric methods for functional supervised learning*

**Operator-valued Kernels for Learning from Functional Response Data**, [16]

In this paper we consider the problems of supervised classification and regression in the case where attributes and labels are functions: a data is represented by a set of functions, and the label is also a function. We focus on the use of reproducing kernel Hilbert space theory to learn from such functional data. Basic concepts and properties of kernel-based learning are extended to include the estimation of function-valued functions. In this setting, the representer theorem is restated, a set of rigorously defined infinite-dimensional operator-valued kernels that can be valuably applied when the data are functions is described, and a learning algorithm for nonlinear functional data analysis is introduced. The methodology is illustrated through speech and audio signal processing experiments.

### 7.3.3. *Differential privacy*

**On the Differential Privacy of Bayesian Inference**, [51]

We study how to communicate findings of Bayesian inference to third parties, while preserving the strong guarantee of differential privacy. Our main contributions are four different algorithms for private Bayesian inference on proba-bilistic graphical models. These include two mechanisms for adding noise to the Bayesian updates, either directly to the posterior parameters, or to their Fourier transform so as to preserve update consistency. We also utilise a recently introduced posterior sampling mechanism, for which we prove bounds for the specific but general case of discrete Bayesian networks; and we introduce a maximum-a-posteriori private mechanism. Our analysis includes utility and privacy bounds, with a novel focus on the influence of graph structure on privacy. Worked examples and experiments with Bayesian naïve Bayes and Bayesian linear regression illustrate the application of our mechanisms.

**Algorithms for Differentially Private Multi-Armed Bandits**, [50]

We present differentially private algorithms for the stochastic Multi-Armed Bandit (MAB) problem. This is a problem for applications such as adaptive clinical trials, experiment design, and user-targeted advertising where private information is connected to individual rewards. Our major contribution is to show that there exist $(\epsilon, \delta)$ differentially private variants of Upper Confidence Bound algorithms which have optimal regret, $O(\epsilon^{-1} + \log T)$. This is a significant improvement over previous results, which only achieve poly-log regret $O(\epsilon^{-2} \log^2 T)$, because of our use of a novel interval-based mechanism. We also substantially improve the bounds of previous family of algorithms which use a continual release mechanism. Experiments clearly validate our theoretical bounds.

## 7.4. Applications

### 7.4.1. *Spoken Dialogue Systems*

**Compact and Interpretable Dialogue State Representation with Genetic Sparse Distributed Memory**, [28]

t User satisfaction is often considered as the objective that should be achieved by spoken dialogue systems. This is why, the reward function of Spoken Dialogue Systems (SDS) trained by Reinforcement Learning (RL) is often designed to reflect user satisfaction. To do so, the state space representation should be based on features capturing user satisfaction characteristics such as the mean speech recognition confidence score for instance. On the other hand, for deployment in industrial systems, there is a need for state representations that are understandable by system engineers. In this paper, we propose to represent the state space using a Genetic Sparse Distributed Memory. This is a state aggregation method computing state prototypes which are selected so as to lead to the best linear representation of the value function in RL. To do so, previous work on Genetic Sparse Distributed Memory for classification is adapted to the Reinforcement Learning task and a new way of building the prototypes is proposed. The approach is tested on a corpus of dialogues collected with an appointment scheduling system. The results are compared to a grid-based linear parametrisation. It is shown that learning is accelerated and made more memory efficient. It is also shown that the framework is calable in that it is possible to include many dialogue features in the representation, interpret the resulting policy and identify the most important dialogue features.

**A Stochastic Model for Computer-Aided Human-Human Dialogue**, [24]

In this paper we introduce a novel model for computer-aided human-human dialogue. In this context, the computer aims at improving the outcome of a human-human task-oriented dialogue by intervening during the course of the interaction. While dialogue state and topic tracking in human-human dialogue have already been studied, few work has been devoted to the sequential part of the problem, where the impact of the system's actions on the future of the conversation is taken into account. This paper addresses this issue by first modelling human-human dialogue as a Markov Reward Process. The task of purposely taking part into the conversation is then optimised within the Linearly Solvable Markov Decision Process framework. Utterances of the Conversational Agent are seen as perturbations in this process, which aim at satisfying the user's long-term goals while keeping the conversation natural. Finally, results obtained by simulation suggest that such an approach is suitable for computer-aided human-human dialogue and is a first step towards three-party dialogue.

**Learning Dialogue Dynamics with the Method of Moments**, [25]

In this paper, we introduce a novel framework to encode the dynamics of dialogues into a probabilistic graphical model. Traditionally, Hidden Markov Models (HMMs) would be used to address this problem, involving a first step of hand-crafting to build a dialogue model (e.g. defining potential hidden states) followed by applying expectation-maximisation (EM) algorithms to refine it. Recently, an alternative class of algorithms based on the Method of Moments (MoM) has proven successful in avoiding issues of the EM-like algorithms such as convergence towards local optima, tractability issues, initialization issues or the lack of theoretical guarantees. In this work, we show that dialogues may be modeled by SP-RFA, a class of graphical models efficiently learnable within the MoM and directly usable in planning algorithms (such as reinforcement learning). Experiments are led on the Ubuntu corpus and dialogues are considered as sequences of dialogue acts, represented by their Latent Dirichlet Allocation (LDA) and Latent Semantic Analysis (LSA). We show that a MoM-based algorithm can learn a compact model of sequences of such acts.

## 7.4.2. *Software development*

**Mutation-Based Graph Inference for Fault Localization**, [45]

We present a new fault localization algorithm, called Vautrin, built on an approximation of causality based on call graphs. The approximation of causality is done using software mutants. The key idea is that if a mutant is killed by a test, certain call graph edges within a path between the mutation point and the failing test are likely causal. We evaluate our approach on the fault localization benchmark by Steimann et al. totaling 5,836 faults. The causal graphs are extracted from 88,732 nodes connected by 119,531 edges. Vautrin improves the fault localization effectiveness for all subjects of the benchmark. Considering the wasted effort at the method level, a classical fault localization evaluation metric, the improvement ranges from 3

**A Large-scale Study of Call Graph-based Impact Prediction using Mutation Testing**, [21]

In software engineering, impact analysis consists in predicting the software elements (e.g. modules, classes, methods) potentially impacted by a change in the source code. Impact analysis is required to optimize the testing effort. In this paper, we propose a framework to predict error propagation. Based on 10 open-source Java projects and 5 classical mutation operators, we create 17000 mutants and study how the error they introduce propagates. This framework enables us to analyze impact prediction based on four types of call graph. Our results show that the sophistication indeed increases completeness of impact prediction. However, and surprisingly to us, the most basic call graph gives the highest trade-off between precision and recall for impact prediction.

**A Learning Algorithm for Change Impact Prediction**, [44]

Change impact analysis (CIA) consists in predicting the impact of a code change in a software application. In this paper, the artifacts that are considered for CIA are methods of object-oriented software; the change under study is a change in the code of the method, the impact is the test methods that fail because of the change that has been performed. We propose LCIP, a learning algorithm that learns from past impacts to predict future impacts. To evaluate LCIP, we consider Java software applications that are strongly tested. We simulate 6000 changes and their actual impact through code mutations, as done in mutation testing. We find that LCIP can predict the impact with a precision of 74

<span style="color:red">SPIRALS Project-Team</span>

# 7. New Results

## 7.1. Change Impact Analysis

In [21], we have proposed a novel evaluation technique for change impact analysis (CIA). CIA is a prediction problem that, given a source code element in a program, determines the other source code elements impacted if one changes this original source code element. Given the large size of the element space in complex programs, this prediction requires a trade-off between different dimensions: precision, completeness, time. The novelty of the result lies in the use of mutation analysis to study simultaneously these three dimensions. This result is backed by an empirical evaluation performed on 10 open-source Java programs and 5 mutation operators, which enabled to generate 17,000 mutants and study how the error they introduce propagates. This result has been achieved in the context of the PhD thesis, defended in November 2016, of Vicenzo Musco [15].

## 7.2. Learning Power Models for Distributed and Virtualized Environments

Energy efficiency is a major concern for modern ICT infrastructures. The a priori estimation of the level of energy consumed by a given service is a difficult problem given the intricate nature of hardware and software that are involved. Consequently, even before considering saving, measuring the exact amount of energy consumed by a given software service or process is required. Over the last few years, a dozen of ad hoc power models have been proposed in the literature. Nevertheless they cannot cope with the constant evolution of software and hardware architecture. We have therefore defined and implemented a toolkit that automatically learns the power models of a given architecture, independently of the features and the complexity it exhibits. This toolkit considers traditional distributed environment as well as virtualized, cloud-based ones. This result has been achieved in the context of the PhD thesis, defended in November 2016, of Maxime Colmant [11].

## 7.3. Crowdmining to Increase the Quality of Software Systems

Modern software systems, especially in the open source world, are more and more part of ecosystems where large quantities of data about these systems are available. These data may come for example from application stores (e.g. Google Play Store or Apple Store for mobile applications), forges (e.g. GitHub), or from the usage conditions experienced by users of these software systems. This large amount of data enables to unlock some specific challenges where knowledge about the software systems can be automatically mined and learnt. In this domain, we obtained new results on the mining of mobile software antipatterns on a crowd of mobile applications and their versions to study their impact on resource consumption [32]. This result has been achieved in the context of the PhD thesis, defended in November 2016, of Geoffrey Hecht [13]. We also consider the crowd of mobile devices and users to detect and reproduce application crashes in the wild. By leveraging our results in the domain of in-breath monitoring, we use the APISENSE® platform (see Section 6.1 ) to collect extended crash reports that can be aggregated to infer the minimal execution path that lead to a crash [28]. This result has been achieved in the context of the PhD thesis, defended in December 2016, of María Gomez Lacruz [12]. These results are also in relation with our activities in the context of the SOMCA associated team (see Section 9.4 ).

## 7.4. Self-Optimization of Virtualized Environments

Elasticity is a major property of virtualized computing environments. In this domain, we especially work at the infrastructure and platform levels of a cloud computing system where we obtained two results that enable to better self-optimize the consumed resources. At the infrastructure level, we proposed CloudGC, a new middleware service for suspending, resuming, and recycling idle virtual machines. The algorithm has been implemented on top of the OpenStack cloud operating system. At the platform level, we proposed a new self-balancing approach to dynamically optimize the performance of the Hadoop framework for the distributed storage and processing of large data sets. These results have been achieved in the context of the PhD thesis, defended in December 2016, of Bo Zhang [16].