

*Inria*

RESEARCH CENTER  
Saclay - Île-de-France

FIELD

Activity Report 2019

# Section Scientific Foundations

Edition: 2020-03-21



ALGORITHMICS, PROGRAMMING, SOFTWARE AND ARCHITECTURE	
1. COMETE Project-Team	4
2. DATASHAPE Project-Team	6
3. DEDUCTEAM Project-Team	8
4. GRACE Project-Team	9
5. MEXICO Project-Team	12
6. PARSIFAL Project-Team	17
7. SPECFUN Project-Team	20
8. TOCCATA Project-Team	25
APPLIED MATHEMATICS, COMPUTATION AND SIMULATION	
9. CELESTE Project-Team	29
10. COMMANDS Project-Team	31
11. DEFI Project-Team	33
12. DISCO Project-Team	37
13. GAMMA Project-Team (section vide)	40
14. POEMS Project-Team	41
15. RANDOPT Project-Team	43
16. TAU Project-Team	48
17. TROPICAL Project-Team	55
DIGITAL HEALTH, BIOLOGY AND EARTH	
18. LIFEWARE Project-Team	57
19. M3DISIM Project-Team	61
20. OPIS Project-Team	62
21. PARIETAL Project-Team	64
22. XPOP Project-Team	68
NETWORKS, SYSTEMS AND SERVICES, DISTRIBUTED COMPUTING	
23. TRIBE Project-Team	73
PERCEPTION, COGNITION AND INTERACTION	
24. AVIZ Project-Team	75
25. CEDAR Project-Team	79
26. EX-SITU Project-Team	81
27. ILDA Project-Team	82
28. PETRUS Project-Team	86

## COMETE Project-Team

### 3. Research Program

#### 3.1. Probability and information theory

**Participants:** Konstantinos Chatzikokolakis, Catuscia Palamidessi, Marco Romanelli, Anna Pazzi.

Much of the research of Comète focuses on security and privacy. In particular, we are interested in the problem of the leakage of secret information through public observables.

Ideally we would like systems to be completely secure, but in practice this goal is often impossible to achieve. Therefore, we need to reason about the amount of information leaked, and the utility that it can have for the adversary, i.e. the probability that the adversary is able to exploit such information.

The recent tendency is to use an information theoretic approach to model the problem and define the leakage in a quantitative way. The idea is to consider the system as an information-theoretic *channel*. The input represents the secret, the output represents the observable, and the correlation between the input and output (*mutual information*) represents the information leakage.

Information theory depends on the notion of entropy as a measure of uncertainty. From the security point of view, this measure corresponds to a particular model of attack and a particular way of estimating the security threat (vulnerability of the secret). Most of the proposals in the literature use Shannon entropy, which is the most established notion of entropy in information theory. We, however, consider also other notions, in particular Rényi min-entropy, which seems to be more appropriate for security in common scenarios like one-try attacks.

#### 3.2. Expressiveness of Concurrent Formalisms

**Participants:** Catuscia Palamidessi, Frank Valencia.

We study computational models and languages for distributed, probabilistic and mobile systems, with a particular attention to expressiveness issues. We aim at developing criteria to assess the expressive power of a model or formalism in a distributed setting, to compare existing models and formalisms, and to define new ones according to an intended level of expressiveness, also taking into account the issue of (efficient) implementability.

#### 3.3. Concurrent constraint programming

**Participants:** Frank Valencia, Santiago Quintero.

Concurrent constraint programming (ccp) is a well established process calculus for modeling systems where agents interact by posting and asking information in a store, much like in users interact in *social networks*. This information is represented as first-order logic formulae, called constraints, on the shared variables of the system (e.g.,  $X > 42$ ). The most distinctive and appealing feature of ccp is perhaps that it unifies in a single formalism the operational view of processes based upon process calculi with a declarative one based upon first-order logic. It also has an elegant denotational semantics that interprets processes as closure operators (over the set of constraints ordered by entailment). In other words, any ccp process can be seen as an idempotent, increasing, and monotonic function from stores to stores. Consequently, ccp processes can be viewed as: computing agents, formulae in the underlying logic, and closure operators. This allows ccp to benefit from the large body of techniques of process calculi, logic and domain theory.

Our research in ccp develops along the following two lines:

1. **(a)** The study of a bisimulation semantics for ccp. The advantage of bisimulation, over other kinds of semantics, is that it can be efficiently verified.
2. **(b)** The extension of ccp with constructs to capture emergent systems such as those in social networks and cloud computing.

### **3.4. Model checking**

**Participants:** Konstantinos Chatzikokolakis, Catuscia Palamidessi.

Model checking addresses the problem of establishing whether a given specification satisfies a certain property. We are interested in developing model-checking techniques for verifying concurrent systems of the kind explained above. In particular, we focus on security and privacy, i.e., on the problem of proving that a given system satisfies the intended security or privacy properties. Since the properties we are interested in have a probabilistic nature, we use probabilistic automata to model the protocols. A challenging problem is represented by the fact that the interplay between nondeterminism and probability, which in security presents subtleties that cannot be handled with the traditional notion of a scheduler,

## **DATASHAPE Project-Team**

### **3. Research Program**

#### **3.1. Algorithmic aspects of topological and geometric data analysis**

TDA requires to construct and manipulate appropriate representations of complex and high dimensional shapes. A major difficulty comes from the fact that the complexity of data structures and algorithms used to approximate shapes rapidly grows as the dimensionality increases, which makes them intractable in high dimensions. We focus our research on simplicial complexes which offer a convenient representation of general shapes and generalize graphs and triangulations. Our work includes the study of simplicial complexes with good approximation properties and the design of compact data structures to represent them.

In low dimensions, effective shape reconstruction techniques exist that can provide precise geometric approximations very efficiently and under reasonable sampling conditions. Extending those techniques to higher dimensions as is required in the context of TDA is problematic since almost all methods in low dimensions rely on the computation of a subdivision of the ambient space. A direct extension of those methods would immediately lead to algorithms whose complexities depend exponentially on the ambient dimension, which is prohibitive in most applications. A first direction to by-pass the curse of dimensionality is to develop algorithms whose complexities depend on the intrinsic dimension of the data (which most of the time is small although unknown) rather than on the dimension of the ambient space. Another direction is to resort to cruder approximations that only captures the homotopy type or the homology of the sampled shape. The recent theory of persistent homology provides a powerful and robust tool to study the homology of sampled spaces in a stable way.

#### **3.2. Statistical aspects of topological and geometric data analysis**

The wide variety of larger and larger available data - often corrupted by noise and outliers - requires to consider the statistical properties of their topological and geometric features and to propose new relevant statistical models for their study.

There exist various statistical and machine learning methods intending to uncover the geometric structure of data. Beyond manifold learning and dimensionality reduction approaches that generally do not allow to assert the relevance of the inferred topological and geometric features and are not well-suited for the analysis of complex topological structures, set estimation methods intend to estimate, from random samples, a set around which the data is concentrated. In these methods, that include support and manifold estimation, principal curves/manifolds and their various generalizations to name a few, the estimation problems are usually considered under losses, such as Hausdorff distance or symmetric difference, that are not sensitive to the topology of the estimated sets, preventing these tools to directly infer topological or geometric information.

Regarding purely topological features, the statistical estimation of homology or homotopy type of compact subsets of Euclidean spaces, has only been considered recently, most of the time under the quite restrictive assumption that the data are randomly sampled from smooth manifolds.

In a more general setting, with the emergence of new geometric inference tools based on the study of distance functions and algebraic topology tools such as persistent homology, computational topology has recently seen an important development offering a new set of methods to infer relevant topological and geometric features of data sampled in general metric spaces. The use of these tools remains widely heuristic and until recently there were only a few preliminary results establishing connections between geometric inference, persistent homology and statistics. However, this direction has attracted a lot of attention over the last three years. In particular, stability properties and new representations of persistent homology information have led to very promising results to which the DATASHAPE members have significantly contributed. These preliminary results open many perspectives and research directions that need to be explored.

Our goal is to build on our first statistical results in TDA to develop the mathematical foundations of Statistical Topological and Geometric Data Analysis. Combined with the other objectives, our ultimate goal is to provide a well-founded and effective statistical toolbox for the understanding of topology and geometry of data.

### 3.3. Topological approach for multimodal data processing

Due to their geometric nature, multimodal data (images, video, 3D shapes, etc.) are of particular interest for the techniques we develop. Our goal is to establish a rigorous framework in which data having different representations can all be processed, mapped and exploited jointly. This requires adapting our tools and sometimes developing entirely new or specialized approaches.

The choice of multimedia data is motivated primarily by the fact that the amount of such data is steadily growing (with e.g. video streaming accounting for nearly two thirds of peak North-American Internet traffic, and almost half a billion images being posted on social networks each day), while at the same time it poses significant challenges in designing informative notions of (dis)-similarity as standard metrics (e.g. Euclidean distances between points) are not relevant.

### 3.4. Experimental research and software development

We develop a high quality open source software platform called GUDHI which is becoming a reference in geometric and topological data analysis in high dimensions. The goal is not to provide code tailored to the numerous potential applications but rather to provide the central data structures and algorithms that underlie applications in geometric and topological data analysis.

The development of the GUDHI platform also serves to benchmark and optimize new algorithmic solutions resulting from our theoretical work. Such development necessitates a whole line of research on software architecture and interface design, heuristics and fine-tuning optimization, robustness and arithmetic issues, and visualization. We aim at providing a full programming environment following the same recipes that made up the success story of the CGAL library, the reference library in computational geometry.

Some of the algorithms implemented on the platform will also be interfaced to other software platform, such as the R software<sup>0</sup> for statistical computing, and languages such as Python in order to make them usable in combination with other data analysis and machine learning tools. A first attempt in this direction has been done with the creation of an R package called TDA in collaboration with the group of Larry Wasserman at Carnegie Mellon University (Inria Associated team CATS) that already includes some functionalities of the GUDHI library and implements some joint results between our team and the CMU team. A similar interface with the Python language is also considered a priority. To go even further towards helping users, we will provide utilities that perform the most common tasks without requiring any programming at all.

---

<sup>0</sup><https://www.r-project.org/>

---

## DEDUCTEAM Project-Team

### 3. Research Program

#### 3.1. Logical Frameworks

A thesis, which is at the root of our research effort, is that logical systems should be expressed as theories in a logical framework. As a consequence, proof-checking systems should not be focused on one theory, such as Simple type theory, Martin-Löf's type theory, or the Calculus of constructions, but should be theory independent. On the more theoretical side, the proof search algorithms, or the algorithmic interpretation of proofs should not depend on the theory in which proofs are expressed, but this theory should just be a parameter. This is for instance expressed in the title of our invited talk at ICALP 2012: *A theory independent Curry-De Bruijn-Howard correspondence* [25].

Various limits of Predicate logic have led to the development of various families of logical frameworks:  $\lambda$ -prolog and Isabelle have allowed terms containing free variables, the Edinburgh logical framework has allowed proofs to be expressed as  $\lambda$ -terms, Pure type systems have allowed propositions to be considered as terms, and Deduction modulo theory has allowed theories to be defined not only with axioms, but also with computation rules.

The  $\lambda\Pi$ -calculus modulo theory, that is implemented in the system DEDUKTI and that is a synthesis of the Edinburgh logical framework and of Deduction modulo theory, subsumes them all. Part of our research effort is focused on improving the  $\lambda\Pi$ -calculus modulo theory, for instance allowing to define congruences with associative and commutative rewriting. Another part of our research effort is focused on the automatic analysis of theories to prove their confluence, termination, and consistency either by pencil and paper proofs or automatically [4].

#### 3.2. Interoperability and proof encyclopaediae

Using a single prover to check proofs coming from different systems naturally leads to investigate how these proofs can be translated from one theory to another and used in a system different from the system in which they have been developed. This issue is of prime importance because developments in proof systems are getting bigger and, unlike other communities in computer science, the proof checking community has given little effort in the direction of standardization and interoperability.

For each proof, independently of the system in which it has been developed, we should be able to identify the systems in which it can be expressed. For instance, we have shown that many proofs developed in the MATITA prover did not use the full strength of the logic of MATITA and could be exported, for instance, to the systems of the HOL family, that are based on a weaker logic.

Rather than importing proofs from one system, transforming them, and exporting them to another system, we can use the same tools to develop system-independent proof encyclopaedia called Logipedia. In such a library, each proof is labeled with the theories in which it can be expressed and so with the systems in which it can be used.

#### 3.3. Interactive theorem proving

If our main goal with DEDUKTI is to import, transform, and export proofs developed in other systems, we also want to investigate how DEDUKTI can be used as the basis of an interactive theorem prover. This leads to two new scientific questions: first, how much can a tactic system be theory independent, and then how does rewriting extends the possibility to write tactics.

This has led to the development of a new version of DEDUKTI, which supports metavariables. Several tactics have been developed for this system, which are intended to help a human user to write proofs in our system instead of writing proof terms by hand. This work is a continuation of the previous work the team did on DEMON, which was an extension of DEDUKTI, whereas the support for interactive theorem proving is now native in DEDUKTI.



## GRACE Project-Team

### 3. Research Program

#### 3.1. Algorithmic Number Theory

**Participants:** Luca de Feo, François Morain, Benjamin Smith, Mathilde de La Morinerie, Antonin Leroux, Guénaél Renault.

Algorithmic Number Theory is concerned with replacing special cases with general algorithms to solve problems in number theory. In the Grace project, it appears in three main threads:

- fundamental algorithms for integers and polynomials (including primality and factorization);
- algorithms for finite fields (including discrete logarithms);
- algorithms for algebraic curves.

Clearly, we use computer algebra in many ways. Research in cryptology has motivated a renewed interest in Algorithmic Number Theory in recent decades—but the fundamental problems still exist *per se*. Indeed, while algorithmic number theory application in cryptanalysis is epitomized by applying factorization to breaking RSA public key, many other problems, are relevant to various area of computer science. Roughly speaking, the problems of the cryptological world are of bounded size, whereas Algorithmic Number Theory is also concerned with asymptotic results.

#### 3.2. Arithmetic Geometry: Curves and their Jacobians

**Participants:** Luca de Feo, François Morain, Benjamin Smith, Mathilde de La Morinerie, Antonin Leroux.

Theme: Arithmetic Geometry: Curves and their Jacobians *Arithmetic Geometry* is the meeting point of algebraic geometry and number theory: that is, the study of geometric objects defined over arithmetic number systems (such as the integers and finite fields). The fundamental objects for our applications in both coding theory and cryptology are curves and their Jacobians over finite fields.

An algebraic *plane curve*  $\mathcal{X}$  over a field  $\mathbf{K}$  is defined by an equation

$$\mathcal{X} : F_{\mathcal{X}}(x, y) = 0 \quad \text{where } F_{\mathcal{X}} \in \mathbf{K}[x, y].$$

(Not every curve is planar—we may have more variables, and more defining equations—but from an algorithmic point of view, we can always reduce to the plane setting.) The *genus*  $g_{\mathcal{X}}$  of  $\mathcal{X}$  is a non-negative integer classifying the essential geometric complexity of  $\mathcal{X}$ ; it depends on the degree of  $F_{\mathcal{X}}$  and on the number of singularities of  $\mathcal{X}$ . The curve  $\mathcal{X}$  is associated in a functorial way with an algebraic group  $J_{\mathcal{X}}$ , called the *Jacobian* of  $\mathcal{X}$ . The group  $J_{\mathcal{X}}$  has a geometric structure: its elements correspond to points on a  $g_{\mathcal{X}}$ -dimensional projective algebraic group variety. Typically, we do not compute with the equations defining this projective variety: there are too many of them, in too many variables, for this to be convenient. Instead, we use fast algorithms based on the representation in terms of classes of formal sums of points on  $\mathcal{X}$ .

The simplest curves with nontrivial Jacobians are curves of genus 1, known as *elliptic curves*; they are typically defined by equations of the form  $y^2 = x^3 + Ax + B$ . Elliptic curves are particularly important given their central role in public-key cryptography over the past two decades. Curves of higher genus are important in both cryptography and coding theory.

#### 3.3. Curve-Based cryptology

**Participants:** Luca de Feo, François Morain, Benjamin Smith, Mathilde de La Morinerie, Antonin Leroux.

Theme: Curve-Based Cryptology

Jacobians of curves are excellent candidates for cryptographic groups when constructing efficient instances of public-key cryptosystems. Diffie–Hellman key exchange is an instructive example.

Suppose Alice and Bob want to establish a secure communication channel. Essentially, this means establishing a common secret *key*, which they will then use for encryption and decryption. Some decades ago, they would have exchanged this key in person, or through some trusted intermediary; in the modern, networked world, this is typically impossible, and in any case completely unscalable. Alice and Bob may be anonymous parties who want to do e-business, for example, in which case they cannot securely meet, and they have no way to be sure of each other’s identities. Diffie–Hellman key exchange solves this problem. First, Alice and Bob publicly agree on a cryptographic group  $G$  with a generator  $P$  (of order  $N$ ); then Alice secretly chooses an integer  $a$  from  $[1..N]$ , and sends  $aP$  to Bob. In the meantime, Bob secretly chooses an integer  $b$  from  $[1..N]$ , and sends  $bP$  to Alice. Alice then computes  $a(bP)$ , while Bob computes  $b(aP)$ ; both have now computed  $abP$ , which becomes their shared secret key. The security of this key depends on the difficulty of computing  $abP$  given  $P$ ,  $aP$ , and  $bP$ ; this is the Computational Diffie–Hellman Problem (CDHP). In practice, the CDHP corresponds to the Discrete Logarithm Problem (DLP), which is to determine  $a$  given  $P$  and  $aP$ .

This simple protocol has been in use, with only minor modifications, since the 1970s. The challenge is to create examples of groups  $G$  with a relatively compact representation and an efficiently computable group law, and such that the DLP in  $G$  is hard (ideally approaching the exponential difficulty of the DLP in an abstract group). The Pohlig–Hellman reduction shows that the DLP in  $G$  is essentially only as hard as the DLP in its largest prime-order subgroup. We therefore look for compact and efficient groups of prime order.

The classic example of a group suitable for the Diffie–Hellman protocol is the multiplicative group of a finite field  $\mathbf{F}_q$ . There are two problems that render its usage somewhat less than ideal. First, it has too much structure: we have a subexponential Index Calculus attack on the DLP in this group, so while it is very hard, the DLP falls a long way short of the exponential difficulty of the DLP in an abstract group. Second, there is only one such group for each  $q$ : its subgroup treillis depends only on the factorization of  $q - 1$ , and requiring  $q - 1$  to have a large prime factor eliminates many convenient choices of  $q$ .

This is where Jacobians of algebraic curves come into their own. First, elliptic curves and Jacobians of genus 2 curves do not have a subexponential index calculus algorithm: in particular, from the point of view of the DLP, a generic elliptic curve is currently *as strong as* a generic group of the same size. Second, they provide some diversity: we have many degrees of freedom in choosing curves over a fixed  $\mathbf{F}_q$ , with a consequent diversity of possible cryptographic group orders. Furthermore, an attack which leaves one curve vulnerable may not necessarily apply to other curves. Third, viewing a Jacobian as a geometric object rather than a pure group allows us to take advantage of a number of special features of Jacobians. These features include efficiently computable pairings, geometric transformations for optimised group laws, and the availability of efficiently computable non-integer endomorphisms for accelerated encryption and decryption.

### 3.4. Algebraic Coding Theory

**Participants:** Daniel Augot, Alain Couvreur, Françoise Levy-Dit-Vehel, Maxime Roméas, Sarah Bordage, Adrien Hauteville, Isabella Panaccione.

Theme: Coding theory

Coding Theory studies originated with the idea of using redundancy in messages to protect against noise and errors. The last decade of the 20th century has seen the success of so-called iterative decoding methods, which enable us to get very close to the Shannon capacity. The capacity of a given channel is the best achievable transmission rate for reliable transmission. The consensus in the community is that this capacity is more easily reached with these iterative and probabilistic methods than with algebraic codes (such as Reed–Solomon codes).

However, algebraic coding is useful in settings other than the Shannon context. Indeed, the Shannon setting is a random case setting, and promises only a vanishing error probability. In contrast, the algebraic Hamming approach is a worst case approach: under combinatorial restrictions on the noise, the noise can be adversarial, with strictly zero errors.

These considerations are renewed by the topic of list decoding after the breakthrough of Guruswami and Sudan at the end of the nineties. List decoding relaxes the uniqueness requirement of decoding, allowing a small list of candidates to be returned instead of a single codeword. List decoding can reach a capacity close to the Shannon capacity, with zero failure, with small lists, in the adversarial case. The method of Guruswami and Sudan enabled list decoding of most of the main algebraic codes: Reed–Solomon codes and Algebraic–Geometry (AG) codes and new related constructions “capacity-achieving list decodable codes”. These results open the way to applications against adversarial channels, which correspond to worst case settings in the classical computer science language.

Another avenue of our studies is AG codes over various geometric objects. Although Reed–Solomon codes are the best possible codes for a given alphabet, they are very limited in their length, which cannot exceed the size of the alphabet. AG codes circumvent this limitation, using the theory of algebraic curves over finite fields to construct long codes over a fixed alphabet. The striking result of Tsfasman–Vladut–Zink showed that codes better than random codes can be built this way, for medium to large alphabets. Disregarding the asymptotic aspects and considering only finite length, AG codes can be used either for longer codes with the same alphabet, or for codes with the same length with a smaller alphabet (and thus faster underlying arithmetic).

From a broader point of view, wherever Reed–Solomon codes are used, we can substitute AG codes with some benefits: either beating random constructions, or beating Reed–Solomon codes which are of bounded length for a given alphabet.

Another area of Algebraic Coding Theory with which we are more recently concerned is the one of Locally Decodable Codes. After having been first theoretically introduced, those codes now begin to find practical applications, most notably in cloud-based remote storage systems.

## MEXICO Project-Team

### 3. Research Program

#### 3.1. Concurrency

**Participants:** Thomas Chatain, Philippe Dague, Stefan Haar, Serge Haddad, Stefan Schwoon.

Concurrency; Semantics; Automatic Control ; Diagnosis ; Verification

**Concurrency:** Property of systems allowing some interacting processes to be executed in parallel.

**Diagnosis:** The process of deducing from a partial observation of a system aspects of the internal states or events of that system; in particular, *fault diagnosis* aims at determining whether or not some non-observable fault event has occurred.

**Conformance Testing:** Feeding dedicated input into an implemented system  $IS$  and deducing, from the resulting output of  $I$ , whether  $I$  respects a formal specification  $S$ .

##### 3.1.1. Introduction

It is well known that, whatever the intended form of analysis or control, a *global* view of the system state leads to overwhelming numbers of states and transitions, thus slowing down algorithms that need to explore the state space. Worse yet, it often blurs the mechanics that are at work rather than exhibiting them. Conversely, respecting concurrency relations avoids exhaustive enumeration of interleavings. It allows us to focus on ‘essential’ properties of non-sequential processes, which are expressible with causal precedence relations. These precedence relations are usually called causal (partial) orders. Concurrency is the explicit absence of such a precedence between actions that do not have to wait for one another. Both causal orders and concurrency are in fact essential elements of a specification. This is especially true when the specification is constructed in a distributed and modular way. Making these ordering relations explicit requires to leave the framework of state/interleaving based semantics. Therefore, we need to develop new dedicated algorithms for tasks such as conformance testing, fault diagnosis, or control for distributed discrete systems. Existing solutions for these problems often rely on centralized sequential models which do not scale up well.

##### 3.1.2. Diagnosis

**Participants:** Stefan Haar, Serge Haddad, Stefan Schwoon, Philippe Dague, Lina Ye.

*Fault Diagnosis* for discrete event systems is a crucial task in automatic control. Our focus is on *event oriented* (as opposed to *state oriented*) model-based diagnosis, asking e.g. the following questions: given a - potentially large - *alarm pattern* formed of observations,

- what are the possible *fault scenarios* in the system that *explain* the pattern ?
- Based on the observations, can we deduce whether or not a certain - invisible - fault has actually occurred ?

Model-based diagnosis starts from a discrete event model of the observed system - or rather, its relevant aspects, such as possible fault propagations, abstracting away other dimensions. From this model, an extraction or unfolding process, guided by the observation, produces recursively the explanation candidates.

In asynchronous partial-order based diagnosis with Petri nets [45], [46], [47], one unfolds the *labelled product* of a Petri net model  $\mathcal{N}$  and an observed alarm pattern  $\mathcal{A}$ , also in Petri net form. We obtain an acyclic net giving partial order representation of the behaviors compatible with the alarm pattern. A recursive online procedure filters out those runs (*configurations*) that explain *exactly*  $\mathcal{A}$ . The Petri-net based approach generalizes to dynamically evolving topologies, in dynamical systems modeled by graph grammars, see [34]

### 3.1.2.1. Observability and Diagnosability

Diagnosis algorithms have to operate in contexts with low observability, i.e., in systems where many events are invisible to the supervisor. Checking *observability* and *diagnosability* for the supervised systems is therefore a crucial and non-trivial task in its own right. Analysis of the relational structure of occurrence nets allows us to check whether the system exhibits sufficient visibility to allow diagnosis. Developing efficient methods for both verification of *diagnosability checking* under concurrency, and the *diagnosis* itself for distributed, composite and asynchronous systems, is an important field for *MExICO*. In 2019, a new property, manifestability, weaker than diagnosability (dual in some sense to opacity) has been studied in the context of automata and timed automata.

### 3.1.2.2. Distribution

Distributed computation of unfoldings allows one to factor the unfolding of the global system into smaller *local* unfoldings, by local supervisors associated with sub-networks and communicating among each other. In [46], [36], elements of a methodology for distributed computation of unfoldings between several supervisors, underwritten by algebraic properties of the category of Petri nets have been developed. Generalizations, in particular to Graph Grammars, are still to be done.

Computing diagnosis in a distributed way is only one aspect of a much vaster topic, that of *distributed diagnosis* (see [43], [49]). In fact, it involves a more abstract and often indirect reasoning to conclude whether or not some given invisible fault has occurred. Combination of local scenarios is in general not sufficient: the global system may have behaviors that do not reveal themselves as faulty (or, dually, non-faulty) on any local supervisor's domain (compare [33], [39]). Rather, the local diagnosers have to join all *information* that is available to them locally, and then deduce collectively further information from the combination of their views. In particular, even the *absence* of fault evidence on all peers may allow to deduce fault occurrence jointly, see [51], [52]. Automating such procedures for the supervision and management of distributed and locally monitored asynchronous systems is a long-term goal to which *MExICO* hopes to contribute.

### 3.1.3. Hybrid Systems

**Participants:** Philippe Dague, Lina Ye, Serge Haddad.

Hybrid systems constitute a model for cyber-physical systems which integrates continuous-time dynamics (modes) governed by differential equations, and discrete transitions which switch instantaneously from one mode to another. Thanks to their ease of programming, hybrid systems have been integrated to power electronics systems, and more generally in cyber-physical systems. In order to guarantee that such systems meet their specifications, classical methods consist in finitely abstracting the systems by discretization of the (infinite) state space, and deriving automatically the appropriate mode control from the specification using standard graph techniques.

Diagnosability of hybrid systems has also been studied through an abstraction / refinement process in terms of timed automata.

### 3.1.4. Contextual Nets

**Participant:** Stefan Schwoon.

Assuring the correctness of concurrent systems is notoriously difficult due to the many unforeseeable ways in which the components may interact and the resulting state-space explosion. A well-established approach to alleviate this problem is to model concurrent systems as Petri nets and analyse their unfoldings, essentially an acyclic version of the Petri net whose simpler structure permits easier analysis [44].

However, Petri nets are inadequate to model concurrent read accesses to the same resource. Such situations often arise naturally, for instance in concurrent databases or in asynchronous circuits. The encoding tricks typically used to model these cases in Petri nets make the unfolding technique inefficient. Contextual nets, which explicitly do model concurrent read accesses, address this problem. Their accurate representation of concurrency makes contextual unfoldings up to exponentially smaller in certain situations. An abstract algorithm for contextual unfoldings was first given in [35]. In recent work, we further studied this subject

from a theoretical and practical perspective, allowing us to develop concrete, efficient data structures and algorithms and a tool (Cunf) that improves upon existing state of the art. This work led to the PhD thesis of César Rodríguez in 2014 .

Contextual unfoldings deal well with two sources of state-space explosion: concurrency and shared resources. Recently, we proposed an improved data structure, called *contextual merged processes* (CMP) to deal with a third source of state-space explosion, i.e. sequences of choices. The work on CMP [53] is currently at an abstract level. In the short term, we want to put this work into practice, requiring some theoretical groundwork, as well as programming and experimentation.

Another well-known approach to verifying concurrent systems is *partial-order reduction*, exemplified by the tool SPIN. Although it is known that both partial-order reduction and unfoldings have their respective strengths and weaknesses, we are not aware of any conclusive comparison between the two techniques. Spin comes with a high-level modeling language having an explicit notion of processes, communication channels, and variables. Indeed, the reduction techniques implemented in Spin exploit the specific properties of these features. On the other side, while there exist highly efficient tools for unfoldings, Petri nets are a relatively general low-level formalism, so these techniques do not exploit properties of higher language features. Our work on contextual unfoldings and CMPs represents a first step to make unfoldings exploit richer models. In the long run, we wish raise the unfolding technique to a suitable high-level modelling language and develop appropriate tool support.

## 3.2. Management of Quantitative Behavior

**Participants:** Thomas Chatain, Stefan Haar, Serge Haddad.

### 3.2.1. Introduction

Besides the logical functionalities of programs, the *quantitative* aspects of component behavior and interaction play an increasingly important role.

- *Real-time* properties cannot be neglected even if time is not an explicit functional issue, since transmission delays, parallelism, etc, can lead to time-outs striking, and thus change even the logical course of processes. Again, this phenomenon arises in telecommunications and web services, but also in transport systems.
- In the same contexts, *probabilities* need to be taken into account, for many diverse reasons such as unpredictable functionalities, or because the outcome of a computation may be governed by race conditions.
- Last but not least, constraints on *cost* cannot be ignored, be it in terms of money or any other limited resource, such as memory space or available CPU time.

Traditional mainframe systems were proprietary and (essentially) localized; therefore, impact of delays, unforeseen failures, etc. could be considered under the control of the system manager. It was therefore natural, in verification and control of systems, to focus on *functional* behavior entirely.

With the increase in size of computing system and the growing degree of compositionality and distribution, quantitative factors enter the stage:

- calling remote services and transmitting data over the web creates *delays*;
- remote or non-proprietary components are not “deterministic”, in the sense that their behavior is uncertain.

*Time* and *probability* are thus parameters that management of distributed systems must be able to handle; along with both, the *cost* of operations is often subject to restrictions, or its minimization is at least desired. The mathematical treatment of these features in distributed systems is an important challenge, which *MEICO* is addressing; the following describes our activities concerning probabilistic and timed systems. Note that cost optimization is not a current activity but enters the picture in several intended activities.

### 3.2.2. Probabilistic distributed Systems

**Participants:** Stefan Haar, Serge Haddad.

#### 3.2.2.1. Non-sequential probabilistic processes

Practical fault diagnosis requires to select explanations of *maximal likelihood*. For partial-order based diagnosis, this leads therefore to the question what the probability of a given partially ordered execution is. In Benveniste et al. [38], [31], we presented a model of stochastic processes, whose trajectories are partially ordered, based on local branching in Petri net unfoldings; an alternative and complementary model based on Markov fields is developed in [48], which takes a different view on the semantics and overcomes the first model's restrictions on applicability.

Both approaches abstract away from real time progress and randomize choices in *logical* time. On the other hand, the relative speed - and thus, indirectly, the real-time behavior of the system's local processes - are crucial factors determining the outcome of probabilistic choices, even if non-determinism is absent from the system.

In another line of research [40] we have studied the likelihood of occurrence of non-sequential runs under random durations in a stochastic Petri net setting. It remains to better understand the properties of the probability measures thus obtained, to relate them with the models in logical time, and exploit them e.g. in *diagnosis*.

#### 3.2.2.2. Distributed Markov Decision Processes

**Participant:** Serge Haddad.

Distributed systems featuring non-deterministic and probabilistic aspects are usually hard to analyze and, more specifically, to optimize. Furthermore, high complexity theoretical lower bounds have been established for models like partially observed Markovian decision processes and distributed partially observed Markovian decision processes. We believe that these negative results are consequences of the choice of the models rather than the intrinsic complexity of problems to be solved. Thus we plan to introduce new models in which the associated optimization problems can be solved in a more efficient way. More precisely, we start by studying connection protocols weighted by costs and we look for online and offline strategies for optimizing the mean cost to achieve the protocol. We have been cooperating on this subject with the SUMO team at Inria Rennes; in the joint work [32]; there, we strive to synthesize for a given MDP a control so as to guarantee a specific stationary behavior, rather than - as is usually done - so as to maximize some reward.

### 3.2.3. Large scale probabilistic systems

Addressing large-scale probabilistic systems requires to face state explosion, due to both the discrete part and the probabilistic part of the model. In order to deal with such systems, different approaches have been proposed:

- Restricting the synchronization between the components as in queuing networks allows to express the steady-state distribution of the model by an analytical formula called a product-form [37].
- Some methods that tackle with the combinatory explosion for discrete-event systems can be generalized to stochastic systems using an appropriate theory. For instance symmetry based methods have been generalized to stochastic systems with the help of aggregation theory [42].
- At last simulation, which works as soon as a stochastic operational semantic is defined, has been adapted to perform statistical model checking. Roughly speaking, it consists to produce a confidence interval for the probability that a random path fulfills a formula of some temporal logic [54].

We want to contribute to these three axes: (1) we are looking for product-forms related to systems where synchronization are more involved (like in Petri nets [6]); (2) we want to adapt methods for discrete-event systems that require some theoretical developments in the stochastic framework and, (3) we plan to address some important limitations of statistical model checking like the expressiveness of the associated logic and the handling of rare events.

### 3.2.4. Real time distributed systems

Nowadays, software systems largely depend on complex timing constraints and usually consist of many interacting local components. Among them, railway crossings, traffic control units, mobile phones, computer servers, and many more safety-critical systems are subject to particular quality standards. It is therefore becoming increasingly important to look at networks of timed systems, which allow real-time systems to operate in a distributed manner.

Timed automata are a well-studied formalism to describe reactive systems that come with timing constraints. For modeling distributed real-time systems, networks of timed automata have been considered, where the local clocks of the processes usually evolve at the same rate [50] [41]. It is, however, not always adequate to assume that distributed components of a system obey a global time. Actually, there is generally no reason to assume that different timed systems in the networks refer to the same time or evolve at the same rate. Any component is rather determined by local influences such as temperature and workload.

#### 3.2.4.1. Implementation of Real-Time Concurrent Systems

**Participants:** Thomas Chatain, Stefan Haar, Serge Haddad.

This was one of the tasks of the ANR ImpRo.

Formal models for real-time systems, like timed automata and time Petri nets, have been extensively studied and have proved their interest for the verification of real-time systems. On the other hand, the question of using these models as specifications for designing real-time systems raises some difficulties. One of those comes from the fact that the real-time constraints introduce some artifacts and because of them some syntactically correct models have a formal semantics that is clearly unrealistic. One famous situation is the case of Zeno executions, where the formal semantics allows the system to do infinitely many actions in finite time. But there are other problems, and some of them are related to the distributed nature of the system. These are the ones we address here.

One approach to implementability problems is to formalize either syntactical or behavioral requirements about what should be considered as a reasonable model, and reject other models. Another approach is to adapt the formal semantics such that only realistic behaviors are considered.

These techniques are preliminaries for dealing with the problem of implementability of models. Indeed implementing a model may be possible at the cost of some transformation, which make it suitable for the target device. By the way these transformations may be of interest for the designer who can now use high-level features in a model of a system or protocol, and rely on the transformation to make it implementable.

We aim at formalizing and automating translations that preserve both the timed semantics and the concurrent semantics. This effort is crucial for extending concurrency-oriented methods for logical time, in particular for exploiting partial order properties. In fact, validation and management - in a broad sense - of distributed systems is not realistic *in general* without understanding and control of their real-time dependent features; the link between real-time and logical-time behaviors is thus crucial for many aspects of *MEXICO*'s work.



## PARSIFAL Project-Team

### 3. Research Program

#### 3.1. General overview

There are two broad approaches for computational specifications. In the *computation as model* approach, computations are encoded as mathematical structures containing nodes, transitions, and state. Logic is used to *describe* these structures, that is, the computations are used as models for logical expressions. Intensional operators, such as the modals of temporal and dynamic logics or the triples of Hoare logic, are often employed to express propositions about the change in state.

The *computation as deduction* approach, in contrast, expresses computations logically, using formulas, terms, types, and proofs as computational elements. Unlike the model approach, general logical apparatus such as cut-elimination or automated deduction becomes directly applicable as tools for defining, analyzing, and animating computations. Indeed, we can identify two main aspects of logical specifications that have been very fruitful:

- *Proof normalization*, which treats the state of a computation as a proof term and computation as normalization of the proof terms. General reduction principles such as  $\beta$ -reduction or cut-elimination are merely particular forms of proof normalization. Functional programming is based on normalization [51], and normalization in different logics can justify the design of new and different functional programming languages [32].
- *Proof search*, which views the state of a computation as a structured collection of formulas, known as a *sequent*, and proof search in a suitable sequent calculus as encoding the dynamics of the computation. Logic programming is based on proof search [55], and different proof search strategies can be used to justify the design of new and different logic programming languages [54].

While the distinction between these two aspects is somewhat informal, it helps to identify and classify different concerns that arise in computational semantics. For instance, confluence and termination of reductions are crucial considerations for normalization, while unification and strategies are important for search. A key challenge of computational logic is to find means of uniting or reorganizing these apparently disjoint concerns.

An important organizational principle is structural proof theory, that is, the study of proofs as syntactic, algebraic and combinatorial objects. Formal proofs often have equivalences in their syntactic representations, leading to an important research question about *canonicity* in proofs – when are two proofs “essentially the same?” The syntactic equivalences can be used to derive normal forms for proofs that illuminate not only the proofs of a given formula, but also its entire proof search space. The celebrated *focusing* theorem of Andreoli [34] identifies one such normal form for derivations in the sequent calculus that has many important consequences both for search and for computation. The combinatorial structure of proofs can be further explored with the use of *deep inference*; in particular, deep inference allows access to simple and manifestly correct cut-elimination procedures with precise complexity bounds.

Type theory is another important organizational principle, but most popular type systems are generally designed for either search or for normalization. To give some examples, the Coq system [60] that implements the Calculus of Inductive Constructions (CIC) is designed to facilitate the expression of computational features of proofs directly as executable functional programs, but general proof search techniques for Coq are rather primitive. In contrast, the Twelf system [57] that is based on the LF type theory (a subsystem of the CIC), is based on relational specifications in canonical form (*i.e.*, without redexes) for which there are sophisticated automated reasoning systems such as meta-theoretic analysis tools, logic programming engines, and inductive theorem provers. In recent years, there has been a push towards combining search and normalization in the same type-theoretic framework. The Beluga system [58], for example, is an extension of the LF type theory with a purely computational meta-framework where operations on inductively defined LF objects can be expressed as functional programs.

The Parsifal team investigates both the search and the normalization aspects of computational specifications using the concepts, results, and insights from proof theory and type theory.

### 3.2. Inductive and co-inductive reasoning

The team has spent a number of years in designing a strong new logic that can be used to reason (inductively and co-inductively) on syntactic expressions containing bindings. This work is based on earlier work by McDowell, Miller, and Tiu [53] [52] [56] [61], and on more recent work by Gacek, Miller, and Nadathur [41] [40]. The Parsifal team, along with our colleagues in Minneapolis, Canberra, Singapore, and Cachan, have been building two tools that exploit the novel features of this logic. These two systems are the following.

- Abella, which is an interactive theorem prover for the full logic.
- Bedwyr, which is a model checker for the “finite” part of the logic.

We have used these systems to provide formalize reasoning of a number of complex formal systems, ranging from programming languages to the  $\lambda$ -calculus and  $\pi$ -calculus.

Since 2014, the Abella system has been extended with a number of new features. A number of new significant examples have been implemented in Abella and an extensive tutorial for it has been written [1].

### 3.3. Developing a foundational approach to defining proof evidence

The team is developing a framework for defining the semantics of proof evidence. With this framework, implementers of theorem provers can output proof evidence in a format of their choice: they will only need to be able to formally define that evidence’s semantics. With such semantics provided, proof checkers can then check alleged proofs for correctness. Thus, anyone who needs to trust proofs from various provers can put their energies into designing trustworthy checkers that can execute the semantic specification.

In order to provide our framework with the flexibility that this ambitious plan requires, we have based our design on the most recent advances within the theory of proofs. For a number of years, various team members have been contributing to the design and theory of *focused proof systems* [35] [37] [38] [39] [43] [49] [50] and we have adopted such proof systems as the corner stone for our framework.

We have also been working for a number of years on the implementation of computational logic systems, involving, for example, both unification and backtracking search. As a result, we are also building an early and reference implementation of our semantic definitions.

### 3.4. Deep inference

Deep inference [44], [46] is a novel methodology for presenting deductive systems. Unlike traditional formalisms like the sequent calculus, it allows rewriting of formulas deep inside arbitrary contexts. The new freedom for designing inference rules creates a richer proof theory. For example, for systems using deep inference, we have a greater variety of normal forms for proofs than in sequent calculus or natural deduction systems. Another advantage of deep inference systems is the close relationship to category-theoretic proof theory. Due to the deep inference design one can directly read off the morphism from the derivations. There is no need for a counter-intuitive translation.

The following research problems are investigated by members of the Parsifal team:

- Find deep inference system for richer logics. This is necessary for making the proof theoretic results of deep inference accessible to applications as they are described in the previous sections of this report.
- Investigate the possibility of focusing proofs in deep inference. As described before, focusing is a way to reduce the non-determinism in proof search. However, it is well investigated only for the sequent calculus. In order to apply deep inference in proof search, we need to develop a theory of focusing for deep inference.

### 3.5. Proof nets, atomic flows, and combinatorial proofs

*Proof nets* graph-like presentations of sequent calculus proofs such that all "trivial rule permutations" are quotiented away. Ideally the notion of proof net should be independent from any syntactic formalism, but most notions of proof nets proposed in the past were formulated in terms of their relation to the sequent calculus. Consequently we could observe features like "boxes" and explicit "contraction links". The latter appeared not only in Girard's proof nets [42] for linear logic but also in Robinson's proof nets [59] for classical logic. In this kind of proof nets every link in the net corresponds to a rule application in the sequent calculus.

Only recently, due to the rise of deep inference, new kinds of proof nets have been introduced that take the formula trees of the conclusions and add additional "flow-graph" information (see e.g., [48][2] leading to the notion of *atomic flow* and [45]). On one side, this gives new insights in the essence of proofs and their normalization. But on the other side, all the known correctness criteria are no longer available.

*Combinatorial proofs* [47] are another form syntax-independent proof presentation which separates the multiplicative from the additive behaviour of classical connectives.

The following research questions investigated by members of the Parsifal team:

- Finding (for classical and intuitionistic logic) a notion of canonical proof presentation that is deductive, i.e., can effectively be used for doing proof search.
- Studying the normalization of proofs using atomic flows and combinatorial proofs, as they simplify the normalization procedure for proofs in deep inference, and additionally allow to get new insights in the complexity of the normalization.
- Studying the size of proofs in the combinatorial proof formalism.

### 3.6. Cost Models and Abstract Machines for Functional Programs

In the *proof normalization* approach, computation is usually reformulated as the evaluation of functional programs, expressed as terms in a variation over the  $\lambda$ -calculus. Thanks to its higher-order nature, this approach provides very concise and abstract specifications. Its strength is however also its weakness: the abstraction from physical machines is pushed to a level where it is no longer clear how to measure the complexity of an algorithm.

Models like Turing machines or RAM rely on atomic computational steps and thus admit quite obvious cost models for time and space. The  $\lambda$ -calculus instead relies on a single non-atomic operation,  $\beta$ -reduction, for which costs in terms of time and space are far from evident.

Nonetheless, it turns out that the number of  $\beta$ -steps is a reasonable time cost model, i.e., it is polynomially related to those of Turing machines and RAM. For the special case of *weak evaluation* (i.e., reducing only  $\beta$ -steps that are not under abstractions)—which is used to model functional programming languages—this is a relatively old result due to Blleloch and Greiner [36] (1995). It is only very recently (2014) that the strong case—used in the implementation models of proof assistants—has been solved by Accattoli and Dal Lago [33].

With the recent recruitment of Accattoli, the team's research has expanded in this direction. The topics under investigations are:

1. *Complexity of Abstract Machines.* Bounding and comparing the overhead of different abstract machines for different evaluation schemas (weak/strong call-by-name/value/need  $\lambda$ -calculi) with respect to the cost model. The aim is the development of a complexity-aware theory of the implementation of functional programs.
2. *Reasonable Space Cost Models.* Essentially nothing is known about reasonable space cost models. It is known, however, that environment-based execution model—which are the mainstream technology for functional programs—do not provide an answer. We are exploring the use of the non-standard implementation models provided by Girard's Geometry of Interaction to address this question.

## SPECFUN Project-Team

### 3. Research Program

#### 3.1. Studying special functions by computer algebra

Computer algebra manipulates symbolic representations of exact mathematical objects in a computer, in order to perform computations and operations like simplifying expressions and solving equations for “closed-form expressions”. The manipulations are often fundamentally of algebraic nature, even when the ultimate goal is analytic. The issue of efficiency is a particular one in computer algebra, owing to the extreme swell of the intermediate values during calculations.

Our view on the domain is that research on the algorithmic manipulation of special functions is anchored between two paradigms:

- adopting linear differential equations as the right data structure for special functions,
- designing efficient algorithms in a complexity-driven way.

It aims at four kinds of algorithmic goals:

- algorithms combining functions,
- functional equations solving,
- multi-precision numerical evaluations,
- guessing heuristics.

This interacts with three domains of research:

- computer algebra, meant as the search for quasi-optimal algorithms for exact algebraic objects,
- symbolic analysis/algebraic analysis;
- experimental mathematics (combinatorics, mathematical physics, ...).

This view is made explicit in the present section.

##### 3.1.1. *Equations as a data structure*

Numerous special functions satisfy linear differential and/or recurrence equations. Under a mild technical condition, the existence of such equations induces a finiteness property that makes the main properties of the functions decidable. We thus speak of *D-finite functions*. For example, 60 % of the chapters in the handbook [18] describe D-finite functions. In addition, the class is closed under a rich set of algebraic operations. This makes linear functional equations just the right data structure to encode and manipulate special functions. The power of this representation was observed in the early 1990s [70], leading to the design of many algorithms in computer algebra. Both on the theoretical and algorithmic sides, the study of D-finite functions shares much with neighbouring mathematical domains: differential algebra, D-module theory, differential Galois theory, as well as their counterparts for recurrence equations.

##### 3.1.2. *Algorithms combining functions*

Differential/recurrence equations that define special functions can be recombined [70] to define: additions and products of special functions; compositions of special functions; integrals and sums involving special functions. Zeilberger’s fast algorithm for obtaining recurrences satisfied by parametrised binomial sums was developed in the early 1990s already [71]. It is the basis of all modern definite summation and integration algorithms. The theory was made fully rigorous and algorithmic in later works, mostly by a group in RISC (Linz, Austria) and by members of the team [59], [67], [35], [33], [34], [54]. The past ÉPI Algorithms contributed several implementations (*gfun* [62], *Mgfun* [35]).

### 3.1.3. Solving functional equations

Encoding special functions as defining linear functional equations postpones some of the difficulty of the problems to a delayed solving of equations. But at the same time, solving (for special classes of functions) is a sub-task of many algorithms on special functions, especially so when solving in terms of polynomial or rational functions. A lot of work has been done in this direction in the 1990s; more intensively since the 2000s, solving differential and recurrence equations in terms of special functions has also been investigated.

### 3.1.4. Multi-precision numerical evaluation

A major conceptual and algorithmic difference exists for numerical calculations between data structures that fit on a machine word and data structures of arbitrary length, that is, *multi-precision* arithmetic. When multi-precision floating-point numbers became available, early works on the evaluation of special functions were just promising that “most” digits in the output were correct, and performed by heuristically increasing precision during intermediate calculations, without intended rigour. The original theory has evolved in a twofold way since the 1990s: by making computable all constants hidden in asymptotic approximations, it became possible to guarantee a *prescribed* absolute precision; by employing state-of-the-art algorithms on polynomials, matrices, etc, it became possible to have evaluation algorithms in a time complexity that is linear in the output size, with a constant that is not more than a few units. On the implementation side, several original works exist, one of which (*NumGfun* [58]) is used in our DDMF.

### 3.1.5. Guessing heuristics

“Differential approximation”, or “Guessing”, is an operation to get an ODE likely to be satisfied by a given approximate series expansion of an unknown function. This has been used at least since the 1970s and is a key stone in spectacular applications in experimental mathematics [32]. All this is based on subtle algorithms for Hermite–Padé approximants [22]. Moreover, guessing can at times be complemented by proven quantitative results that turn the heuristics into an algorithm [30]. This is a promising algorithmic approach that deserves more attention than it has received so far.

### 3.1.6. Complexity-driven design of algorithms

The main concern of computer algebra has long been to prove the feasibility of a given problem, that is, to show the existence of an algorithmic solution for it. However, with the advent of faster and faster computers, complexity results have ceased to be of theoretical interest only. Nowadays, a large track of works in computer algebra is interested in developing fast algorithms, with time complexity as close as possible to linear in their output size. After most of the more pervasive objects like integers, polynomials, and matrices have been endowed with fast algorithms for the main operations on them [41], the community, including ourselves, started to turn its attention to differential and recurrence objects in the 2000s. The subject is still not as developed as in the commutative case, and a major challenge remains to understand the combinatorics behind summation and integration. On the methodological side, several paradigms occur repeatedly in fast algorithms: “divide and conquer” to balance calculations, “evaluation and interpolation” to avoid intermediate swell of data, etc. [27].

## 3.2. Trusted computer-algebra calculations

### 3.2.1. Encyclopedias

Handbooks collecting mathematical properties aim at serving as reference, therefore trusted, documents. The decision of several authors or maintainers of such knowledge bases to move from paper books [18], [20], [63] to websites and wikis<sup>0</sup> allows for a more collaborative effort in proof reading. Another step toward further confidence is to manage to generate the content of an encyclopedia by computer-algebra programs, as is the case with the Wolfram Functions Site<sup>0</sup> or DDMF<sup>0</sup>. Yet, due to the lingering doubts about computer-algebra systems, some encyclopedias propose both cross-checking by different systems and handwritten companion paper proofs of their content<sup>0</sup>. As of today, there is no encyclopedia certified with formal proofs.

<sup>0</sup>for instance <http://dlmf.nist.gov/> for special functions or <http://oeis.org/> for integer sequences

<sup>0</sup><http://functions.wolfram.com/>

<sup>0</sup><http://ddmf.msr-inria.inria.fr/1.9.1/ddmf>

### ***3.2.2. Computer algebra and symbolic logic***

Several attempts have been made in order to extend existing computer-algebra systems with symbolic manipulations of logical formulas. Yet, these works are more about extending the expressivity of computer-algebra systems than about improving the standards of correctness and semantics of the systems. Conversely, several projects have addressed the communication of a proof system with a computer-algebra system, resulting in an increased automation available in the proof system, to the price of the uncertainty of the computations performed by this oracle.

### ***3.2.3. Certifying systems for computer algebra***

More ambitious projects have tried to design a new computer-algebra system providing an environment where the user could both program efficiently and elaborate formal and machine-checked proofs of correctness, by calling a general-purpose proof assistant like the Coq system. This approach requires a huge manpower and a daunting effort in order to re-implement a complete computer-algebra system, as well as the libraries of formal mathematics required by such formal proofs.

### ***3.2.4. Semantics for computer algebra***

The move to machine-checked proofs of the mathematical correctness of the output of computer-algebra implementations demands a prior clarification about the often implicit assumptions on which the presumably correctly implemented algorithms rely. Interestingly, this preliminary work, which could be considered as independent from a formal certification project, is seldom precise or even available in the literature.

### ***3.2.5. Formal proofs for symbolic components of computer-algebra systems***

A number of authors have investigated ways to organize the communication of a chosen computer-algebra system with a chosen proof assistant in order to certify specific components of the computer-algebra systems, experimenting various combinations of systems and various formats for mathematical exchanges. Another line of research consists in the implementation and certification of computer-algebra algorithms inside the logic [66], [46], [55] or as a proof-automation strategy. Normalization algorithms are of special interest when they allow to check results possibly obtained by an external computer-algebra oracle [38]. A discussion about the systematic separation of the search for a solution and the checking of the solution is already clearly outlined in [52].

### ***3.2.6. Formal proofs for numerical components of computer-algebra systems***

Significant progress has been made in the certification of numerical applications by formal proofs. Libraries formalizing and implementing floating-point arithmetic as well as large numbers and arbitrary-precision arithmetic are available. These libraries are used to certify floating-point programs, implementations of mathematical functions and for applications like hybrid systems.

## **3.3. Machine-checked proofs of formalized mathematics**

To be checked by a machine, a proof needs to be expressed in a constrained, relatively simple formal language. Proof assistants provide facilities to write proofs in such languages. But, as merely writing, even in a formal language, does not constitute a formal proof just per se, proof assistants also provide a proof checker: a small and well-understood piece of software in charge of verifying the correctness of arbitrarily large proofs. The gap between the low-level formal language a machine can check and the sophistication of an average page of mathematics is conspicuous and unavoidable. Proof assistants try to bridge this gap by offering facilities, like notations or automation, to support convenient formalization methodologies. Indeed, many aspects, from the logical foundation to the user interface, play an important role in the feasibility of formalized mathematics inside a proof assistant.

---

<sup>0</sup><http://129.81.170.14/~vhm/Table.html>

### 3.3.1. Logical foundations and proof assistants

While many logical foundations for mathematics have been proposed, studied, and implemented, type theory is the one that has been more successfully employed to formalize mathematics, to the notable exception of the Mizar system [56], which is based on set theory. In particular, the calculus of construction (CoC) [36] and its extension with inductive types (CIC) [37], have been studied for more than 20 years and been implemented by several independent tools (like Lego, Matita, and Agda). Its reference implementation, Coq [64], has been used for several large-scale formalizations projects (formal certification of a compiler back-end; four-color theorem). Improving the type theory underlying the Coq system remains an active area of research. Other systems based on different type theories do exist and, whilst being more oriented toward software verification, have been also used to verify results of mainstream mathematics (prime-number theorem; Kepler conjecture).

### 3.3.2. Computations in formal proofs

The most distinguishing feature of CoC is that computation is promoted to the status of rigorous logical argument. Moreover, in its extension CIC, we can recognize the key ingredients of a functional programming language like inductive types, pattern matching, and recursive functions. Indeed, one can program effectively inside tools based on CIC like Coq. This possibility has paved the way to many effective formalization techniques that were essential to the most impressive formalizations made in CIC.

Another milestone in the promotion of the computations-as-proofs feature of Coq has been the integration of compilation techniques in the system to speed up evaluation. Coq can now run realistic programs in the logic, and hence easily incorporates calculations into proofs that demand heavy computational steps.

Because of their different choice for the underlying logic, other proof assistants have to simulate computations outside the formal system, and indeed fewer attempts to formalize mathematical proofs involving heavy calculations have been made in these tools. The only notable exception, which was finished in 2014, the Kepler conjecture, required a significant work to optimize the rewriting engine that simulates evaluation in Isabelle/HOL.

### 3.3.3. Large-scale computations for proofs inside the Coq system

Programs run and proved correct inside the logic are especially useful for the conception of automated decision procedures. To this end, inductive types are used as an internal language for the description of mathematical objects by their syntax, thus enabling programs to reason and compute by case analysis and recursion on symbolic expressions.

The output of complex and optimized programs external to the proof assistant can also be stamped with a formal proof of correctness when their result is easier to *check* than to *find*. In that case one can benefit from their efficiency without compromising the level of confidence on their output at the price of writing and certify a checker inside the logic. This approach, which has been successfully used in various contexts, is very relevant to the present research project.

### 3.3.4. Relevant contributions from the Mathematical Component libraries

Representing abstract algebra in a proof assistant has been studied for long. The libraries developed by the MathComp project for the proof of the Odd Order Theorem provide a rather comprehensive hierarchy of structures; however, they originally feature a large number of instances of structures that they need to organize. On the methodological side, this hierarchy is an incarnation of an original work [40] based on various mechanisms, primarily type inference, typically employed in the area of programming languages. A large amount of information that is implicit in handwritten proofs, and that must become explicit at formalization time, can be systematically recovered following this methodology.

Small-scale reflection [43] is another methodology promoted by the MathComp project. Its ultimate goal is to ease formal proofs by systematically dealing with as many bureaucratic steps as possible, by automated computation. For instance, as opposed to the style advocated by Coq's standard library, decidable predicates are systematically represented using computable boolean functions: comparison on integers is expressed as

program, and to state that  $a \leq b$  one compares the output of this program run on  $a$  and  $b$  with *true*. In many cases, for example when  $a$  and  $b$  are values, one can prove or disprove the inequality by pure computation.

The MathComp library was consistently designed after uniform principles of software engineering. These principles range from simple ones, like naming conventions, to more advanced ones, like generic programming, resulting in a robust and reusable collection of formal mathematical components. This large body of formalized mathematics covers a broad panel of algebraic theories, including of course advanced topics of finite group theory, but also linear algebra, commutative algebra, Galois theory, and representation theory. We refer the interested reader to the online documentation of these libraries [65], which represent about 150,000 lines of code and include roughly 4,000 definitions and 13,000 theorems.

Topics not addressed by these libraries and that might be relevant to the present project include real analysis and differential equations. The most advanced work of formalization on these domains is available in the HOL-Light system [48], [49], [50], although some existing developments of interest [25], [57] are also available for Coq. Another aspect of the MathComp libraries that needs improvement, owing to the size of the data we manipulate, is the connection with efficient data structures and implementations, which only starts to be explored.

### **3.3.5. User interaction with the proof assistant**

The user of a proof assistant describes the proof he wants to formalize in the system using a textual language. Depending on the peculiarities of the formal system and the applicative domain, different proof languages have been developed. Some proof assistants promote the use of a declarative language, when the Coq and Matita systems are more oriented toward a procedural style.

The development of the large, consistent body of MathComp libraries has prompted the need to design an alternative and coherent language extension for the Coq proof assistant [45], [44], enforcing the robustness of proof scripts to the numerous changes induced by code refactoring and enhancing the support for the methodology of small-scale reflection.

The development of large libraries is quite a novelty for the Coq system. In particular any long-term development process requires the iteration of many refactoring steps and very little support is provided by most proof assistants, with the notable exception of Mizar [61]. For the Coq system, this is an active area of research.



## TOCCATA Project-Team

### 3. Research Program

#### 3.1. Research Program

##### 3.1.1. Panorama of Deductive Verification

There are two main families of approaches for deductive verification. Methods in the first family build on top of mathematical proof assistants (e.g., Coq, Isabelle) in which both the model and the program are encoded; the proof that the program meets its specification is typically conducted in an interactive way using the underlying proof construction engine. Methods from the second family proceed by the design of standalone tools taking as input a program in a particular programming language (e.g., C, Java) specified with a dedicated annotation language (e.g., ACSL [49], JML [56]) and automatically producing a set of mathematical formulas (the *verification conditions*) which are typically proved using automatic provers (e.g., Z3 [70], Alt-Ergo [58], CVC4 [48]).

The first family of approaches usually offers a higher level of assurance than the second, but also demands more work to perform the proofs (because of their interactive nature) and makes them less easy to adopt by industry. Moreover, they generally do not allow to directly analyze a program written in a mainstream programming language like Java or C. The second kind of approaches has benefited in the past years from the tremendous progress made in SAT and SMT solving techniques, allowing more impact on industrial practices, but suffers from a lower level of trust: in all parts of the proof chain (the model of the input programming language, the VC generator, the back-end automatic prover), potential errors may appear, compromising the guarantee offered. Moreover, while these approaches are applied to mainstream languages, they usually support only a subset of their features.

##### 3.1.2. Overall Goals of the Toccata Project

One of our original skills is the ability to conduct proofs by using automatic provers and proof assistants at the same time, depending on the difficulty of the program, and specifically the difficulty of each particular verification condition. We thus believe that we are in a good position to propose a bridge between the two families of approaches of deductive verification presented above. Establishing this bridge is one of the goals of the Toccata project: we want to provide methods and tools for deductive program verification that can offer both a high amount of proof automation and a high guarantee of validity. Indeed, an axis of research of Toccata is the development of languages, methods and tools that are themselves formally proved correct.

In industrial applications, numerical calculations are very common (e.g. control software in transportation). Typically they involve floating-point numbers. Some of the members of Toccata have an internationally recognized expertise on deductive program verification involving floating-point computations. Our past work includes a new approach for proving behavioral properties of numerical C programs using Frama-C/Jessie [47], various examples of applications of that approach [54], the use of the Gappa solver for proving numerical algorithms [68], an approach to take architectures and compilers into account when dealing with floating-point programs [55], [66]. We also contributed to the Handbook of Floating-Point Arithmetic [65]. A representative case study is the analysis and the proof of both the method error and the rounding error of a numerical analysis program solving the one-dimension acoustic wave equation [52] [51]. Our experience led us to a conclusion that verification of numerical programs can benefit a lot from combining automatic and interactive theorem proving [53], [54], [59]. Verification of numerical programs is another main axis of Toccata.

Our scientific programme detailed below is structured into four axes:

1. Foundations and spreading of deductive program verification;
2. Reasoning on mutable memory in program verification;
3. Verification of Computer Arithmetic;
4. Spreading Formal Proofs.

Let us conclude with more general considerations about our agenda of the next four years: we want to keep on

- with general audience actions;
- industrial transfer, in particular through an extension of the perimeter of the ProofInUse joint lab.

### 3.2. Foundations and spreading of deductive program verification

Permanent researchers: S. Conchon, J.-C. Filliâtre, C. Marché, G. Melquiond, A. Paskevich

This axis covers the central theme of the team: deductive verification, from the point of view of its foundations but also our will to spread its use in software development. The general motto we want to defend is “deductive verification for the masses”. A non-exhaustive list of subjects we want to address is as follows.

- The verification of general-purpose algorithms and data structures: the challenge is to discover adequate invariants to obtain a proof, in the most automatic way as possible, in the continuation of the current VOCaL project and the various case studies presented in Axis 4 below.
- Uniform approaches to obtain correct-by-construction programs and libraries, in particular by automatic extraction of executable code (in OCaml, C, CakeML, etc.) from verified programs, and including innovative general methods like advanced ghost code, ghost monitoring, etc.
- Automated reasoning dedicated to deductive verification, so as to improve proof automation; improved combination of interactive provers and fully automated ones, proof by reflection.
- Improved feedback in case of proof failures: based on generation of counterexamples, or symbolic execution, or possibly randomized techniques à la quickcheck.
- Reduction of the trusted computing base in our toolchains: production of certificates from automatic proofs, for goal transformations (like those done by Why3), and from the generation of VCs

A significant part of the work achieved in this axis is related to the Why3 toolbox and its ecosystem, displayed on Figure 1. The boxes in red background correspond to the tools we develop in the Toccata team.

### 3.3. Reasoning on mutable memory in program verification

Permanent researchers: J.-C. Filliâtre, C. Marché, G. Melquiond, A. Paskevich

This axis concerns specifically the techniques for reasoning on programs where aliasing is the central issue. It covers the methods based on type-based alias analysis and related memory models, on specific program logics such as separation logics, and extended model-checking. It concerns the application on analysis of C or C++ codes, on Ada codes involving pointers, but also concurrent programs in general. The main topics planned are:

- The study of advanced type systems dedicated to verification, for controlling aliasing, and their use for obtaining easier-to-prove verification conditions. Modern typing system in the style of Rust, involving ownership and borrowing, will be considered.
- The design of front-ends of Why3 for the proofs of programs where aliasing cannot be fully controlled statically, via adequate memory models, aiming in particular at extraction to C; and also for concurrent programs.
- The continuation of fruitful work on concurrent parameterized systems, and its corresponding specific SMT-based model-checking.
- Concurrent programming on weak memory models, on one hand as an extension of parameterized systems above, but also in the specific context of OCaml multicore (<http://ocamlabs.io/doc/multicore.html>).
- In particular in the context of the ProofInUse joint lab, design methods for Ada, C, C++ or Java using memory models involving fine-grain analysis of pointers. Rust programs could be considered as well.

### 3.4. Verification of Computer Arithmetic

Permanent researchers: S. Boldo, C. Marché, G. Melquiond

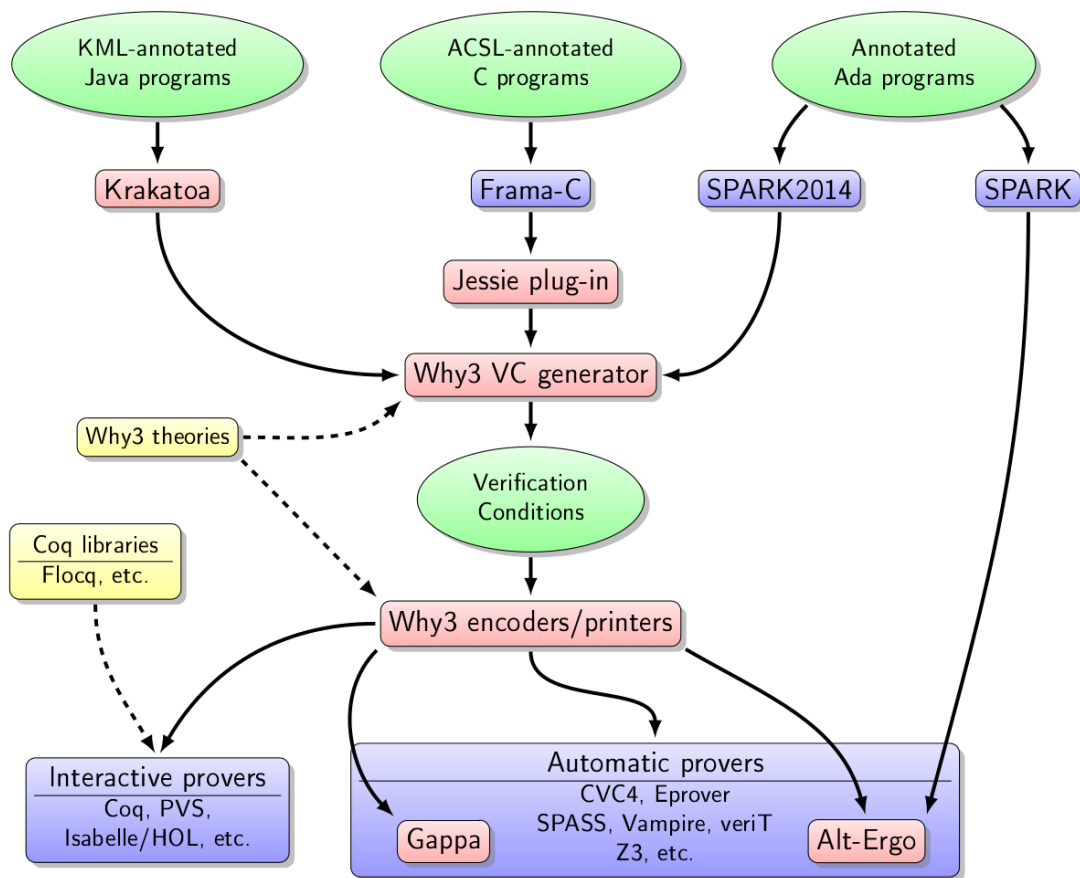


Figure 1. The Why3 ecosystem

We of course want to keep this axis which is a major originality of Toccata. The main topics of the next 4 years will be:

- Fundamental studies concerning formalization of floating-point computations, algorithms, and error analysis. Related to numerical integration, we will develop the relationships between mathematical stability and floating-point stability of numerical schemes.
- A significant effort dedicated to verification of numerical programs written in C, Ada, C++. This involves combining specifications in real numbers and computation in floating-point, and underlying automated reasoning techniques with floating-point numbers and real numbers. A new approach we have in mind concerns some variant of symbolic execution of both code and specifications involving real numbers.
- We have not yet studied embedded systems. Our approach is first to tackle numerical filters. This requires more results on fixed-point arithmetic and a careful study of overflows.
- Also a specific focus on arbitrary precision integer arithmetic, in the continuation of the ongoing PhD thesis of R. Rieu-Helft.

### **3.5. Spreading Formal Proofs**

Permanent researchers: S. Boldo, S. Conchon, J.-C. Filliâtre, C. Marché, G. Melquiond, A. Paskevich

This axis covers applications in general. The applications we currently have in mind are:

- Hybrid Systems, i.e., systems mixing discrete and continuous transitions. This theme covers many aspects such as general techniques for formally reasoning of differential equations, and extending SMT-based reasoning. The challenge is to support both abstract mathematical reasoning and concrete program execution (e.g., using floating-point representation). Hybrid systems will be a common effort with other members of the future laboratory joint with LSV of ENS Cachan.
- Applied mathematics, in the continuation of the current efforts towards verification of Finite Element Method. It has only been studied in the mathematical point of view during this period. We plan to also consider their floating-point behavior and a demanding application is that of molecular simulation exhibited in the new EMC2 project. The challenge here is both in the mathematics to be formalized, in the numerical errors that have never been studied (and that may be huge in specific cases), and in the size of the programs, which requires that our tools scale.
- Continuation of our work on analysis of shell scripts. The challenge is to be able to analyze a large number of scripts (more than 30,000 in the corpus of Debian packages installation scripts) in an automatic manner. An approach that will be considered is some form of symbolic execution.
- Explore proof tools for mathematics, in particular automated reasoning for real analysis (application: formalization of the weak Goldbach conjecture), and in number theory.
- Obtain and distribute verified OCaml libraries, as expected outcome of the VOCaL project.
- Formalization of abstract interpretation and WP calculi: in the continuation of the former project Verasco, and an ongoing project proposal joint with CEA List. The difficulty of achieving full verification of such tools will be mitigated by use of certificate techniques.

## CELESTE Project-Team

### 3. Research Program

#### 3.1. General presentation

Our objectives correspond to four major challenges of machine learning where mathematical statistics have a key role. First, any machine learning procedure depends on hyperparameters that must be chosen, and many procedures are available for any given learning problem: both are an estimator selection problem. Second, with high-dimensional and/or large data, the computational complexity of algorithms must be taken into account differently, leading to possible trade-offs between statistical accuracy and complexity, for machine learning procedures themselves as well as for estimator selection procedures. Third, real data are almost always corrupted partially, making it necessary to provide learning (and estimator selection) procedures that are robust to outliers and heavy tails, while being able to handle large datasets. Fourth, science currently faces a reproducibility crisis, making it necessary to provide statistical inference tools (p-values, confidence regions) for assessing the significance of the output of any learning algorithm (including the tuning of its hyperparameters), in a computationally efficient way.

#### 3.2. Estimator selection

An important goal of CELESTE is to build and study procedures that can deal with general estimators (especially those actually used in practice, which often rely on some optimization algorithm), such as cross-validation and Lepski's method. In order to be practical, estimator selection procedures must be fully data-driven (that is, not relying on any unknown quantity), computationally tractable (especially in the high-dimensional setting, for which specific procedures must be developed) and robust to outliers (since most real data sets include a few outliers). CELESTE aims at providing a precise theoretical analysis (for new and existing popular estimator selection procedures), that explains as well as possible their observed behaviour in practice.

#### 3.3. Relating statistical accuracy to computational complexity

When several learning algorithms are available, with increasing computational complexity and statistical performance, which one should be used, given the amount of data and the computational power available? This problem has emerged as a key question induced by the challenge of analyzing large amounts of data – the “big data” challenge. CELESTE wants to tackle the major challenge of understanding the time-accuracy trade-off, which requires providing new statistical analyses of machine learning procedures – as they are done in practice, including optimization algorithms – that are *precise enough* in order to account for differences of performance observed in practice, leading to general conclusions that can be trusted more generally. For instance, we study the performance of ensemble methods combined with subsampling, which is a common strategy for handling big data; examples include random forests and median-of-means algorithms.

#### 3.4. Robustness to outliers and heavy tails (with tractable algorithms)

The classical theory of robustness in statistics has recently received a lot of attention in the machine learning community. The reason is simple: large datasets are easily corrupted, due to – for instance – storage and transmission issues, and most learning algorithms are highly sensitive to dataset corruption. For example, the lasso can be completely misled by the presence of even a single outlier in a dataset. A major challenge in robust learning is to provide computationally tractable estimators with optimal subgaussian guarantees. A second important challenge in robust learning is to deal with datasets where every  $(x_i, y_i)$  is slightly corrupted. In large-dimensional data, every single data point  $x_i$  is likely to have several corrupted coordinates, and no estimator currently has strong theoretical guarantees for such data. A third important challenge is that of

robust estimator selection or aggregation. Even if several robust estimators can be built, the final aggregation or selection step in a user's routine is usually based on empirical means. This is not robust, and may damage the global performance of the procedure. Instead, we can consider more sophisticated types of aggregation of the base robust estimators built so far. A convenient framework to do so is called adversarial learning (also known as: prediction of individual sequences). Here, data is not assumed to be stochastic, and it could even be chosen by an adversary.

### **3.5. Statistical inference: (multiple) tests and confidence regions (including post-selection)**

CELESTE considers the problems of quantifying the uncertainty of predictions or estimations (thanks to confidence intervals) and of providing significance levels ( $p$ -values, corrected for multiplicity if needed) for each "discovery" made by a learning algorithm. This is an important practical issue when performing feature selection – one then speaks of post-selection inference – change-point detection or outlier detection, to name but a few. We tackle it in particular through a collaboration with the Parietal team (Inria Saclay) and LBBE (CNRS), with applications in neuroimaging and genomics.

## COMMANDS Project-Team

### 3. Research Program

#### 3.1. Historical aspects

The roots of deterministic optimal control are the “classical” theory of the calculus of variations, illustrated by the work of Newton, Bernoulli, Euler, and Lagrange (whose famous multipliers were introduced in [24]), with improvements due to the “Chicago school”, Bliss [16] during the first part of the 20th century, and by the notion of relaxed problem and generalized solution (Young [29]).

*Trajectory optimization* really started with the spectacular achievement done by Pontryagin’s group [28] during the fifties, by stating, for general optimal control problems, nonlocal optimality conditions generalizing those of Weierstrass. This motivated the application to many industrial problems (see the classical books by Bryson and Ho [20], Leitmann [26], Lee and Markus [25], Ioffe and Tihomirov [23]).

*Dynamic programming* was introduced and systematically studied by R. Bellman during the fifties. The HJB equation, whose solution is the value function of the (parameterized) optimal control problem, is a variant of the classical Hamilton-Jacobi equation of mechanics for the case of dynamics parameterized by a control variable. It may be viewed as a differential form of the dynamic programming principle. This nonlinear first-order PDE appears to be well-posed in the framework of *viscosity solutions* introduced by Crandall and Lions [21]. The theoretical contributions in this direction did not cease growing, see the books by Barles [14] and Bardi and Capuzzo-Dolcetta [13].

#### 3.2. Trajectory optimization

The so-called *direct methods* consist in an optimization of the trajectory, after having discretized time, by a nonlinear programming solver that possibly takes into account the dynamic structure. So the two main problems are the choice of the discretization and the nonlinear programming algorithm. A third problem is the possibility of refinement of the discretization once after solving on a coarser grid.

In the *full discretization approach*, general Runge-Kutta schemes with different values of control for each inner step are used. This allows to obtain and control high orders of precision, see Hager [22], Bonnans [17]. In the *indirect* approach, the control is eliminated thanks to Pontryagin’s maximum principle. One has then to solve the two-points boundary value problem (with differential variables state and costate) by a single or multiple shooting method. The questions are here the choice of a discretization scheme for the integration of the boundary value problem, of a (possibly globalized) Newton type algorithm for solving the resulting finite dimensional problem in  $\mathbb{R}^n$  ( $n$  is the number of state variables), and a methodology for finding an initial point.

#### 3.3. Hamilton-Jacobi-Bellman approach

This approach consists in calculating the value function associated with the optimal control problem, and then synthesizing the feedback control and the optimal trajectory using Pontryagin’s principle. The method has the great particular advantage of reaching directly the global optimum, which can be very interesting when the problem is not convex.

*Optimal stochastic control problems* occur when the dynamical system is uncertain. A decision typically has to be taken at each time, while realizations of future events are unknown (but some information is given on their distribution of probabilities). In particular, problems of economic nature deal with large uncertainties (on prices, production and demand). Specific examples are the portfolio selection problems in a market with risky and non-risky assets, super-replication with uncertain volatility, management of power resources (dams, gas). Air traffic control is another example of such problems.

For solving stochastic control problems, we studied the so-called Generalized Finite Differences (GFD), that allow to choose at any node, the stencil approximating the diffusion matrix up to a certain threshold [19]. Determining the stencil and the associated coefficients boils down to a quadratic program to be solved at each point of the grid, and for each control. This is definitely expensive, with the exception of special structures where the coefficients can be computed at low cost. For two dimensional systems, we designed a (very) fast algorithm for computing the coefficients of the GFD scheme, based on the Stern-Brocot tree [18].



## **DEFI Project-Team**

### **3. Research Program**

#### **3.1. Research Program**

The research activity of our team is dedicated to the design, analysis and implementation of efficient numerical methods to solve inverse and shape/topological optimization problems, eventually including system uncertainties, in connection with wave imaging, structural design, non-destructive testing and medical imaging modalities. We are particularly interested in the development of fast methods that are suited for real-time applications and/or large scale problems. These goals require to work on both the physical and the mathematical models involved and indeed a solid expertise in related numerical algorithms. A part of the research activity is also devoted to take into account system uncertainties in the solving of inverse/optimization problems. At the interface of physics, mathematics, and computer science, Uncertainty Quantification (UQ) focuses on the development of frameworks and methods to characterize uncertainties in predictive computations. Uncertainties and errors arise at different stages of the numerical simulation. First, errors are introduced due to the physical simplifications in the mathematical modeling of the system investigated; other errors come from the numerical resolution of the mathematical model, due in particular to finite discretization and computations with finite accuracy and tolerance; finally, errors are due a limited knowledge of input quantities (parameters) appearing in the definition of the numerical model being solved.

This section intends to give a general overview of our research interests and themes. We choose to present them through the specific academic example of inverse scattering problems (from inhomogeneities), which is representative of foreseen developments on both inversion and (topological) optimization methods. The practical problem would be to identify an inclusion from measurements of diffracted waves that result from the interaction of the sought inclusion with some (incident) waves sent into the probed medium. Typical applications include biomedical imaging where using micro-waves one would like to probe the presence of pathological cells, or imaging of urban infrastructures where using ground penetrating radars (GPR) one is interested in finding the location of buried facilities such as pipelines or waste deposits. This kind of applications requires in particular fast and reliable algorithms.

By “imaging” we refer to the inverse problem where the concern is only the location and the shape of the inclusion, while “identification” may also indicate getting informations on the inclusion physical parameters.

Both problems (imaging and identification) are non linear and ill-posed (lack of stability with respect to measurements errors if some careful constrains are not added). Moreover, the unique determination of the geometry or the coefficients is not guaranteed in general if sufficient measurements are not available. As an example, in the case of anisotropic inclusions, one can show that an appropriate set of data uniquely determine the geometry but not the material properties.

These theoretical considerations (uniqueness, stability) are not only important in understanding the mathematical properties of the inverse problem, but also guide the choice of appropriate numerical strategies (which information can be stably reconstructed) and also the design of appropriate regularization techniques. Moreover, uniqueness proofs are in general constructive proofs, i.e. they implicitly contain a numerical algorithm to solve the inverse problem, hence their importance for practical applications. The sampling methods introduced below are one example of such algorithms.

A large part of our research activity is dedicated to numerical methods applied to the first type of inverse problems, where only the geometrical information is sought. In its general setting the inverse problem is very challenging and no method can provide universally satisfying solution (respecting the balance cost-precision-stability). This is why in the majority of the practically employed algorithms, some simplification of the underlying mathematical model is used, according to the specific configuration of the imaging experiment. The most popular ones are geometric optics (the Kirchhoff approximation) for high frequencies and weak scattering (the Born approximation) for small contrasts or small obstacles. They actually give full satisfaction

for a wide range of applications as attested by the large success of existing imaging devices (radar, sonar, ultrasound, X-ray tomography, etc.), that rely on one of these approximations.

In most cases, the used simplification result in a linearization of the inverse problem and therefore is usually valid only if the latter is weakly non-linear. The development of simplified models and the improvement of their efficiency is still a very active research area. With that perspective, we are particularly interested in deriving and studying higher order asymptotic models associated with small geometrical parameters such as: small obstacles, thin coatings, wires, periodic media, .... Higher order models usually introduce some non linearity in the inverse problem, but are in principle easier to handle from the numerical point of view than in the case of the exact model.

A larger part of our research activity is dedicated to algorithms that avoid the use of such approximations and that are efficient where classical approaches fail: i.e. roughly speaking when the non linearity of the inverse problem is sufficiently strong. This type of configuration is motivated by the applications mentioned below, and occurs as soon as the geometry of the unknown media generates non negligible multiple scattering effects (multiply-connected and closely spaces obstacles) or when the used frequency is in the so-called resonant region (wave-length comparable to the size of the sought medium). It is therefore much more difficult to deal with and requires new approaches. Our ideas to tackle this problem is mainly motivated and inspired by recent advances in shape and topological optimization methods and in so-called sampling methods.

Sampling methods are fast imaging solvers adapted to multi-static data (multiple receiver-transmitter pairs) at a fixed frequency. Even if they do not use any linearization the forward model, they rely on computing the solutions to a set of linear problems of small size, that can be performed in a completely parallel procedure. Our team has already a solid expertise in these methods applied to electromagnetic 3-D problems. The success of such approaches was their ability to provide a relatively quick algorithm for solving 3-D problems without any need for a priori knowledge on the physical parameters of the targets. These algorithms solve only the imaging problem, in the sense that only the geometrical information is provided.

Despite the large efforts already spent in the development of this type of methods, either from the algorithmic point of view or the theoretical one, numerous questions are still open. These attractive new algorithms also suffer from the lack of experimental validations, due to their relatively recent introduction. We also would like to invest on this side by developing collaborations with engineering research groups that have experimental facilities. From the practical point of view, the most potential limitation of sampling methods would be the need of a large amount of data to achieve a reasonable accuracy. On the other hand, optimization methods do not suffer from this constrain but they require good initial guess to ensure convergence and reduce the number of iterations. Therefore it seems natural to try to combine the two class of methods in order to calibrate the balance between cost and precision.

Among various shape optimization methods, the Level Set method seems to be particularly suited for such a coupling. First, because it shares similar mechanism as sampling methods: the geometry is captured as a level set of an “indicator function” computed on a cartesian grid. Second, because the two methods do not require any a priori knowledge on the topology of the sought geometry. Beyond the choice of a particular method, the main question would be to define in which way the coupling can be achieved. Obvious strategies consist in using one method to pre-process (initialization) or post-process (find the level set) the other. But one can also think of more elaborate ones, where for instance a sampling method can be used to optimize the choice of the incident wave at each iteration step. The latter point is closely related to the design of so called “focusing incident waves” (which are for instance the basis of applications of the time-reversal principle). In the frequency regime, these incident waves can be constructed from the eigenvalue decomposition of the data operator used by sampling methods. The theoretical and numerical investigations of these aspects are still not completely understood for electromagnetic or elastodynamic problems.

Other topological optimization methods, like the homogenization method or the topological gradient method, can also be used, each one provides particular advantages in specific configurations. It is evident that the development of these methods is very suited to inverse problems and provide substantial advantage compared to classical shape optimization methods based on boundary variation. Their applications to inverse problems has not been fully investigated. The efficiency of these optimization methods can also be increased for adequate

asymptotic configurations. For instance small amplitude homogenization method can be used as an efficient relaxation method for the inverse problem in the presence of small contrasts. On the other hand, the topological gradient method has shown to perform well in localizing small inclusions with only one iteration.

A broader perspective would be the extension of the above mentioned techniques to time-dependent cases. Taking into account data in time domain is important for many practical applications, such as imaging in cluttered media, the design of absorbing coatings or also crash worthiness in the case of structural design.

For the identification problem, one would like to also have information on the physical properties of the targets. Of course optimization methods is a tool of choice for these problems. However, in some applications only a qualitative information is needed and obtaining it in a cheaper way can be performed using asymptotic theories combined with sampling methods. We also refer here to the use of so called transmission eigenvalues as qualitative indicators for non destructive testing of dielectrics.

We are also interested in parameter identification problems arising in diffusion-type problems. Our research here is mostly motivated by applications to the imaging of biological tissues with the technique of Diffusion Magnetic Resonance Imaging (DMRI). Roughly speaking DMRI gives a measure of the average distance travelled by water molecules in a certain medium and can give useful information on cellular structure and structural change when the medium is biological tissue. In particular, we would like to infer from DMRI measurements changes in the cellular volume fraction occurring upon various physiological or pathological conditions as well as the average cell size in the case of tumor imaging. The main challenges here are 1) correctly model measured signals using diffusive-type time-dependent PDEs 2) numerically handle the complexity of the tissues 3) use the first two to identify physically relevant parameters from measurements. For the last point we are particularly interested in constructing reduced models of the multiple-compartment Bloch-Torrey partial differential equation using homogenization methods.

The Team devotes a large effort focused on the formulation, implementation and validation of numerical methods for using scientific computing to drive experiments and available data (coming from models, simulation and experiments) by taking into account the system uncertainty. The team is also invested in exploiting the intimate relationship between optimisation and UQ to make Optimisation Under Uncertainty (OUU) tractable. A part of these activities is declined to the simulation of high-fidelity models for fluids, in three main fields, aerospace, energy and environment.

The Team is working on developing original UQ representations and algorithms to deal with complex and large scale models, having high dimensional input parameters with complex influences. We are organizing our core research activities along different methodological UQ developments related to the challenges discussed above. Obviously, some efforts are shared by different initiatives or projects, and some of them include the continuous improvement of the non-intrusive methods constituting our software libraries. These actions are not detailed in the following, to focus the presentation on more innovative aspects, but we mentioned nonetheless the continuous developments and incorporation in our libraries of advanced sparse grid methods, sparsity promoting strategies and low rank methods.

An effort is dedicated to the efficient construction of surrogate models that are central in both forward and backward UQ problems, aiming at large-scale simulations relevant to engineering applications, with high dimensional input parameters.

Sensitivity analyses and other forward UQ problems (e.g., estimation of failure probabilities, rare events, . . .) depends on the input uncertainty model. Most often, for convenience or because of the lack of data, the independence of the uncertain inputs is assumed. In the Team, we are investigating approaches dedicated to a) the construction of uncertainty models that integrate the available information and expert knowledge(s) in a consistent and objective fashion. To this end, several mathematical frameworks are already available, e.g the maximum entropy principle, likelihood maximization and moment matching methods, but their application to real engineering problems remains scarce and their systematic use raises multiple challenges, both to construct the uncertainty model and to solve the related UQ problems (forward and backward). Because of the importance of the available data and expertise to build the model, the contributions of the Team in these areas depend on the needs and demands of end-users and industrial partners.

To mitigate computational complexity, the Team is exploring multi-fidelity approaches in the context of expensive simulations. We combine predictions of models with different levels of discretizations and physical simplifications to construct, at a controlled cost, reliable surrogate models of simulation outputs or directly objective functions and possibly constraints, to enable the resolution of robust optimization and stochastic inverse problems. Again, one difficulty to be addressed by the Team is the design of the computer experiments to obtain the best multi-fidelity model at the lowest cost (of for a prescribed computational budgets), with respect to the end use of the model. The last point is particularly challenging as it calls for accuracy for output values that are usually unknown a priori but must be estimated as the model construction proceeds.

## DISCO Project-Team

### 3. Research Program

#### 3.1. Analysis of interconnected systems

The major questions considered are those of the characterization of the stability (also including the problems of sensitivity compared to the variations of the parameters) and the determination of stabilizing controllers of interconnected dynamic systems. In many situations, the dynamics of the interconnections can be naturally modelled by systems with delays (constant, distributed or time-varying delays) possibly of fractional order. In other cases, partial differential equations (PDE) models can be better represented or approximated by using systems with delays. Our expertise on this subject, on both time and frequency domain methods, allows us to challenge difficult problems (e.g. systems with an infinite number of unstable poles).

- Robust stability of linear systems

Within an interconnection context, lots of phenomena are modelled directly or after an approximation by delay systems. These systems may have constant delays, time-varying delays, distributed delays ...

For various infinite-dimensional systems, particularly delay and fractional systems, input-output and time-domain methods are jointly developed in the team to characterize stability. This research is developed at four levels: analytic approaches ( $H_\infty$ -stability, BIBO-stability, robust stability, robustness metrics) [1], [2], [6], [7], symbolic computation approaches (SOS methods are used for determining easy-to-check conditions which guarantee that the poles of a given linear system are not in the closed right half-plane, certified CAD techniques), numerical approaches (root-loci, continuation methods) and by means of softwares developed in the team [6], [7].

- Robustness/fragility of biological systems

Deterministic biological models describing, for instance, species interactions, are frequently composed of equations with important disturbances and poorly known parameters. To evaluate the impact of the uncertainties, we use the techniques of designing of global strict Lyapunov functions or functional developed in the team.

However, for other biological systems, the notion of robustness may be different and this question is still in its infancy (see, e.g. [57]). Unlike engineering problems where a major issue is to maintain stability in the presence of disturbances, a main issue here is to maintain the system response in the presence of disturbances. For instance, a biological network is required to keep its functioning in case of a failure of one of the nodes in the network. The team, which has a strong expertise in robustness for engineering problems, aims at contributing at the development of new robustness metrics in this biological context.

#### 3.2. Stabilization of interconnected systems

- Linear systems: Analytic and algebraic approaches are considered for infinite-dimensional linear systems studied within the input-output framework.

In the recent years, the Youla-Kučera parametrization (which gives the set of all stabilizing controllers of a system in terms of its coprime factorizations) has been the cornerstone of the success of the  $H_\infty$ -control since this parametrization allows one to rewrite the problem of finding the optimal stabilizing controllers for a certain norm such as  $H_\infty$  or  $H_2$  as affine, and thus, convex problem.

A central issue studied in the team is the computation of such factorizations for a given infinite-dimensional linear system as well as establishing the links between stabilizability of a system for a certain norm and the existence of coprime factorizations for this system. These questions are fundamental for robust stabilization problems [1], [2].

We also consider simultaneous stabilization since it plays an important role in the study of reliable stabilization, i.e. in the design of controllers which stabilize a finite family of plants describing a system during normal operating conditions and various failed modes (e.g. loss of sensors or actuators, changes in operating points). Moreover, we investigate strongly stabilizable systems, namely systems which can be stabilized by stable controllers, since they have a good ability to track reference inputs and, in practice, engineers are reluctant to use unstable controllers especially when the system is stable.

- Nonlinear systems

In any physical systems a feedback control law has to account for limitation stemming from safety, physical or technological constraints. Therefore, any realistic control system analysis and design has to account for these limitations appearing mainly from sensors and actuators nonlinearities and from the regions of safe operation in the state space. This motivates the study of linear systems with more realistic, thus complex, models of actuators. These constraints appear as nonlinearities as saturation and quantization in the inputs of the system [10].

The project aims at developing robust stabilization theory and methods for important classes of nonlinear systems that ensure good controller performance under uncertainty and time delays. The main techniques include techniques called backstepping and forwarding, constructions of strict Lyapunov functions through so-called "strictification" approaches [4] and construction of Lyapunov-Krasovskii functionals [5], [6], [7] or Lyapunov functionals for PDE systems [9].

### 3.3. Synthesis of reduced complexity controllers

- PID controllers

Even though the synthesis of control laws of a given complexity is not a new problem, it is still open, even for finite-dimensional linear systems. Our purpose is to search for good families of "simple" (e.g. low order) controllers for infinite-dimensional dynamical systems. Within our approach, PID candidates are first considered in the team [2], [60].

For interconnected systems appearing in teleoperation applications, such as the steer-by-wire, Proportional-Derivative laws are simple control strategies allowing to reproduce the efforts in both ends of the teleoperation system. However, due to delays introduced in the communication channels these strategies may result in loss of closed loop stability or in performance degradation when compared to the system with a mechanical link (no communication channel). In this context we search for non-linear proportional and derivative gains to improve performance. This is assessed in terms of reduction of overshoot and guaranteed convergence rates.

- Delayed feedback

Control systems often operate in the presence of delays, primarily due to the time it takes to acquire the information needed for decision-making, to create control decisions and to execute these decisions. Commonly, such a time delay induces desynchronizing and/or destabilizing effects on the dynamics. However, some recent studies have emphasized that the delay may have a stabilizing effect in the control design. In particular, the closed-loop stability may be guaranteed precisely by the existence of the delay. The interest of considering such control laws lies in the simplicity of the controller as well as in its easy practical implementation. It is intended by the team members to provide a unified approach for the design of such stabilizing control laws for finite and infinite dimensional plants [3], [8].

- Finite Time and Interval Observers for nonlinear systems

We aim to develop techniques of construction of output feedbacks relying on the design of observers. The objectives pertain to the design of robust control laws which converge in finite time, the construction of intervals observers which ensure that the solutions belong to guaranteed intervals, continuous/discrete observers for systems with discrete measurements and observers for systems with switches.

Finally, the development of algorithms based on both symbolic computation and numerical methods, and their implementations in dedicated Scilab/Matlab/Maple toolboxes are important issues in the project.

**GAMMA Project-Team (section vide)**



## POEMS Project-Team

### 3. Research Program

#### 3.1. Expertises

The activity of the team is oriented towards the design, the analysis and the numerical approximation of mathematical models for all types of problems involving wave propagation phenomena, in mechanics, physics and engineering sciences. Let us briefly describe our core business and current expertise, in order to clarify the new challenges that we want to address in the short and long terms.

Typically, our works are based on *boundary value problems* established by physicists to model the propagation of waves in various situations. The basic ingredient is a partial differential equation of the hyperbolic type, whose prototype is the scalar wave equation, or the Helmholtz equation if time-periodic solutions are considered. More generally, we systematically consider both the transient problem, in the time domain, and the time-harmonic problem, in the frequency domain. Let us mention that, even if different waves share a lot of common properties, the transition from the scalar acoustic equation to the vectorial electromagnetism and elastodynamics systems raises a lot of mathematical and numerical difficulties, and requires a specific expertise.

A notable particularity of the problems that we consider is that they are generally set in *unbounded domains*: for instance, for radar applications, it is necessary to simulate the interaction of the electromagnetic waves with the airplane only, without any complex environment perturbing the wave phenomena. This raises an intense research activity, both from a theoretical and a numerical point of view. There exist several approaches which all consist in rewriting the problem (or an approximation of it) in a bounded domain, the new formulation being well-suited for classical mathematical and numerical techniques.

One class of methods consists in applying an appropriate condition on some boundary enclosing the zone of interest. In the frequency domain, one can use a non-local transparent condition, which can be expressed by a convolution with a Green function like in integral equation techniques, or by a modal decomposition when a separation of variables is applicable. But for explicit schemes in the time domain, local radiation conditions at a finite distance are generally preferred (constructed as local approximations at various orders of the exact non-local condition). A second class of methods consists in surrounding the computational domain by so called *Perfectly Matched absorbing Layers* (PML), which are very popular because they are easy to implement. POEMS members have provided several contributions to these two classes of methods for more than twenty-five years. Among them, one can mention the understanding of the instability of PMLs in anisotropic media and in dispersive media, the derivation of transparent boundary conditions in periodic media or the improvement of Fast Multipole techniques for elastodynamic integral equations.

In addition to more classical domains of applied mathematics that we are led to use (variational analysis and functional analysis, interpolation and approximation theory, linear algebra of large systems, etc...), we have acquired a deep expertise in *spectral theory*. Indeed, the analysis of wave phenomena is intimately linked to the study of some associated spectral problems. Acoustic resonance frequencies of a cavity correspond to the eigenvalues of a selfadjoint Laplacian operator, modal solutions in a waveguide correspond to a spectral problem set in the cross section. In these two examples, if the cavity or the cross-section is unbounded, a part of the spectrum is a continuum. Again, POEMS has produced several contributions in this field. In particular, a large number of significant results have been obtained for the existence or non-existence of guided modes in open waveguides and of trapped modes in infinite domains.

To end this far from exhaustive presentation of our main expertise domains, let us mention the *asymptotic techniques* with respect to some small scale appearing in the model: it can be the wavelength compared to the size of the scatterer, or on the contrary, the scale of the scatterer compared to the wavelength, it can be the scale of some microstructure in a composite material or the width of a thin layer or a thin tube. In each case, the objective, in order to avoid the use of costly meshes, is to derive effective simplified models. Our

specificity here is that we can combine skills in physics, mathematics and numerics: in particular, we take care of the mathematical properties of the effective model, which are used to ensure the robustness of the numerical method, and also to derive error estimates with respect to the small parameter. There has been a lot of contributions of POEMS to this topic, going from the modeling of electromagnetic coatings to the justification of models for piezoelectric sensors. Let us mention that effective models for small scatterers and thin coatings have been used to improve imaging techniques that we are developing (topological gradient, time reversal or sampling techniques).

### 3.2. New domains

In order to consider more and more challenging problems (involving non-deterministic, large-scale and more realistic models), we decided recently to enlarge our domain of expertise in three directions.

Firstly, we want to reinforce our activity on *efficient solvers for large-scale wave propagation problems*. Since its inception, POEMS has frequently contributed to the development and the analysis of numerical methods that permit the fast solution of large-scale problems, such as high-order finite element methods, boundary elements methods and domain decomposition methods. Nevertheless, implementing these methods in parallel programming environments and dealing with large-scale benchmarks have generally not been done by the team. We want to continue our activities on these methods and, in a more comprehensive approach, we will incorporate modern algebraic strategies and high-performance computing (HPC) aspects in our methodology. In collaboration with academic or industrial partners, we would like to address industrial-scale benchmarks to assess the performance of our approaches. We believe that taking all these aspects into consideration will allow us to design more efficient wave-specific computational tools for large-scale simulations.

Secondly, up to now, *probabilistic methods* were outside the expertise of POEMS team, restricting us to deterministic approaches for wave propagation problems. We however firmly believe in the importance and usefulness of addressing uncertainty and randomness inherent to many propagation phenomena. Randomness may occur in the description of complex propagation media (for example in the modeling of ultrasound waves in concrete for the simulation of non-destructive testing experiments) or of data uncertainties. To quantify the effect of such uncertainties on the design, behavior, performance or reliability of many systems is then a natural goal in diverse fields of application.

Thirdly and lastly, we wish to develop and strengthen collaborations allowing a *closer interaction between our mathematical, modeling and computing activities and physical experiments*, where the latter may either provide reality checks on existing models or strongly affect the choice of modeling assumptions. Within our typical domain of activities, we can mention three areas for which such considerations are highly relevant. One is musical acoustics, where POEMS has made several well-recognized contributions dealing with the simulation of musical instruments. Another area is inverse problems, whose very purpose is to extract useful information from actual measurements with the help of (propagation) models. This is a core of our partnership with CEA on ultrasonic Non Destructive Testing. A third area is the modelling of effective (acoustic or electromagnetic) metamaterials, where predictions based on homogenized models have to be confirmed by experiments.

## RANDOPT Project-Team

### 3. Research Program

#### 3.1. Introduction

The lines of research we intend to pursue is organized along four axis namely developing novel theoretical framework, developing novel algorithms, setting novel standards in scientific experimentation and benchmarking and applications.

#### 3.2. Developing Novel Theoretical Frameworks for Analyzing and Designing Adaptive Stochastic Algorithms

Stochastic black-box algorithms typically optimize **non-convex, non-smooth functions**. This is possible because the algorithms rely on weak mathematical properties of the underlying functions: the algorithms do not use the derivatives—hence the function does not need to be differentiable—and, additionally, often do not use the exact function value but instead how the objective function ranks candidate solutions (such methods are sometimes called function-value-free). (To illustrate a comparison-based update, consider an algorithm that samples  $\lambda$  (with  $\lambda$  an even integer) candidate solutions from a multivariate normal distribution. Let  $x_1, \dots, x_\lambda$  in  $\mathbb{R}^n$  denote those  $\lambda$  candidate solutions at a given iteration. The solutions are evaluated on the function  $f$  to be minimized and ranked from the best to the worse:

$$f(x_{1:\lambda}) \leq \dots \leq f(x_{\lambda:\lambda}) .$$

In the previous equation  $i:\lambda$  denotes the index of the sampled solution associated to the  $i$ -th best solution. The new mean of the Gaussian vector from which new solutions will be sampled at the next iteration can be updated as

$$m \leftarrow \frac{1}{\lambda} \sum_{i=1}^{\lambda/2} x_{i:\lambda} .$$

The previous update moves the mean towards the  $\lambda/2$  best solutions. Yet the update is only based on the ranking of the candidate solutions such that the update is the same if  $f$  is optimized or  $g \circ f$  where  $g : \text{Im}(f) \rightarrow \mathbb{R}$  is strictly increasing. Consequently, such algorithms are invariant with respect to strictly increasing transformations of the objective function. This entails that they are robust and their performances generalize well.)

Additionally, adaptive stochastic optimization algorithms typically have a **complex state space** which encodes the parameters of a probability distribution (e.g. mean and covariance matrix of a Gaussian vector) and other state vectors. This state-space is a **manifold**. While the algorithms are Markov chains, the complexity of the state-space makes that **standard Markov chain theory tools do not directly apply**. The same holds with tools stemming from stochastic approximation theory or Ordinary Differential Equation (ODE) theory where it is usually assumed that the underlying ODE (obtained by proper averaging and limit for learning rate to zero) has its critical points inside the search space. In contrast, in the cases we are interested in, the **critical points of the ODEs are at the boundary of the domain**.

Last, since we aim at developing theory that on the one hand allows to analyze the main properties of state-of-the-art methods and on the other hand is useful for algorithm design, we need to be careful not to use simplifications that would allow a proof to be done but would not capture the important properties of the algorithms. With that respect one tricky point is to develop **theory that accounts for invariance properties**.

To face those specific challenges, we need to develop novel theoretical frameworks exploiting invariance properties and accounting for peculiar state-spaces. Those frameworks should allow researchers to analyze one of the core properties of adaptive stochastic methods, namely **linear convergence** on the widest possible class of functions.

We are planning to approach the question of linear convergence from three different complementary angles, using three different frameworks:

- the Markov chain framework where the convergence derives from the analysis of the stability of a normalized Markov chain existing on scaling-invariant functions for translation and scale-invariant algorithms [18]. This framework allows for a fine analysis where the exact convergence rate can be given as an implicit function of the invariant measure of the normalized Markov chain. Yet it requires the objective function to be scaling-invariant. The stability analysis can be particularly tricky as the Markov chain that needs to be studied writes as  $\Phi_{t+1} = F(\Phi_t, W_{t+1})$  where  $\{W_t : t > 0\}$  are independent identically distributed and  $F$  is typically discontinuous because the algorithms studied are comparison-based. This implies that practical tools for analyzing a standard property like irreducibility, that rely on investigating the stability of underlying deterministic control models [34], cannot be used. Additionally, the construction of a drift to prove ergodicity is particularly delicate when the state space includes a (normalized) covariance matrix as it is the case for analyzing the CMA-ES algorithm.
- The stochastic approximation or ODE framework. Those are standard techniques to prove the convergence of stochastic algorithms when an algorithm can be expressed as a stochastic approximation of the solution of a mean field ODE [20], [21], [32]. What is specific and induces difficulties for the algorithms we aim at analyzing is the **non-standard state-space** since the ODE variables correspond to the state-variables of the algorithm (e.g.  $\mathbb{R}^n \times \mathbb{R}_{>0}$  for step-size adaptive algorithms,  $\mathbb{R}^n \times \mathbb{R}_{>0} \times S_{++}^n$  where  $S_{++}^n$  denotes the set of positive definite matrices if a covariance matrix is additionally adapted). Consequently, the ODE can have many critical points at the boundary of its definition domain (e.g. all points corresponding to  $\sigma_t = 0$  are critical points of the ODE) which is not typical. Also we aim at proving **linear convergence**, for that it is crucial that the learning rate does not decrease to zero which is non-standard in ODE method.
- The direct framework where we construct a global Lyapunov function for the original algorithm from which we deduce bounds on the hitting time to reach an  $\epsilon$ -ball of the optimum. For this framework as for the ODE framework, we expect that the class of functions where we can prove linear convergence are composite of  $g \circ f$  where  $f$  is differentiable and  $g : \text{Im}(f) \rightarrow \mathbb{R}$  is strictly increasing and that we can show convergence to a local minimum.

We expect those frameworks to be complementary in the sense that the assumptions required are different. Typically, the ODE framework should allow for proofs under the assumptions that learning rates are small enough while it is not needed for the Markov chain framework. Hence this latter framework captures better the real dynamics of the algorithm, yet under the assumption of scaling-invariance of the objective functions. Also, we expect some overlap in terms of function classes that can be studied by the different frameworks (typically convex-quadratic functions should be encompassed in the three frameworks). By studying the different frameworks in parallel, we expect to gain synergies and possibly understand what is the most promising approach for solving the holy grail question of the linear convergence of CMA-ES. We foresee for instance that similar approaches like the use of Foster-Lyapunov drift conditions are needed in all the frameworks and that intuition can be gained on how to establish the conditions from one framework to another one.

### 3.3. Algorithmic developments

We are planning on developing algorithms in the subdomains with strong practical demand for better methods of constrained, multiobjective, large-scale and expensive optimization.

Many of the algorithm developments, we propose, rely on the CMA-ES method. While this seems to restrict our possibilities, we want to emphasize that CMA-ES became a *family of methods* over the years that nowadays include various techniques and developments from the literature to handle non-standard optimization problems (noisy, large-scale, ...). The core idea of all CMA-ES variants—namely the mechanism to adapt a Gaussian distribution—has furthermore been shown to derive naturally from first principles with only minimal assumptions in the context of derivative-free black-box stochastic optimization [35], [25]. This is a strong justification for relying on the CMA-ES premises while new developments naturally include new techniques typically borrowed from other fields. While CMA-ES is now a full family of methods, for visibility reasons, we continue to refer often to “the CMA-ES algorithm”.

### 3.3.1. *Constrained optimization*

Many (real-world) optimization problems have constraints related to technical feasibility, cost, etc. Constraints are classically handled in the black-box setting either via rejection of solutions violating the constraints—which can be quite costly and even lead to quasi-infinite loops—or by penalization with respect to the distance to the feasible domain (if this information can be extracted) or with respect to the constraint function value [22]. However, the penalization coefficient is a sensitive parameter that needs to be adapted in order to achieve a robust and general method [23]. Yet, **the question of how to handle properly constraints is largely unsolved**. The latest constraints handling for CMA-ES is an ad-hoc technique driven by many heuristics [23]. Also, it is particularly only recently that it was pointed out that **linear convergence properties should be preserved** when addressing constraint problems [16].

Promising approaches though, rely on using augmented Lagrangians [16], [17]. The augmented Lagrangian, here, is the objective function optimized by the algorithm. Yet, it depends on coefficients that are adapted online. The adaptation of those coefficients is the difficult part: the algorithm should be stable and the adaptation efficient. We believe that the theoretical frameworks developed (particularly the Markov chain framework) will be useful to understand how to design the adaptation mechanisms. Additionally, the question of invariance will also be at the core of the design of the methods: augmented Lagrangian approaches break the invariance to monotonic transformation of the objective functions, yet understanding the maximal invariance that can be achieved seems to be an important step towards understanding what adaptation rules should satisfy.

### 3.3.2. *Large-scale Optimization*

In the large-scale setting, we are interested to optimize problems with the order of  $10^3$  to  $10^4$  variables. For one to two orders of magnitude more variables, we will talk about a “very large-scale” setting.

In this context, algorithms with a quadratic scaling (internal and in terms of number of function evaluations needed to optimize the problem) cannot be afforded. In CMA-ES-type algorithms, we typically need to restrict the model of the covariance matrix to have only a linear number of parameters to learn such that the algorithms scale linearly in terms of internal complexity, memory and number of function evaluations to solve the problem. The main challenge is thus to have rich enough models for which we can efficiently design proper adaptation mechanisms. Some first large-scale variants of CMA-ES have been derived. They include the online adaptation of the complexity of the model [15], [14]. Yet so far they fail to optimize functions whose Hessian matrix has some small eigenvalues (say around  $10^{-4}$ ) some eigenvalues equal to 1 and some very large eigenvalue (say around  $10^4$ ), that is functions whose level sets have short and long axis.

Another direction, we want to pursue, is exploring the use of large-scale variants of CMA-ES to solve reinforcement learning problems [36].

Last, we are interested to investigate the very-large-scale setting. One approach consists in doing optimization in subspaces. This entails the efficient identification of relevant spaces and the restriction of the optimization to those subspaces.

### 3.3.3. *Multiobjective Optimization*

Multiobjective optimization, i.e., the simultaneous optimization of multiple objective functions, differs from single-objective optimization in particular in its optimization goal. Instead of aiming at converging to the

solution with the best possible function value, in multiobjective optimization, a set of solutions<sup>0</sup> is sought. This set, called Pareto-set, contains all trade-off solutions in the sense of Pareto-optimality—no solution exists that is better in *all* objectives than a Pareto-optimal one. Because converging towards a set differs from converging to a single solution, it is no surprise that we might lose many good convergence properties if we directly apply search operators from single-objective methods. However, this is what has typically been done so far in the literature. Indeed, most of the research in stochastic algorithms for multiobjective optimization focused instead on the so called selection part, that decides which solutions should be kept during the optimization—a question that can be considered as solved for many years in the case of single-objective stochastic adaptive methods.

We therefore aim at rethinking search operators and adaptive mechanisms to improve existing methods. We expect that we can obtain orders of magnitude better convergence rates for certain problem types if we choose the right search operators. We typically see two angles of attack: On the one hand, we will study methods based on scalarizing functions that transform the multiobjective problem into a set of single-objective problems. Those single-objective problems can then be solved with state-of-the-art single-objective algorithms. Classical methods for multiobjective optimization fall into this category, but they all solve multiple single-objective problems subsequently (from scratch) instead of dynamically changing the scalarizing function during the search. On the other hand, we will improve on currently available population-based methods such as the first multiobjective versions of the CMA-ES. Here, research is needed on an even more fundamental level such as trying to understand success probabilities observed during an optimization run or how we can introduce non-elitist selection (the state of the art in single-objective stochastic adaptive algorithms) to increase robustness regarding noisy evaluations or multi-modality. The challenge here, compared to single-objective algorithms, is that the quality of a solution is not anymore independent from other sampled solutions, but can potentially depend on all known solutions (in the case of three or more objective functions), resulting in a more noisy evaluation as the relatively simple function-value-based ranking within single-objective optimizers.

### 3.3.4. Expensive Optimization

In the so-called expensive optimization scenario, a single function evaluation might take several minutes or even hours in a practical setting. Hence, the available budget in terms of number of function evaluation calls to find a solution is very limited in practice. To tackle such expensive optimization problems, it is needed to exploit the first few function evaluations in the best way. To this end, typical methods couple the learning of a surrogate (or meta-model) of the expensive objective function with traditional optimization algorithms.

In the context of expensive optimization and CMA-ES, which usually shows its full potential when the number  $n$  of variables is not too small (say larger than 3) and if the number of available function evaluations is about  $100n$  or larger, several research directions emerge. The two main possibilities to integrate meta-models into the search with CMA-ES type algorithms are (i) the successive injection of the minimum of a learned meta-model at each time step into the learning of CMA-ES's covariance matrix and (ii) the use of a meta-model to predict the internal ranking of solutions. While for the latter, first results exist, the former idea is entirely unexplored for now. In both cases, a fundamental question is which type of meta-model (linear, quadratic, Gaussian Process, ...) is the best choice for a given number of function evaluations (as low as one or two function evaluations) and at which time the type of the meta-model shall be switched.

## 3.4. Setting novel standards in scientific experimentation and benchmarking

Numerical experimentation is needed as a complement to theory to test novel ideas, hypotheses, the stability of an algorithm, and/or to obtain quantitative estimates. Optimally, theory and experimentation go hand in hand, jointly guiding the understanding of the mechanisms underlying optimization algorithms. Though performing numerical experimentation on optimization algorithms is crucial and a common task, it is non-trivial and easy to fall in (common) pitfalls as stated by J. N. Hooker in his seminal paper [27].

In the RandOpt team we aim at raising the standards for both scientific experimentation and benchmarking.

<sup>0</sup>Often, this set forms a manifold of dimension one smaller than the number of objectives.

On the experimentation aspect, we are convinced that there is common ground over how scientific experimentation should be done across many (sub-)domains of optimization, in particular with respect to the visualization of results, testing extreme scenarios (parameter settings, initial conditions, etc.), how to conduct understandable and small experiments, how to account for invariance properties, performing scaling up experiments and so forth. We therefore want to formalize and generalize these ideas in order to make them known to the entire optimization community with the final aim that they become standards for experimental research.

Extensive numerical benchmarking, on the other hand, is a compulsory task for evaluating and comparing the performance of algorithms. It puts algorithms to a standardized test and allows to make recommendations which algorithms should be used preferably in practice. To ease this part of optimization research, we have been developing the Comparing Continuous Optimizers platform (COCO) since 2007 which allows to automatize the tedious task of benchmarking. It is a game changer in the sense that the freed time can now be spent on the scientific part of algorithm design (instead of implementing the experiments, visualization, statistical tests, etc.) and it opened novel perspectives in algorithm testing. COCO implements a thorough, well-documented methodology that is based on the above mentioned general principles for scientific experimentation.

Also due to the freely available data from 300+ algorithms benchmarked with the platform, COCO became a quasi-standard for single-objective, noiseless optimization benchmarking. It is therefore natural to extend the reach of COCO towards other subdomains (particularly constrained optimization, many-objective optimization) which can benefit greatly from an automated benchmarking methodology and standardized tests without (much) effort. This entails particularly the design of novel test suites and rethinking the methodology for measuring performance and more generally evaluating the algorithms. Particularly challenging is the design of scalable non-trivial testbeds for constrained optimization where one can still control where the solutions lies. Other optimization problem types, we are targeting are expensive problems (and the Bayesian optimization community in particular, see our AESOP project), optimization problems in machine learning (for example parameter tuning in reinforcement learning), and the collection of real-world problems from industry.

Another aspect of our future research on benchmarking is to investigate the large amounts of benchmarking data, we collected with COCO during the years. Extracting information about the influence of algorithms on the best performing portfolio, clustering algorithms of similar performance, or the automated detection of anomalies in terms of good/bad behavior of algorithms on a subset of the functions or dimensions are some of the ideas here.

Last, we want to expand the focus of COCO from automatized (large) benchmarking experiments towards everyday experimentation, for example by allowing the user to visually investigate algorithm internals on the fly or by simplifying the set up of algorithm parameter influence studies.

## TAU Project-Team

### 3. Research Program

#### 3.1. Toward Good AI

As discussed by [141], the topic of ethical AI was non-existent until 2010, was laughed at in 2016, and became a hot topic in 2017 as the AI disruptivity with respect to the fabric of life (travel, education, entertainment, social networks, politics, to name a few) became unavoidable [138], together with its expected impacts on the nature and amount of jobs. As of now, it seems that the risk of a new AI Winter might arise from legal<sup>0</sup> and societal<sup>0</sup> issues. While privacy is now recognized as a civil right in Europe, it is feared that the GAFAM, BATX and others can already capture a sufficient fraction of human preferences and their dynamics to achieve their commercial and other goals, and build a Brave New Big Brother (BNBB, a system that is openly beneficial to many, covertly nudging, and possibly dictatorial).

The ambition of TAU is to mitigate the BNBB risk along several intricated dimensions, and build i) causal and explainable models; ii) fair data and models; iii) provably robust models.

##### 3.1.1. Causal modeling and biases

**Participants:** Isabelle Guyon, Michèle Sebag, Philippe Caillou, Paola Tubaro

**PhD:** Diviyam Kalainathan

**Collaboration:** Olivier Goudet (Université d'Angers), David Lopez-Paz (Facebook)

The extraction of causal models, a long goal of AI [139], [117], [140], became a strategic issue as the usage of learned models gradually shifted from *prediction* to *prescription* in the last years. This evolution, following Auguste Comte's vision of science (*Savoir pour prévoir, afin de pouvoir*) indeed reflects the exuberant optimism about AI: Knowledge enables Prediction; Prediction enables Control. However, although predictive models can be based on correlations, prescriptions can only be based on causal models<sup>0</sup>.

Among the research applications concerned with causal modeling, predictive modeling or collaborative filtering at TAU are all projects described in section 4.1 (see also Section 3.4), studying the relationships between: i) the educational background of persons and the job openings (FUI project JobAgile and DataIA project Vadore); ii) the quality of life at work and the economic performance indicators of the enterprises (ISN Lidex project Amiqap) [119]; iii) the nutritional items bought by households (at the level of granularity of the barcode) and their health status, as approximated from their body-mass-index (IRS UPSaclay Nutriperso); iv) the actual offer of restaurants and their scores on online rating systems. In these projects, a wealth of data is available (though hardly sufficient for applications ii), iii and iv))) and there is little doubt that these data reflect the imbalances and biases of the world as is, ranging from gender to racial to economical prejudices. Preventing the learned models from perpetuating such biases is essential to deliver an AI endowed with common decency.

In some cases, the bias is known; for instance, the cohorts in the Nutriperso study are more well-off than the average French population, and the Kantar database includes explicit weights to address this bias through importance sampling. In other cases, the bias is only guessed; for instance, the companies for which Secafi data are available hardly correspond to a uniform sample as these data have been gathered upon the request of the company trade union.

<sup>0</sup>For instance, the (fictitious) plea challenge proposed to law students in Oct. 2018 considered a chain reaction pileup occurred among autonomous and humanly operated vehicles on a highway.

<sup>0</sup>For instance related to information bubbles and nudge [100], [155].

<sup>0</sup>One can predict that it rains based on the presence of umbrellas in the street; but one cannot induce rainfall by going out with an umbrella. Likewise, the presence of books/tablets at home and the good scores of children at school are correlated; but offering books/tablets to all children might fail to improve their scores *per se*, if both good scores and books are explained by a so-called confounder variable, like the presence of adults versed in books/tablets at home.



### 3.1.2. Robustness of Learned Models

**Participants:** Guillaume Charpiat, Marc Schoenauer, Michèle Sebag

**PhDs:** Julien Girard, Marc Nabhan, Nizham Makhoud

**Collaboration:** Zakarian Chihani (CEA); Hiba Hage, and Yves Tourbier (Renault); Jérôme Kodjabachian (Thalès THERESIS)

Due to their outstanding performances, deep neural networks and more generally machine learning-based decision making systems, referred to as MLs in the following, have been raising hopes in the recent years to achieve breakthroughs in critical systems, ranging from autonomous vehicles to defense. The main pitfall for such applications lies in the lack of guarantees for MLs robustness.

Specifically, MLs are used when the mainstream software design process does not apply, that is, when no formal specification of the target software behavior is available and/or when the system is embedded in an open unpredictable world. The extensive body of knowledge developed to deliver guarantees about mainstream software – ranging from formal verification, model checking and abstract interpretation to testing, simulation and monitoring – thus does not directly apply either. Another weakness of MLs regards their dependency to the amount and quality of the training data, as their performances are sensitive to slight perturbations of the data distribution. Such perturbations can occur naturally due to domain or concept drift (e.g. due to a change in light intensity or a scratch on a camera lens); they can also result from intentional malicious attacks, a.k.a adversarial examples [156].

These downsides, currently preventing the dissemination of MLs in safety-critical systems (SCS), call for a considerable amount of research, in order to understand when and to which extent an MLs can be certified to provide the desired level of guarantees.

Julien Girard's PhD (CEA scholarship), started in Oct. 2018, co-supervised by Guillaume Charpiat and Zakaria Chihani (CEA), is devoted to the extension of abstract interpretation to deep neural nets, and the formal characterization of the transition kernel from input to output space achieved by a DNN (robustness by design, coupled with formally assessing the coverage of the training set). This approach is tightly related to the inspection and opening of black-box models, aimed to characterize the patterns in the input instances responsible for a decision – another step toward explainability.

On the other hand, experimental validation of MLs, akin statistical testing, also faces three limitations: i) real-world examples are notoriously insufficient to ensure a good coverage in general; ii) for this reason, simulated examples are extensively used; but their use raises the *reality gap* issue [128] of the distance between real and simulated worlds; iii) independently, the real-world is naturally subject to domain shift (e.g. due to the technical improvement and/or aging of sensors). Our collaborations with Renault tackle such issues in the context of the autonomous vehicle (see Section 7.1.3).

## 3.2. Hybridizing numerical modeling and learning systems

**Participants:** Alessandro Bucci, Guillaume Charpiat, Cécile Germain, Isabelle Guyon, Marc Schoenauer, Michèle Sebag

**PhD:** Théophile Sanchez, Loris Felardos, Wenzhuo Liu

In sciences and engineering, human knowledge is commonly expressed in closed form, through equations or mechanistic models characterizing how a natural or social phenomenon, or a physical device, will behave/evolve depending on its environment and external stimuli, under some assumptions and up to some approximations. The field of numerical engineering, and the simulators based on such mechanistic models, are at the core of most approaches to understand and analyze the world, from solid mechanics to computational fluid dynamics, from chemistry to molecular biology, from astronomy to population dynamics, from epidemiology and information propagation in social networks to economy and finance.

Most generally, numerical engineering supports the simulation, and when appropriate the optimization and control<sup>0</sup> of the phenomenons under study, although several sources of discrepancy might adversely affect the results, ranging from the underlying assumptions and simplifying hypotheses in the models, to systematic experiment errors to statistical measurement errors (not to mention numerical issues). This knowledge and know-how are materialized in millions of lines of code, capitalizing the expertise of academic and industrial labs. These softwares have been steadily extended over decades, modeling new and more fine-grained effects through layered extensions, making them increasingly harder to maintain, extend and master. Another difficulty is that complex systems most often resort to hybrid (pluridisciplinary) models, as they involve many components interacting along several time and space scales, hampering their numerical simulation.

At the other extreme, machine learning offers the opportunity to model phenomenons from scratch, using any available data gathered through experiments or simulations. Recent successes of machine learning in computer vision, natural language processing and games to name a few, have demonstrated the power of such agnostic approaches and their efficiency in terms of prediction [123], inverse problem solving [170], and sequential decision making [162], [81], despite their lack of any "semantic" understanding of the universe. Even before these successes, Anderson's claim was that *the data deluge [might make] the scientific method obsolete* [70], as if a reasonable option might be to throw away the existing equational or software bodies of knowledge, and let Machine Learning rediscover all models from scratch. Such a claim is hampered among others by the fact that not all domains offer a wealth of data, as any academic involved in an industrial collaboration around data has discovered.

Another approach will be considered in TAU, investigating how existing mechanistic models and related simulators can be partnered with ML algorithms: i) to achieve the same goals with the same methods with a gain of accuracy or time; ii) to achieve new goals; iii) to achieve the same goals with new methods.

**Toward more robust numerical engineering:** In domains where satisfying mechanistic models and simulators are available, ML can contribute to improve their accuracy or usability. A first direction is to refine or extend the models and simulators to better fit the empirical evidence. The goal is to finely account for the different biases and uncertainties attached to the available knowledge and data, distinguishing the different types of *known unknowns*. Such *known unknowns* include the model hyper-parameters (coefficients), the systematic errors due to e.g., experiment imperfections, and the statistical errors due to e.g., measurement errors. A second approach is based on learning a surrogate model for the phenomenon under study that incorporate domain knowledge from the mechanistic model (or its simulation). See Section 7.5 for case studies.

A related direction, typically when considering black-box simulators, aims to learn a model of the error, or equivalently, a post-processor of the software. The discrepancy between simulated and empirical results, referred to as *reality gap* [128], can be tackled in terms of domain adaptation [74], [99]. Specifically, the source domain here corresponds to the simulated phenomenon, offering a wealth of inexpensive data, and the target domain corresponds to the actual phenomenon, with rare and expensive data; the goal is to devise accurate target models using the source data and models.

**Extending numerical engineering:** ML, using both experimental and numerical data, can also be used to tackle new goals, that are beyond the current state-of-the-art of standard approaches. Inverse problems are such goals, identifying the parameters or the initial conditions of phenomenons for which the model is not differentiable, or amenable to the adjoint state method.

A slightly different kind of inverse problem is that of recovering the ground truth when only noisy data is available. This problem can be formulated as a search for the simplest model explaining the data. The question then becomes to formulate and efficiently exploit such a simplicity criterion.

Another goal can be to model the distribution of given quantiles for some system: The challenge is to exploit available data to train a generative model, aimed at sampling the target quantiles.

---

<sup>0</sup>Note that the causal nature of mechanistic models is established from prior knowledge and experimentations.

Examples tackled in TAU are detailed in Section 7.5 . Note that the "Cracking the Glass Problem", described in Section 7.2.3 is yet another instance of a similar problem.

**Data-driven numerical engineering** : Finally, ML can also be used to sidestep numerical engineering limitations in terms of scalability, or to build a simulator emulating the resolution of the (unknown) mechanistic model from data, or to revisit the formal background.

When the mechanistic model is known and sufficiently accurate, it can be used to train a deep network on an arbitrary set of (space,time) samples, resulting in a meshless numerical approximation of the model [151], supporting by construction *differentiable programming* [125].

When no mechanistic model is sufficiently efficient, the model must be identified from the data only. Genetic programming has been used to identify systems of ODEs [149], through the identification of invariant quantities from data, as well as for the direct identification of control commands of nonlinear complex systems, including some chaotic systems [88]. Another recent approach uses two deep neural networks, one for the state of the system, the other for the equation itself [142]. The critical issues for both approaches include the scalability, and the explainability of the resulting models. Such line of research will benefit from TAU unique mixed expertise in Genetic Programming and Deep Learning.

Finally, in the realm of signal processing (SP), the question is whether and how deep networks can be used to revisit mainstream feature extraction based on Fourier decomposition, wavelet and scattering transforms [76]. E. Bartenlian's PhD (started Oct. 2018), co-supervised by M. Sebag and F. Pascal (Centrale-Supélec), focusing on musical audio-to-score translation [150], inspects the effects of supervised training, taking advantage from the fact that convolution masks can be initialized and analyzed in terms of frequency.

### 3.3. Learning to learn

According to Ali Rahimi's test of times award speech at NIPS 17, the current ML algorithms *have become a form of alchemy*. Competitive testing and empirical breakthroughs gradually become mandatory for a contribution to be acknowledged; an increasing part of the community adopts trials and errors as main scientific methodology, and theory is lagging behind practice. This style of progress is typical of technological and engineering revolutions for some; others ask for consolidated and well-understood theoretical advances, saving the time wasted in trying to build upon hardly reproducible results.

Basically, while practical achievements have often passed the expectations, there exist caveats along three dimensions. Firstly, excellent performances do not imply that the model has captured what was to learn, as shown by the phenomenon of adversarial examples. Following Ian Goodfellow, some well-performing models might be compared to *Clever Hans*, the horse that was able to solve mathematical exercises using non verbal cues from its teacher [116]; it is the purpose of Pillar I. to alleviate the *Clever Hans* trap (section 3.1).

Secondly, some major advances, e.g. related to the celebrated adversarial learning [105], [99], establish proofs of concept more than a sound methodology, where the reproducibility is limited due to i) the computational power required for training (often beyond reach of academic labs); ii) the numerical instabilities (witnessed as random seeds happen to be found in the codes); iii) the insufficiently documented experimental settings. What works, why and when is still a matter of speculation, although better understanding the limitations of the current state of the art is acknowledged to be a priority. After Ali Rahimi again, *simple experiments, simple theorems are the building blocks that help us understand more complicated systems*. Along this line, [135] propose toy examples to demonstrate and understand the defaults of convergence of gradient descent adversarial learning.

Thirdly, and most importantly, the reported achievements rely on carefully tuned learning architectures and hyper-parameters. The sensitivity of the results to the selection and calibration of algorithms has been identified since the end 80s as a key ML bottleneck, and the field of automatic algorithm selection and calibration, referred to as AutoML or Auto-☆ in the following, is at the ML forefront.

TAU aims to contribute to the ML evolution toward a more mature stage along three dimensions. In the short term, the research done in Auto- $\star$  will be pursued (section 3.3.1). In the medium term, an information theoretic perspective will be adopted to capture the data structure and to calibrate the learning algorithm *depending on the nature and amount of the available data* (section 3.3.2). In the longer term, our goal is to leverage the methodologies forged in statistical physics to understand and control the trajectories of complex learning systems (section 3.3.3).

### 3.3.1. Auto-\*

**Participants:** Isabelle Guyon, Marc Schoenauer, Michèle Sebag

**PhD:** Guillaume Doquet, Zhengying Liu, Herilalaina Rakotoarison, Lisheng Sun

**Collaboration:** Olivier Bousquet, André Elisseeff (Google Zurich)

The so-called Auto- $\star$  task, concerned with selecting a (quasi) optimal algorithm and its hyper-parameters depending on the problem instance at hand, remained a key issue in ML for the last three decades [75], as well as in optimization at large [115], including combinatorial optimization and constraint satisfaction [122], [104] and continuous optimization [71]. This issue, tackled by several European projects along the decades, governs the knowledge transfer to industry, due to the shortage of data scientists. It becomes even more crucial as models are more complex and their training requires more computational resources. This has motivated several international challenges devoted to Auto-ML [113] (see also Section 3.4), including the AutoDL challenge series [129] launched in 2019<sup>0</sup> (see also Section 7.6).

Several approaches have been used to tackle Auto- $\star$  in the literature, and TAU has been particularly active in several of them. Meta-learning aims to build a surrogate performance model, estimating the performance of an algorithm configuration on *any* problem instance characterized from its meta-feature values [146], [104], [72], [71], [103]. Collaborative filtering, considering that a problem instance "likes better" an algorithm configuration yielding a better performance, learns to recommend good algorithms to problem instances [153], [137]. Bayesian optimization proceeds by alternatively building a surrogate model of algorithm performances on *the* problem instance at hand, and tackling it [95]. This last approach currently is the prominent one; as shown in [137], the meta-features developed for AutoML are hardly relevant, hampering both meta-learning and collaborative filtering. The design of better features is another long-term research direction, in which TAU has recently been [32], and still is very active. more recent approach used in TAU [40] extends the Bayesian Optimization approach with a Multi-Armed Bandit algorithm to generate the full Machine Learning pipeline, competing with the famed AutoSKLearn [95] (see Section 7.2.1). These results are presented in Section 7.2.1

### 3.3.2. Information theory: adjusting model complexity and data fitting

**Participants:** Guillaume Charpiat, Marc Schoenauer, Michèle Sebag

**PhD:** Corentin Tallec, Pierre Wolinski, Léonard Blier

**Collaboration:** Yann Ollivier (Facebook)

In the 60s, Kolmogorov and Solomonoff provided a well-grounded theory for building (probabilistic) models best explaining the available data [147], [108], that is, the shortest programs able to generate these data. Such programs can then be used to generate further data or to answer specific questions (interpreted as missing values in the data). Deep learning, from this viewpoint, efficiently explores a space of computation graphs, described from its hyperparameters (network structure) and parameters (weights). Network training amounts to optimizing these parameters, namely, navigating the space of computational graphs to find a network, as simple as possible, that explain the past observations well.

This vision is at the core of variational auto-encoders [121], directly optimizing a bound on the Kolmogorov complexity of the dataset. More generally variational methods provide quantitative criteria to identify superfluous elements (edges, units) in a neural network, that can potentially be used for structural optimization of the network (Leonard Blier's PhD, started Oct. 2018).

<sup>0</sup><https://autodl.chalearn.org/neurips2019>

The same principles apply to unsupervised learning, aimed to find the maximum amount of structure hidden in the data, quantified using this information-theoretic criterion.

The known invariances in the data can be exploited to guide the model design (e.g. as translation invariance leads to convolutional structures, or LSTM is shown to enforce the invariance to time affine transformations of the data sequence [157]). Scattering transforms exploit similar principles [76]. A general theory of how to detect *unknown* invariances in the data, however, is currently lacking.

The view of information theory and Kolmogorov complexity suggests that key program operations (composition, recursivity, use of predefined routines) should intervene when searching for a good computation graph. One possible framework for exploring the space of computation graphs with such operations is that of Genetic Programming. It is interesting to see that evolutionary computation appeared in the last two years among the best candidates to explore the space of deep learning structures [145], [126]. Other approaches might proceed by combining simple models into more powerful ones, e.g. using “Context Tree Weighting” [166] or switch distributions [90]. Another option is to formulate neural architecture design as a reinforcement learning problem [73]; the value of the building blocks (predefined routines) might be defined using e.g., Monte-Carlo Tree Search. A key difficulty is the computational cost of retraining neural nets from scratch upon modifying their architecture; an option might be to use neutral initializations to support warm-restart.

### 3.3.3. Analyzing and Learning Complex Systems

**Participants:** Cyril Furtlehner, Aurélien Decelle, François Landes, Michèle Sebag

**PhD:** Giancarlo Fissore

**Collaboration:** Enrico Camporeale (CWI); Jacopo Rocchi (LPTMS Paris Sud), the Simons team: Rahul Chako (post-doc), Andrea Liu (UPenn), David Reichman (Columbia), Giulio Biroli (ENS), Olivier Dauchot (ESPCI), Hufei Han (Symantec).

Methods and criteria from statistical physics have been widely used in ML. In early days, the capacity of Hopfield networks (associative memories defined by the attractors of an energy function) was investigated by using the replica formalism [69]. Restricted Boltzmann machines likewise define a generative model built upon an energy function trained from the data. Along the same lines, Variational Auto-Encoders can be interpreted as systems relating the free energy of the distribution, the information about the data and the entropy (the degree of ignorance about the micro-states of the system) [165]. A key promise of the statistical physics perspective and the Bayesian view of deep learning is to harness the tremendous growth of the model size (billions of weights in recent machine translation networks), and make them sustainable through e.g. posterior drop-out [136], weight quantization and probabilistic binary networks [131]. Such “informational cooling” of a trained deep network can reduce its size by several orders of magnitude while preserving its performance.

Statistical physics is among the key expertises of TAU, originally only represented by Cyril Furtlehner, later strengthened by Aurélien Decelle’s and François Landes’ arrivals in 2014 and 2018. On-going studies are conducted along several directions.

Generative models are most often expressed in terms of a Gibbs distributions  $P[S] = \exp(-E[S])$ , where energy  $E$  involves a sum of building blocks, modelling the interactions among variables. This formalization makes it natural to use mean-field methods of statistical physics and associated inference algorithms to both train and exploit such models. The difficulty is to find a good trade-off between the richness of the structure and the efficiency of mean-field approaches. One direction of research pursued in TAU, [97] in the context of traffic forecasting, is to account for the presence of cycles in the interaction graph, to adapt inference algorithms to such graphs with cycles, while constraining graphs to remain compatible with mean-field inference.

Another direction, explored in TAO/TAU in the recent years, is based on the definition and exploitation of self-consistency properties, enforcing principled divide-and-conquer resolutions. In the particular case of the message-passing Affinity Propagation algorithm for instance [168], self-consistency imposes the invariance of the solution when handled at different scales, thus enabling to characterize the critical value of the penalty and other hyper-parameters in closed form (in the case of simple data distributions) or empirically otherwise [98].

A more recent research direction examines the quantity of information in a (deep) neural net along the random matrix theory framework [78]. It is addressed in Giancarlo Fissore's PhD, and is detailed in Section 7.2.3 .

Finally, we note the recent surge in using ML to address fundamental physics problems: from turbulence to high-energy physics and soft matter as well (with amorphous materials at its core) [19]. TAU's dual expertise in Deep Networks and in statistical physics places it in an ideal position to significantly contribute to this domain and shape the methods that will be used by the physics community in the future. François Landes' recent arrival in the team makes TAU a unique place for such interdisciplinary research, thanks to his collaborators from the *Simons Collaboration Cracking the Glass Problem* (gathering 13 statistical physics teams at the international level). This project is detailed in Section 7.2.3 .

Independently, François Landes is actively collaborating with statistical physicists (Alberto Rosso, LPTMS, Univ. Paris-Saclay) and physicists at the frontier with geophysics (Eugenio Lippiello, Second Univ. of Naples) [20]. A possible CNRS grant (80Prime) may finance a shared PhD, at the frontier between seismicity and ML (Alberto Rosso, Marc Schoenauer and François Landes).

### 3.4. Organisation of Challenges

**Participants:** Cécile Germain, Isabelle Guyon, Marc Schoenauer, Michèle Sebag

Challenges have been an important drive for Machine Learning research for many years, and TAO members have played important roles in the organization of many such challenges: Michèle Sebag was head of the challenge programme in the *Pascal European Network of Excellence* (2005-2013); Isabelle Guyon, as mentioned, was the PI of many challenges ranging from causation challenges [109], to AutoML [110]. The *Higgs challenge* [67], most attended ever Kaggle challenge, was jointly organized by TAO (C. Germain), LAL-IN2P3 (D. Rousseau and B. Kegl) and I. Guyon (not yet at TAO), in collaboration with CERN and Imperial College.

TAU was also particularly implicated with the ChaLearn Looking At People (LAP) challenge series in Computer Vision, in collaboration with the University of Barcelona [92] including the *Job Candidate Screening Coepetition* [91]; the *Real Versus Fake Expressed Emotion Challenge* (ICCV 2017) [163]; the *Large-scale Continuous Gesture Recognition Challenge* (ICCV 2017) [163]; the *Large-scale Isolated Gesture Recognition Challenge* (ICCV 2017) [163].

Other challenges have been organized in 2019, or are planned for the near future, detailed in Section 7.6 . In particular, many of them now run on the Codalab platform, managed by Isabelle Guyon and maintained at LRI.

## TROPICAL Project-Team

### 3. Research Program

#### 3.1. Optimal control and zero-sum games

The dynamic programming approach allows one to analyze one or two-player dynamic decision problems by means of operators, or partial differential equations (Hamilton–Jacobi or Isaacs PDEs), describing the time evolution of the value function, i.e., of the optimal reward of one player, thought of as a function of the initial state and of the horizon. We work especially with problems having long or infinite horizon, modelled by stopping problems, or ergodic problems in which one optimizes a mean payoff per time unit. The determination of optimal strategies reduces to solving nonlinear fixed point equations, which are obtained either directly from discrete models, or after a discretization of a PDE.

**The geometry of solutions of optimal control and game problems** Basic questions include, especially for stationary or ergodic problems, the understanding of existence and uniqueness conditions for the solutions of dynamic programming equations, for instance in terms of controllability or ergodicity properties, and more generally the understanding of the structure of the full set of solutions of stationary Hamilton–Jacobi PDEs and of the set of optimal strategies. These issues are already challenging in the one-player deterministic case, which is an application of choice of tropical methods, since the Lax–Oleinik semigroup, i.e., the evolution semigroup of the Hamilton–Jacobi PDE, is a linear operator in the tropical sense. Recent progress in the deterministic case has been made by combining dynamical systems and PDE techniques (weak KAM theory [72]), and also using metric geometry ideas (abstract boundaries can be used to represent the sets of solutions [86], [4]). The two player case is challenging, owing to the lack of compactness of the analogue of the Lax–Oleinik semigroup and to a richer geometry. The conditions of solvability of ergodic problems for games (for instance, solvability of ergodic Isaacs PDEs), and the representation of solutions are only understood in special cases, for instance in the finite state space case, through tropical geometry and non-linear Perron–Frobenius methods [38], [41], [3].

**Algorithmic aspects: from combinatorial algorithms to the attenuation of the curse of dimensionality**

Our general goal is to push the limits of solvable models by means of fast algorithms adapted to large scale instances. Such instances arise from discrete problems, in which the state space may be so large that it is only accessible through local oracles (for instance, in some web ranking applications, the number of states may be the number of web pages) [73]. They also arise from the discretization of PDEs, in which the number of states grows exponentially with the number of degrees of freedom, according to the “curse of dimensionality”. A first line of research is the development of *new approximation methods for the value function*. So far, classical approximations by linear combinations have been used, as well as approximation by suprema of linear or quadratic forms, which have been introduced in the setting of dual dynamic programming and of the so called “max-plus basis methods” [74]. We believe that more concise or more accurate approximations may be obtained by unifying these methods. Also, some max-plus basis methods have been shown to *attenuate the curse of dimensionality* for very special problems (for instance involving switching) [97], [78]. This suggests that the complexity of control or games problems may be measured by more subtle quantities than the mere number of states, for instance, by some forms of metric entropy (for example, certain large scale problems have a low complexity owing to the presence of decomposition properties, “highway hierarchies”, etc.). A second line of our research is the development of *combinatorial algorithms*, to solve large scale zero-sum two-player problems with discrete state space. This is related to current open problems in algorithmic game theory. In particular, the existence of polynomial-time algorithms for games with ergodic payment is an open question. See e.g. [43] for a polynomial time average complexity result derived by tropical methods. The two lines of research are related, as the understanding of the geometry of solutions allows to develop better approximation or combinatorial algorithms.

### 3.2. Non-linear Perron-Frobenius theory, nonexpansive mappings and metric geometry

Several applications (including population dynamics [10] and discrete event systems [56], [64], [46]) lead to studying classes of dynamical systems with remarkable properties: preserving a cone, preserving an order, or being nonexpansive in a metric. These can be studied by techniques of non-linear Perron-Frobenius theory [3] or metric geometry [11]. Basic issues concern the existence and computation of the “escape rate” (which determines the throughput, the growth rate of the population), the characterizations of stationary regimes (non-linear fixed points), or the study of the dynamical properties (convergence to periodic orbits). Nonexpansive mappings also play a key role in the “operator approach” to zero-sum games, since the one-day operators of games are nonexpansive in several metrics, see [8].

### 3.3. Tropical algebra and convex geometry

The different applications mentioned in the other sections lead us to develop some basic research on tropical algebraic structures and in convex and discrete geometry, looking at objects or problems with a “piecewise-linear” structure. These include the geometry and algorithmics of tropical convex sets [49], [40], tropical semialgebraic sets [52], the study of semi-modules (analogues of vector spaces when the base field is replaced by a semi-field), the study of systems of equations linear in the tropical sense, investigating for instance the analogues of the notions of rank, the analogue of the eigenproblems [42], and more generally of systems of tropical polynomial equations. Our research also builds on, and concerns, classical convex and discrete geometry methods.

### 3.4. Tropical methods applied to optimization, perturbation theory and matrix analysis

Tropical algebraic objects appear as a deformation of classical objects through various asymptotic procedures. A familiar example is the rule of asymptotic calculus,

$$e^{-a/\epsilon} + e^{-b/\epsilon} \asymp e^{-\min(a,b)/\epsilon}, \quad e^{-a/\epsilon} \times e^{-b/\epsilon} = e^{-(a+b)/\epsilon}, \quad (1)$$

when  $\epsilon \rightarrow 0^+$ . Deformations of this kind have been studied in different contexts: large deviations, zero-temperature limits, Maslov’s “dequantization method” [96], non-archimedean valuations, log-limit sets and Viro’s patchworking method [122], etc.

This entails a relation between classical algorithmic problems and tropical algorithmic problems, one may first solve the  $\epsilon = 0$  case (non-archimedean problem), which is sometimes easier, and then use the information gotten in this way to solve the  $\epsilon = 1$  (archimedean) case.

In particular, tropicalization establishes a connection between polynomial systems and piecewise affine systems that are somehow similar to the ones arising in game problems. It allows one to transfer results from the world of combinatorics to “classical” equations solving. We investigate the consequences of this correspondence on complexity and numerical issues. For instance, combinatorial problems can be solved in a robust way. Hence, situations in which the tropicalization is faithful lead to improved algorithms for classical problems. In particular, scalings for the polynomial eigenproblems based on tropical preprocessings have started to be used in matrix analysis [80], [84].

Moreover, the tropical approach has been recently applied to construct examples of linear programs in which the central path has an unexpectedly high total curvature [44], and it has also led to positive polynomial-time average case results concerning the complexity of mean payoff games. Similarly, we are studying semidefinite programming over non-archimedean fields [52], [51], with the goal to better understand complexity issues in classical semidefinite and semi-algebraic programming.



## LIFEWARE Project-Team

### 3. Research Program

#### 3.1. Computational Systems Biology

Bridging the gap between the complexity of biological systems and our capacity to model and **quantitatively predict system behaviors** is a central challenge in systems biology. We believe that a deeper understanding of the concept and theory of biochemical computation is necessary to tackle that challenge. Progress in the theory is necessary for scaling, and enabling the application of static analysis, module identification and decomposition, model reductions, parameter search, and model inference methods to large biochemical reaction systems. A measure of success on this route will be the production of better computational modeling tools for elucidating the complex dynamics of natural biological processes, designing synthetic biological circuits and biosensors, developing novel therapy strategies, and optimizing patient-tailored therapeutics.

Progress on the **coupling of models to data** is also necessary. Our approach based on quantitative temporal logics provides a powerful framework for formalizing experimental observations and using them as formal specification in model building. Key to success is a tight integration between *in vivo* and *in silico* work, and on the mixing of dry and wet experiments, enabled by novel biotechnologies. In particular, the use of microfluidic devices makes it possible to measure behaviors at both single-cell and cell population levels *in vivo*, provided innovative modeling, analysis and control methods are deployed *in silico*.

In synthetic biology, while the construction of simple intracellular circuits has shown feasible, the design of larger, **multicellular systems** is a major open issue. In engineered tissues for example, the behavior results from the subtle interplay between intracellular processes (signal transduction, gene expression) and intercellular processes (contact inhibition, gradient of diffusible molecule), and the question is how should cells be genetically modified such that the desired behavior robustly emerges from cell interactions.

#### 3.2. Chemical Reaction Network (CRN) Theory

Feinberg's chemical reaction network theory and Thomas's influence network analyses provide sufficient and/or necessary structural conditions for the existence of multiple steady states and oscillations in regulatory networks. Those conditions can be verified by static analyzers without knowing kinetic parameter values nor making any simulation. In this domain, most of our work consists in analyzing the interplay between the **structure** (Petri net properties, influence graph, subgraph epimorphisms) and the **dynamics** (Boolean, CTMC, ODE, time scale separations) of biochemical reaction systems. In particular, our study of influence graphs of reaction systems, our generalization of Thomas' conditions of multi-stationarity and Soulé's proof to reaction systems<sup>0</sup>, the inference of reaction systems from ODEs<sup>0</sup>, the computation of structural invariants by constraint programming techniques, and the analysis of model reductions by subgraph epimorphisms now provide solid ground for developing static analyzers, using them on a large scale in systems biology, and elucidating modules.

#### 3.3. Logical Paradigm for Systems Biology

Our group was among the first ones in 2002 to apply **model-checking** methods to systems biology in order to reason on large molecular interaction networks, such as Kohn's map of the mammalian cell cycle (800 reactions over 500 molecules)<sup>0</sup>. The logical paradigm for systems biology that we have subsequently developed for quantitative models can be summarized by the following identifications :

<sup>0</sup>Sylvain Soliman. A stronger necessary condition for the multistationarity of chemical reaction networks. *Bulletin of Mathematical Biology*, 75(11):2289–2303, 2013.

<sup>0</sup>François Fages, Steven Gay, Sylvain Soliman. Inferring reaction systems from ordinary differential equations. *Journal of Theoretical Computer Science (TCS)*, Elsevier, 2015, 599, pp.64–78.

<sup>0</sup>N. Chabrier-Rivier, M. Chiaverini, V. Danos, F. Fages, V. Schächter. Modeling and querying biochemical interaction networks. *Theoretical Computer Science*, 325(1):25–44, 2004.

biological model = transition system  $K$   
 dynamical behavior specification = temporal logic formula  $\phi$   
 model validation = model-checking  $K, s \models \phi$   
 model reduction = sub-model-checking,  $K' \subset K$  s.t.  $K' \models \phi, s \models \phi$   
 model prediction = formula enumeration,  $\phi$  s.t.  $K, s \models \phi$   
 static experiment design = symbolic model-checking, state  $s$  s.t.  $K, s \models \phi$   
 model synthesis = constraint solving  $K?, s \models \phi$   
 dynamic experiment design = constraint solving  $K?, s? \models \phi$

In particular, the definition of a continuous satisfaction degree for **first-order temporal logic** formulae with constraints over the reals, was the key to generalize this approach to quantitative models, opening up the field of model-checking to model optimization<sup>0</sup> This line of research continues with the development of temporal logic patterns with efficient constraint solvers and their generalization to handle stochastic effects.

### 3.4. Computer-Aided Design of CRNs for Synthetic Biology

The continuous nature of many protein interactions leads us to consider models of analog computation, and in particular, the recent results in the theory of analog computability and complexity obtained by Amaury Pouly<sup>0</sup> and Olivier Bournez, establish fundamental links with digital computation. In a paper published last year<sup>0</sup> We have derived from these results the Turing completeness result of elementary CRNs (without polymerization) under the differential semantics, closing a long-standing open problem in CRN theory. The proof of this result shows how computable function over the reals, described by Ordinary Differential Equations, namely by Polynomial Initial Value Problems (PIVP), can be compiled into elementary biochemical reactions, furthermore with a notion of analog computation complexity defined as the length of the trajectory to reach a given precision on the result. This opens a whole research avenue to analyze biochemical circuits in Systems Biology, transform behavioural specifications into biochemical reactions for Synthetic Biology, and compare artificial circuits with natural circuits acquired through evolution, from the novel point of view of analog computation and complexity.

### 3.5. Modeling of Phenotypic Heterogeneity in Cellular Processes

Since nearly two decades, a significant interest has grown for getting a quantitative understanding of the functioning of biological systems at the cellular level. Given their complexity, proposing a model accounting for the observed cell responses, or better, predicting novel behaviors, is now regarded as an essential step to validate a proposed mechanism in systems biology. Moreover, the constant improvement of stimulation and observation tools creates a strong push for the development of methods that provide predictions that are increasingly precise (single cell precision) and robust (complex stimulation profiles).

It is now fully apparent that cells do not respond identically to a same stimulation, even when they are all genetically-identical. This phenotypic heterogeneity plays a significant role in a number of problems ranging from cell resistance to anticancer drug treatments to stress adaptation and bet hedging.

<sup>0</sup>On a continuous degree of satisfaction of temporal logic formulae with applications to systems biology A. Rizk, G. Batt, F. Fages, S. Soliman International Conference on Computational Methods in Systems Biology, 251-268

<sup>0</sup>Amaury Pouly, "Continuous models of computation: from computability to complexity", PhD Thesis, Ecole Polytechnique, Nov. 2015.

<sup>0</sup>Fages, François, Le Guludec, Guillaume and Bournez, Olivier, Pouly, Amaury. Strong Turing Completeness of Continuous Chemical Reaction Networks and Compilation of Mixed Analog-Digital Programs. In CMSB'17: Proceedings of the fifteen international conference on Computational Methods in Systems Biology, pages 108–127, volume 10545 of Lecture Notes in Computer Science. Springer-Verlag, 2017.

Dedicated modeling frameworks, notably **stochastic** modeling frameworks, such as chemical master equations, and **statistic** modeling frameworks, such as ensemble models, are then needed to capture biological variability.

Appropriate mathematical and computational tools should then be employed for the analysis of these models and their calibration to experimental data. One can notably mention **global optimization** tools to search for appropriate parameters within large spaces, **moment closure** approaches to efficiently approximate stochastic models<sup>0</sup>, and (stochastic approximations of) the **expectation maximization** algorithm for the identification of mixed-effects models<sup>0</sup>.

### 3.6. External Control of Cell Processes

External control has been employed since many years to regulate culture growth and other physiological properties. Recently, taking inspiration from developments in synthetic biology, closed loop control has been applied to the regulation of intracellular processes. Such approaches offer unprecedented opportunities to investigate how a cell process dynamical information by maintaining it around specific operating points or driving it out of its standard operating conditions. They can also be used to complement and help the development of synthetic biology through the creation of hybrid systems resulting from the interconnection of in vivo and in silico computing devices.

In collaboration with Pascal Hersen (CNRS MSC lab), we developed a platform for gene expression control that enables to control protein concentrations in yeast cells. This platform integrates microfluidic devices enabling long-term observation and rapid change of the cells environment, microscopy for single cell measurements, and software for real-time signal quantification and model based control. We demonstrated in 2012 that this platform enables controlling the level of a fluorescent protein in cells with unprecedented accuracy and for many cell generations<sup>0</sup>.

More recently, motivated by an analogy with a benchmark control problem, the stabilization of an inverted pendulum, we investigated the possibility to balance a genetic toggle switch in the vicinity of its unstable equilibrium configuration. We searched for solutions to balance an individual cell and even an entire population of heterogeneous cells, each harboring a toggle switch<sup>0</sup>.

Independently, in collaboration with colleagues from IST Austria, we investigated the problem of controlling cells, one at a time, by constructing an integrated optogenetic-enabled microscopy platform. It enables experiments that bridge individual and population behaviors. We demonstrated: (i) population structuring by independent closed-loop control of gene expression in many individual cells, (ii) cell-cell variation control during antibiotic perturbation, (iii) hybrid bio-digital circuits in single cells, and freely specifiable digital communication between individual bacteria<sup>0</sup>.

### 3.7. Constraint Solving and Optimization

Constraint solving and optimization methods are important in our research. On the one hand, static analysis of biochemical reaction networks involves solving hard combinatorial optimization problems, for which **constraint programming** techniques have shown particularly successful, often beating dedicated algorithms

<sup>0</sup>Moment-based inference predicts bimodality in transient gene expression, C. Zechner C, J. Ruess, P. Krenn, S. Pelet, M. Peter, J. Lygeros, and H. Koepl, Proceedings of the National Academy of Sciences USA, 9(5):109(21):8340-5, 2012

<sup>0</sup>What population reveals about individual cell identity: estimation of single-cell models of gene expression in yeast, A. Llamasi, A.M. Gonzalez-Vargas, C. Versari, E. Cinquemani, G. Ferrari-Trecate, P. Hersen, and G. Batt, PLoS Computational Biology, 9(5): e1003056, 2015

<sup>0</sup>Jannis Uhlendorf, Agnès Miermont, Thierry Delaveau, Gilles Charvin, François Fages, Samuel Bottani, Grégory Batt, Pascal Hersen. Long-term model predictive control of gene expression at the population and single-cell levels. Proceedings of the National Academy of Sciences USA, 109(35):14271–14276, 2012.

<sup>0</sup>Jean-Baptiste Lugagne, Sebastian Sosa Carrillo and Melanie Kirch, Agnes Köhler, Gregory Batt and Pascal Hersen. Balancing a genetic toggle switch by real-time feedback control and periodic forcing. Nature Communications, 8(1):1671, 2017.

<sup>0</sup>Remy Chait, Jakob Ruess, Tobias Bergmiller and Gavsper Tkavcik, Cvalin Guet. Shaping bacterial population behavior through computer-interfaced control of individual cells. Nature Communications, 8(1):1535, 2017.

and allowing to solve large instances from model repositories. On the other hand, parameter search and model calibration problems involve similarly solving hard continuous optimization problems, for which **evolutionary algorithms**, and especially the covariance matrix evolution strategy (**CMA-ES**)<sup>0</sup> have been shown to provide best results in our context, for up to 100 parameters. This has been instrumental in building challenging quantitative models, gaining model-based insights, revisiting admitted assumptions, and contributing to biological knowledge<sup>00</sup>.

---

<sup>0</sup>N. Hansen, A. Ostermeier (2001). Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2) pp. 159–195.

<sup>0</sup>Domitille Heitzler, Guillaume Durand, Nathalie Gallay, Aurélien Rizk, Seungki Ahn, Jihee Kim, Jonathan D. Violin, Laurence Dupuy, Christophe Gauthier, Vincent Piketty, Pascale Crépieux, Anne Poupon, Frédérique Clément, François Fages, Robert J. Lefkowitz, Eric Reiter. Competing G protein-coupled receptor kinases balance G protein and  $\beta$ -arrestin signaling. *Molecular Systems Biology*, 8(590), 2012.

<sup>0</sup>Pauline Traynard, Céline Feillet, Sylvain Soliman, Franck Delaunay, François Fages. Model-based Investigation of the Circadian Clock and Cell Cycle Coupling in Mouse Embryonic Fibroblasts: Prediction of RevErb-alpha Up-Regulation during Mitosis. *Biosystems*, 149:59–69, 2016.

## **M3DISIM Project-Team**

### **3. Research Program**

#### **3.1. Multi-scale modeling and coupling mechanisms for biomechanical systems, with mathematical and numerical analysis**

Over the past decade, we have laid out the foundations of a multi-scale 3D model of the cardiac mechanical contraction responding to electrical activation. Several collaborations have been crucial in this enterprise, see below references. By integrating this formulation with adapted numerical methods, we are now able to represent the whole organ behavior in interaction with the blood during complete heart beats. This subject was our first achievement to combine a deep understanding of the underlying physics and physiology and our constant concern of proposing well-posed mathematical formulations and adequate numerical discretizations. In fact, we have shown that our model satisfies the essential thermo-mechanical laws, and in particular the energy balance, and proposed compatible numerical schemes that – in consequence – can be rigorously analyzed, see [6]. In the same spirit, we have formulated a poromechanical model adapted to the blood perfusion in the heart, hence precisely taking into account the large deformation of the mechanical medium, the fluid inertia and moving domain, and so that the energy balance between fluid and solid is fulfilled from the model construction to its discretization, see [7].

#### **3.2. Inverse problems with actual data – Fundamental formulation, mathematical analysis and applications**

A major challenge in the context of biomechanical modeling – and more generally in modeling for life sciences – lies in using the large amount of data available on the system to circumvent the lack of absolute modeling ground truth, since every system considered is in fact patient-specific, with possibly non-standard conditions associated with a disease. We have already developed original strategies for solving this particular type of inverse problems by adopting the observer stand-point. The idea we proposed consists in incorporating to the classical discretization of the mechanical system an estimator filter that can use the data to improve the quality of the global approximation, and concurrently identify some uncertain parameters possibly related to a diseased state of the patient. Therefore, our strategy leads to a coupled model-data system solved similarly to a usual PDE-based model, with a computational cost directly comparable to classical Galerkin approximations. We have already worked on the formulation, the mathematical and numerical analysis of the resulting system – see [5] – and the demonstration of the capabilities of this approach in the context of identification of constitutive parameters for a heart model with real data, including medical imaging, see [3].

## OPIS Project-Team

### 3. Research Program

#### 3.1. Accelerated algorithms for solving high-dimensional continuous optimization problems

Variational problems requiring the estimation of a huge number of variables have now to be tackled, especially in the field of 3D reconstruction/restoration (e.g.  $\geq 10^9$  variables in 3D imaging). In addition to the curse of dimensionality, another difficulty to overcome is that the cost function usually reads as the sum of several loss/regularization terms, possibly composed with large-size linear operators. These terms can be nonsmooth and/or nonconvex, as they may serve to promote the sparsity of the sought solution in some suitable representation (e.g. a frame) or to fulfill some physical constraints. In such a challenging context, there is a strong need for developing fast parallelized optimization algorithms for which sound theoretical guarantees of convergence can be established. We explore deterministic and stochastic approaches based on proximal tools, MM (Majorization-Minimization) strategies, and trust region methods. Because of the versatility of the methods that will be proposed, a wide range of applications in image recovery are considered: parallel MRI, breast tomosynthesis, 3D ultrasound imaging, and two-photon microscopy. For example, in breast tomosynthesis (collaboration with GE Healthcare), 3D breast images have to be reconstructed from a small number of X-ray projections with limited view angles. Our objective is to facilitate the clinical task by developing advanced reconstruction methods allowing micro-calcifications to be highlighted. In two-photon microscopy (collaboration with XLIM), our objective is to provide effective numerical solutions to improve the 3D resolution of the microscope, especially when cheap laser sources are used, with applications to muscle disease screening.

#### 3.2. Optimization over graphs

Graphs and hypergraphs are rich data structures for capturing complex, possibly irregular, dependencies in multidimensional data. Coupled with Markov models, they constitute the backbones of many techniques used in computer vision. Optimization is omnipresent in graph processing. Firstly, it allows the structure of the underlying graph to be inferred from the observed data, when the former is hidden. Second, it permits to develop graphical models based on the prior definition of a meaningful cost function. This leads to powerful nonlinear estimates of variables corresponding to unknown weights on the vertices and/or the edges of the graph. Tasks such as partitioning the graph into subgraphs corresponding to different clusters (e.g., communities in social networks) or graph matching, can effectively be performed within this framework. Finally, graphs by themselves offer flexible structures for formulating and solving optimization problems in an efficient distributed manner. On all these topics, our group has acquired a long-term expertise that we plan to further strengthen. In terms of applications, novel graph mining methods are proposed for gene regulatory and brain network analysis. For example, we plan to develop sophisticated methods for better understanding the gene regulatory network of various microscopic fungi, in order to improve the efficiency of the production of bio-fuels (collaboration with IFP Energies Nouvelles).

#### 3.3. Toward more understandable deep learning

Nowadays, deep learning techniques efficiently solve supervised tasks in classification or regression by utilizing large amounts of labeled data and the powerful high level features that they learn by using the input data. Their good performance has caught the attention of the optimization community since currently these methods offer virtually no guarantee of convergence, stability or generalization. Deep neural networks are optimized through a computationally intensive engineering process via methods based on stochastic gradient descent. These methods are slow and they may not lead to relevant local minima. Thus, more efforts must

be dedicated in order to improve the training of deep neural networks by proposing better optimization algorithms applicable to large-scale datasets. Beyond optimization, incorporating some structure in deep neural networks permits more advanced regularization than the current methods. This should reduce their complexity, as well as allow us to derive some bounds regarding generalization. For example, many signal processing models (e.g. those based on multiscale decompositions) exhibit some strong correspondence with deep learning architectures, yet they do not require as many parameters. One can thus think of introducing some supervision into these models in order to improve their performance on standard benchmarks. A better mathematical understanding of these methods permits to improve them, but also to propose some new models and representations for high-dimensional data. This is particularly interesting in settings such as the diagnosis or prevention of diseases from medical images, because they correspond to critical applications where the made decision is crucial and needs to be interpretable. One of the main applications of this work is to propose robust models for the prediction of the outcome of cancer immunotherapy treatments from multiple and complementary sources of information: images, gene expression data, patient profile, etc (collaboration with Institut Gustave Roussy).

## PARIETAL Project-Team

### 3. Research Program

#### 3.1. Inverse problems in Neuroimaging

Many problems in neuroimaging can be framed as forward and inverse problems. For instance, brain population imaging is concerned with the *inverse problem* that consists in predicting individual information (behavior, phenotype) from neuroimaging data, while the corresponding *forward problem* boils down to explaining neuroimaging data with the behavioral variables. Solving these problems entails the definition of two terms: a loss that quantifies the goodness of fit of the solution (does the model explain the data well enough?), and a regularization scheme that represents a prior on the expected solution of the problem. These priors can be used to enforce some properties on the solutions, such as sparsity, smoothness or being piece-wise constant.

Let us detail the model used in typical inverse problem: Let  $\mathbf{X}$  be a neuroimaging dataset as an  $(n_{subjects}, n_{voxels})$  matrix, where  $n_{subjects}$  and  $n_{voxels}$  are the number of subjects under study, and the image size respectively,  $\mathbf{Y}$  a set of values that represent characteristics of interest in the observed population, written as  $(n_{subjects}, n_{features})$  matrix, where  $n_{features}$  is the number of characteristics that are tested, and  $\mathbf{w}$  an array of shape  $(n_{voxels}, n_{features})$  that represents a set of pattern-specific maps. In the first place, we may consider the columns  $\mathbf{Y}_1, \dots, \mathbf{Y}_{n_{features}}$  of  $\mathbf{Y}$  independently, yielding  $n_{features}$  problems to be solved in parallel:

$$\mathbf{Y}_i = \mathbf{X}\mathbf{w}_i + \epsilon_i, \forall i \in \{1, \dots, n_{features}\},$$

where the vector contains  $\mathbf{w}_i$  is the  $i^{th}$  row of  $\mathbf{w}$ . As the problem is clearly ill-posed, it is naturally handled in a regularized regression framework:

$$\hat{\mathbf{w}}_i = \operatorname{argmin}_{\mathbf{w}_i} \|\mathbf{Y}_i - \mathbf{X}\mathbf{w}_i\|^2 + \Psi(\mathbf{w}_i), \quad (2)$$

where  $\Psi$  is an adequate penalization used to regularize the solution:

$$\Psi(\mathbf{w}; \lambda_1, \lambda_2, \eta_1, \eta_2) = \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|_2 + \eta_1 \|\nabla \mathbf{w}\|_{2,1} + \eta_2 \|\nabla \mathbf{w}\|_{2,2} \quad (3)$$

with  $\lambda_1, \lambda_2, \eta_1, \eta_2 \geq 0$  (this formulation particularly highlights the fact that convex regularizers are norms or quasi-norms). In general, only one or two of these constraints is considered (hence is enforced with a non-zero coefficient):

- When  $\lambda_1 > 0$  only (LASSO), and to some extent, when  $\lambda_1, \lambda_2 > 0$  only (elastic net), the optimal solution  $\mathbf{w}$  is (possibly very) sparse, but may not exhibit a proper image structure; it does not fit well with the intuitive concept of a brain map.
- Total Variation regularization (see Fig. 1) is obtained for  $(\eta_1 > 0)$  only, and typically yields a piece-wise constant solution. It can be associated with Lasso to enforce both sparsity and sparse variations.
- Smooth lasso is obtained with  $(\eta_2 > 0)$  and  $\lambda_1 > 0$  only, and yields smooth, compactly supported spatial basis functions.

Note that, while the qualitative aspect of the solutions are very different, the predictive power of these models is often very close.



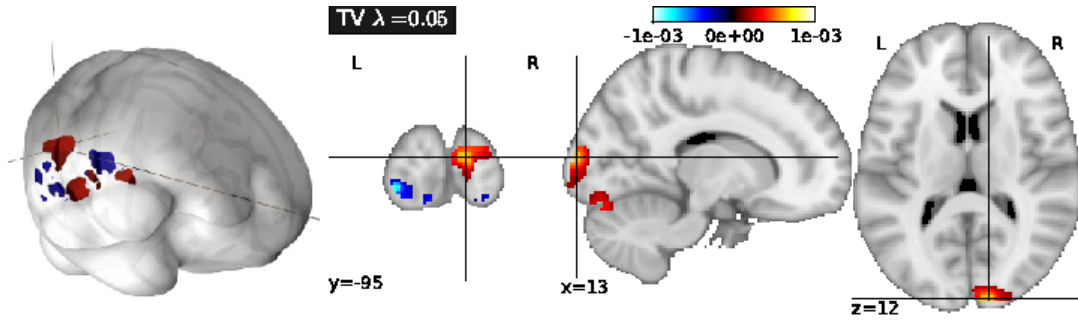


Figure 1. Example of the regularization of a brain map with total variation in an inverse problem. The problem here is to predict the spatial scale of an object presented as a stimulus, given functional neuroimaging data acquired during the presentation of an image. Learning and test are performed across individuals. Unlike other approaches, Total Variation regularization yields a sparse and well-localized solution that also enjoys high predictive accuracy.

The performance of the predictive model can simply be evaluated as the amount of variance in  $\mathbf{Y}_i$  fitted by the model, for each  $i \in \{1, \dots, n_{features}\}$ . This can be computed through cross-validation, by *learning*  $\hat{\mathbf{w}}_i$  on some part of the dataset, and then estimating  $\|\mathbf{Y}_i - \mathbf{X}\hat{\mathbf{w}}_i\|^2$  using the remainder of the dataset.

This framework is easily extended by considering

- *Grouped penalization*, where the penalization explicitly includes a prior clustering of the features, i.e. voxel-related signals, into given groups. This amounts to enforcing structured priors on the solution.
- *Combined penalizations*, i.e. a mixture of simple and group-wise penalizations, that allow some variability to fit the data in different populations of subjects, while keeping some common constraints.
- *Logistic and hinge regression*, where a non-linearity is applied to the linear model so that it yields a probability of classification in a binary classification problem.
- *Robustness to between-subject variability* to avoid the learned model overly reflecting a few outlying particular observations of the training set. Note that noise and deviating assumptions can be present in both  $\mathbf{Y}$  and  $\mathbf{X}$
- *Multi-task learning*: if several target variables are thought to be related, it might be useful to constrain the estimated parameter vector  $\mathbf{w}$  to have a shared support across all these variables. For instance, when one of the variables  $\mathbf{Y}_i$  is not well fitted by the model, the estimation of other variables  $\mathbf{Y}_j, j \neq i$  may provide constraints on the support of  $\mathbf{w}_i$  and thus, improve the prediction of  $\mathbf{Y}_i$ .

$$\mathbf{Y} = \mathbf{X}\mathbf{w} + \epsilon, \quad (4)$$

then

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}=(\mathbf{w}_i), i=1..n_f} \sum_{i=1}^{n_f} \|\mathbf{Y}_i - \mathbf{X}\mathbf{w}_i\|^2 + \lambda \sum_{j=1}^{n_{voxels}} \sqrt{\sum_{i=1}^{n_f} \mathbf{w}_{i,j}^2} \quad (5)$$

### 3.2. Multivariate decompositions

Multivariate decompositions provide a way to model complex data such as brain activation images: for instance, one might be interested in extracting an *atlas of brain regions* from a given dataset, such as regions exhibiting similar activity during a protocol, across multiple protocols, or even in the absence of protocol (during resting-state). These data can often be factorized into spatial-temporal components, and thus can be estimated through *regularized Principal Components Analysis* (PCA) algorithms, which share some common steps with regularized regression.

Let  $\mathbf{X}$  be a neuroimaging dataset written as an  $(n_{subjects}, n_{voxels})$  matrix, after proper centering; the model reads

$$\mathbf{X} = \mathbf{A}\mathbf{D} + \epsilon, \quad (6)$$

where  $\mathbf{D}$  represents a set of  $n_{comp}$  spatial maps, hence a matrix of shape  $(n_{comp}, n_{voxels})$ , and  $\mathbf{A}$  the associated subject-wise loadings. While traditional PCA and independent components analysis (ICA) are limited to reconstructing components  $\mathbf{D}$  within the space spanned by the column of  $\mathbf{X}$ , it seems desirable to add some constraints on the rows of  $\mathbf{D}$ , that represent spatial maps, such as sparsity, and/or smoothness, as it makes the interpretation of these maps clearer in the context of neuroimaging. This yields the following estimation problem:

$$\min_{\mathbf{D}, \mathbf{A}} \|\mathbf{X} - \mathbf{A}\mathbf{D}\|^2 + \Psi(\mathbf{D}) \quad \text{s.t.} \quad \|\mathbf{A}_i\| = 1 \quad \forall i \in \{1..n_{features}\}, \quad (7)$$

where  $(\mathbf{A}_i)$ ,  $i \in \{1..n_{features}\}$  represents the columns of  $\mathbf{A}$ .  $\Psi$  can be chosen such as in Eq. (2) in order to enforce smoothness and/or sparsity constraints.

The problem is not jointly convex in all the variables but each penalization given in Eq (2) yields a convex problem on  $\mathbf{D}$  for  $\mathbf{A}$  fixed, and conversely. This readily suggests an alternate optimization scheme, where  $\mathbf{D}$  and  $\mathbf{A}$  are estimated in turn, until convergence to a local optimum of the criterion. As in PCA, the extracted components can be ranked according to the amount of fitted variance. Importantly, also, estimated PCA models can be interpreted as a probabilistic model of the data, assuming a high-dimensional Gaussian distribution (probabilistic PCA).

Ultimately, the main limitations to these algorithms is the cost due to the memory requirements: holding datasets with large dimension and large number of samples (as in recent neuroimaging cohorts) leads to inefficient computation. To solve this issue, online methods are particularly attractive [1].

### 3.3. Covariance estimation

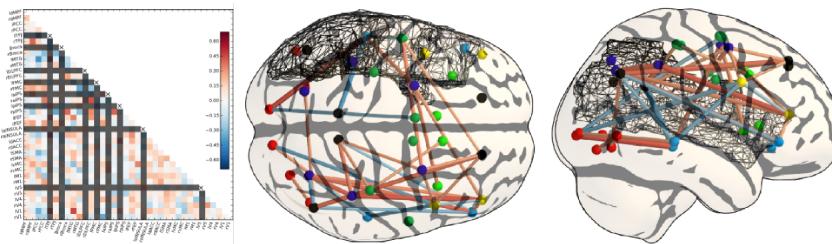
Another important estimation problem stems from the general issue of learning the relationship between sets of variables, in particular their covariance. Covariance learning is essential to model the dependence of these variables when they are used in a multivariate model, for instance to study potential interactions among them and with other variables. Covariance learning is necessary to model latent interactions in high-dimensional observation spaces, e.g. when considering multiple contrasts or functional connectivity data.

The difficulties are two-fold: on the one hand, there is a shortage of data to learn a good covariance model from an individual subject, and on the other hand, subject-to-subject variability poses a serious challenge to the use of multi-subject data. While the covariance structure may vary from population to population, or depending on the input data (activation versus spontaneous activity), assuming some shared structure across problems, such as their sparsity pattern, is important in order to obtain correct estimates from noisy data. Some of the most important models are:

- **Sparse Gaussian graphical models**, as they express meaningful conditional independence relationships between regions, and do improve conditioning/avoid overfit.

- **Decomposable models**, as they enjoy good computational properties and enable intuitive interpretations of the network structure. Whether they can faithfully or not represent brain networks is still an open question.
- **PCA-based regularization of covariance** which is powerful when modes of variation are more important than conditional independence relationships.

Adequate model selection procedures are necessary to achieve the right level of sparsity or regularization in covariance estimation; the natural evaluation metric here is the out-of-sample likelihood of the associated Gaussian model. Another essential remaining issue is to develop an adequate statistical framework to test differences between covariance models in different populations. To do so, we consider different means of parametrizing covariance distributions and how these parametrizations impact the test of statistical differences across individuals.



*Figure 2. Example of functional connectivity analysis: The correlation matrix describing brain functional connectivity in a post-stroke patient (lesion volume outlined as a mesh) is compared to a group of control subjects. Some edges of the graphical model show a significant difference, but the statistical detection of the difference requires a sophisticated statistical framework for the comparison of graphical models.*

## XPOP Project-Team

### 3. Research Program

#### 3.1. Scientific positioning

"Interfaces" is the defining characteristic of XPOP:

**The interface between statistics, probability and numerical methods.** Mathematical modelling of complex biological phenomena require to combine numerical, stochastic and statistical approaches. The CMAP is therefore the right place to be for positioning the team at the interface between several mathematical disciplines.

**The interface between mathematics and the life sciences.** The goal of XPOP is to bring the right answers to the right questions. These answers are mathematical tools (statistics, numerical methods, etc.), whereas the questions come from the life sciences (pharmacology, medicine, biology, etc.). This is why the point of XPOP is not to take part in mathematical projects only, but also pluridisciplinary ones.

**The interface between mathematics and software development.** The development of new methods is the main activity of XPOP. However, new methods are only useful if they end up being implemented in a software tool. On one hand, a strong partnership with Lixoft (the spin-off company who continue developing MONOLIX) allows us to maintaining this positioning. On the other hand, several members of the team are very active in the R community and develop widely used packages.

#### 3.2. The mixed-effects models

Mixed-effects models are statistical models with both fixed effects and random effects. They are well-adapted to situations where repeated measurements are made on the same individual/statistical unit.

Consider first a single subject  $i$  of the population. Let  $y_i = (y_{ij}, 1 \leq j \leq n_i)$  be the vector of observations for this subject. The model that describes the observations  $y_i$  is assumed to be a parametric probabilistic model: let  $p_Y(y_i; \psi_i)$  be the probability distribution of  $y_i$ , where  $\psi_i$  is a vector of parameters.

In a population framework, the vector of parameters  $\psi_i$  is assumed to be drawn from a population distribution  $p_\Psi(\psi_i; \theta)$  where  $\theta$  is a vector of population parameters.

Then, the probabilistic model is the joint probability distribution

$$p(y_i, \psi_i; \theta) = p_Y(y_i | \psi_i) p_\Psi(\psi_i; \theta) \quad (8)$$

To define a model thus consists in defining precisely these two terms.

In most applications, the observed data  $y_i$  are continuous longitudinal data. We then assume the following representation for  $y_i$ :

$$y_{ij} = f(t_{ij}, \psi_i) + g(t_{ij}, \psi_i) \varepsilon_{ij} \quad , \quad 1 \leq i \leq N \quad , \quad 1 \leq j \leq n_i. \quad (9)$$

Here,  $y_{ij}$  is the observation obtained from subject  $i$  at time  $t_{ij}$ . The residual errors ( $\varepsilon_{ij}$ ) are assumed to be standardized random variables (mean zero and variance 1). The residual error model is represented by function  $g$  in model (2).

Function  $f$  is usually the solution to a system of ordinary differential equations (pharmacokinetic/pharmacodynamic models, etc.) or a system of partial differential equations (tumor growth, respiratory system, etc.). This component is a fundamental component of the model since it defines the prediction of the observed kinetics for a given set of parameters.

The vector of individual parameters  $\psi_i$  is usually function of a vector of population parameters  $\psi_{\text{pop}}$ , a vector of random effects  $\eta_i \sim \mathcal{N}(0, \Omega)$ , a vector of individual covariates  $c_i$  (weight, age, gender, ...) and some fixed effects  $\beta$ .

The joint model of  $y$  and  $\psi$  depends then on a vector of parameters  $\theta = (\psi_{\text{pop}}, \beta, \Omega)$ .

### 3.3. Computational Statistical Methods

Central to modern statistics is the use of probabilistic models. To relate these models to data requires the ability to calculate the probability of the observed data: the likelihood function, which is central to most statistical methods and provides a principled framework to handle uncertainty.

The emergence of computational statistics as a collection of powerful and general methodologies for carrying out likelihood-based inference made complex models with non-standard data accessible to likelihood, including hierarchical models, models with intricate latent structure, and missing data.

In particular, algorithms previously developed by POPIX for mixed effects models, and today implemented in several software tools (especially MONOLIX) are part of these methods:

- the adaptive Metropolis-Hastings algorithm allows one to sample from the conditional distribution of the individual parameters  $p(\psi_i | y_i; c_i, \theta)$ ,
- the SAEM algorithm is used to maximize the observed likelihood  $\mathcal{L}(\theta; y) = p(y; \theta)$ ,
- Importance Sampling Monte Carlo simulations provide an accurate estimation of the observed log-likelihood  $\log(\mathcal{L}(\theta; y))$ .

Computational statistics is an area which remains extremely active today. Recently, one can notice that the incentive for further improvements and innovation comes mainly from three broad directions: the high dimensional challenge, the quest for adaptive procedures that can eliminate the cumbersome process of tuning "by hand" the settings of the algorithms and the need for flexible theoretical support, arguably required by all recent developments as well as many of the traditional MCMC algorithms that are widely used in practice.

Working in these three directions is a clear objective for XPOP.

### 3.4. Markov Chain Monte Carlo algorithms

While these Monte Carlo algorithms have turned into standard tools over the past decade, they still face difficulties in handling less regular problems such as those involved in deriving inference for high-dimensional models. One of the main problems encountered when using MCMC in this challenging settings is that it is difficult to design a Markov chain that efficiently samples the state space of interest.

The Metropolis-adjusted Langevin algorithm (MALA) is a Markov chain Monte Carlo (MCMC) method for obtaining random samples from a probability distribution for which direct sampling is difficult. As the name suggests, MALA uses a combination of two mechanisms to generate the states of a random walk that has the target probability distribution as an invariant measure:

1. new states are proposed using Langevin dynamics, which use evaluations of the gradient of the target probability density function;
2. these proposals are accepted or rejected using the Metropolis-Hastings algorithm, which uses evaluations of the target probability density (but not its gradient).

Informally, the Langevin dynamics drives the random walk towards regions of high probability in the manner of a gradient flow, while the Metropolis-Hastings accept/reject mechanism improves the mixing and convergence properties of this random walk.

Several extensions of MALA have been proposed recently by several authors, including fMALA (fast MALA), AMALA (anisotropic MALA), MMALA (manifold MALA), position-dependent MALA (PMALA), ...

MALA and these extensions have demonstrated to represent very efficient alternative for sampling from high dimensional distributions. We therefore need to adapt these methods to general mixed effects models.

### 3.5. Parameter estimation

The Stochastic Approximation Expectation Maximization (SAEM) algorithm has shown to be extremely efficient for maximum likelihood estimation in incomplete data models, and particularly in mixed effects models for estimating the population parameters. However, there are several practical situations for which extensions of SAEM are still needed:

**High dimensional model:** a complex physiological model may have a large number of parameters (in the order of 100). Then several problems arise:

- when most of these parameters are associated with random effects, the MCMC algorithm should be able to sample, for each of the  $N$  individuals, parameters from a high dimensional distribution. Efficient MCMC methods for high dimensions are then required.
- Practical identifiability of the model is not ensured with a limited amount of data. In other words, we cannot expect to be able to properly estimate all the parameters of the model, including the fixed effects and the variance-covariance matrix of the random effects. Then, some random effects should be removed, assuming that some parameters do not vary in the population. It may also be necessary to fix the value of some parameters (using values from the literature for instance). The strategy to decide which parameters should be fixed and which random effects should be removed remains totally empirical. XPOP aims to develop a procedure that will help the modeller to take such decisions.

**Large number of covariates:** the covariate model aims to explain part of the inter-patient variability of some parameters. Classical methods for covariate model building are based on comparisons with respect to some criteria, usually derived from the likelihood (AIC, BIC), or some statistical test (Wald test, LRT, etc.). In other words, the modelling procedure requires two steps: first, all possible models are fitted using some estimation procedure (e.g. the SAEM algorithm) and the likelihood of each model is computed using a numerical integration procedure (e.g. Monte Carlo Importance Sampling); then, a model selection procedure chooses the "best" covariate model. Such a strategy is only possible with a reduced number of covariates, i.e., with a "small" number of models to fit and compare.

As an alternative, we are thinking about a Bayesian approach which consists of estimating simultaneously the covariate model and the parameters of the model in a single run. An (informative or uninformative) prior is defined for each model by defining a prior probability for each covariate to be included in the model. In other words, we extend the probabilistic model by introducing binary variables that indicate the presence or absence of each covariate in the model. Then, the model selection procedure consists of estimating and maximizing the conditional distribution of this sequence of binary variables. Furthermore, a probability can be associated to any of the possible covariate models.

This conditional distribution can be estimated using an MCMC procedure combined with the SAEM algorithm for estimating the population parameters of the model. In practice, such an approach can only deal with a limited number of covariates since the dimension of the probability space to explore increases exponentially with the number of covariates. Consequently, we would like to have methods able to find a small number of variables (from a large starting set) that influence certain parameters in populations of individuals. That means that, instead of estimating the conditional distribution of all the covariate models as described above, the algorithm should focus on the most likely ones.

**Fixed parameters:** it is quite frequent that some individual parameters of the model have no random component and are purely fixed effects. Then, the model may not belong to the exponential family anymore and the original version of SAEM cannot be used as it is. Several extensions exist:

- introduce random effects with decreasing variances for these parameters,
- introduce a prior distribution for these fixed effects,
- apply the stochastic approximation directly on the sequence of estimated parameters, instead of the sufficient statistics of the model.

None of these methods always work correctly. Furthermore, what are the pros and cons of these methods is not clear at all. Then, developing a robust methodology for such model is necessary.

**Convergence toward the global maximum of the likelihood:** convergence of SAEM can strongly depend on the initial guess when the observed likelihood has several local maxima. A kind of simulated annealing version of SAEM was previously developed and implemented in MONOLIX. The method works quite well in most situations but there is no theoretical justification and choosing the settings of this algorithm (i.e. how the temperature decreases during the iterations) remains empirical. A precise analysis of the algorithm could be very useful to better understand why it "works" in practice and how to optimize it.

**Convergence diagnostic:** Convergence of SAEM was theoretically demonstrated under very general hypothesis. Such result is important but of little interest in practice at the time to use SAEM in a finite amount of time, i.e. in a finite number of iterations. Some qualitative and quantitative criteria should be defined in order to both optimize the settings of the algorithm, detect a poor convergence of SAEM and evaluate the quality of the results in order to avoid using them unwisely.

### 3.6. Model building

Defining an optimal strategy for model building is far from easy because a model is the assembled product of numerous components that need to be evaluated and perhaps improved: the structural model, residual error model, covariate model, covariance model, etc.

How to proceed so as to obtain the best possible combination of these components? There is no magic recipe but an effort will be made to provide some qualitative and quantitative criteria in order to help the modeller for building his model.

The strategy to take will mainly depend on the time we can dedicate to building the model and the time required for running it. For relatively simple models for which parameter estimation is fast, it is possible to fit many models and compare them. This can also be done if we have powerful computing facilities available (e.g., a cluster) allowing large numbers of simultaneous runs.

However, if we are working on a standard laptop or desktop computer, model building is a sequential process in which a new model is tested at each step. If the model is complex and requires significant computation time (e.g., when involving systems of ODEs), we are constrained to limit the number of models we can test in a reasonable time period. In this context, it also becomes important to carefully choose the tasks to run at each step.

### 3.7. Model evaluation

Diagnostic tools are recognized as an essential method for model assessment in the process of model building. Indeed, the modeler needs to confront "his" model with the experimental data before concluding that this model is able to reproduce the data and before using it for any purpose, such as prediction or simulation for instance.

The objective of a diagnostic tool is twofold: first we want to check if the assumptions made on the model are valid or not ; then, if some assumptions are rejected, we want to get some guidance on how to improve the model.

As is the usual case in statistics, it is not because this "final" model has not been rejected that it is necessarily the "true" one. All that we can say is that the experimental data does not allow us to reject it. It is merely one of perhaps many models that cannot be rejected.

Model diagnostic tools are for the most part graphical, i.e., visual; we "see" when something is not right between a chosen model and the data it is hypothesized to describe. These diagnostic plots are usually based on the empirical Bayes estimates (EBEs) of the individual parameters and EBEs of the random effects: scatterplots of individual parameters versus covariates to detect some possible relationship, scatterplots of pairs of random effects to detect some possible correlation between random effects, plot of the empirical distribution of the random effects (boxplot, histogram,...) to check if they are normally distributed, ...

The use of EBEs for diagnostic plots and statistical tests is efficient with rich data, i.e. when a significant amount of information is available in the data for recovering accurately all the individual parameters. On the contrary, tests and plots can be misleading when the estimates of the individual parameters are greatly shrunk.

We propose to develop new approaches for diagnosing mixed effects models in a general context and derive formal and unbiased statistical tests for testing separately each feature of the model.

### **3.8. Missing data**

The ability to easily collect and gather a large amount of data from different sources can be seen as an opportunity to better understand many processes. It has already led to breakthroughs in several application areas. However, due to the wide heterogeneity of measurements and objectives, these large databases often exhibit an extraordinary high number of missing values. Hence, in addition to scientific questions, such data also present some important methodological and technical challenges for data analyst.

Missing values occur for a variety of reasons: machines that fail, survey participants who do not answer certain questions, destroyed or lost data, dead animals, damaged plants, etc. Missing values are problematic since most statistical methods can not be applied directly on a incomplete data. Many progress have been made to properly handle missing values. However, there are still many challenges that need to be addressed in the future, that are crucial for the users.

- State of arts methods often consider the case of continuous or categorical data whereas real data are very often mixed. The idea is to develop a multiple imputation method based on a specific principal component analysis (PCA) for mixed data. Indeed, PCA has been used with success to predict (impute) the missing values. A very appealing property is the ability of the method to handle very large matrices with large amount of missing entries.
- The asymptotic regime underlying modern data is not any more to consider that the sample size increases but that both number of observations and number of variables are very large. In practice first experiments showed that the coverage properties of confidence areas based on the classical methods to estimate variance with missing values varied widely. The asymptotic method and the bootstrap do well in low-noise setting, but can fail when the noise level gets high or when the number of variables is much greater than the number of rows. On the other hand, the jackknife has good coverage properties for large noisy examples but requires a minimum number of variables to be stable enough.
- Inference with missing values is usually performed under the assumption of "Missing at Random" (MAR) values which means that the probability that a value is missing may depend on the observed data but does not depend on the missing value itself. In real data and in particular in data coming from clinical studies, both "Missing Non at Random" (MNAR) and MAR values occur. Taking into account in a proper way both types of missing values is extremely challenging but is worth investigating since the applications are extremely broad.

It is important to stress that missing data models are part of the general incomplete data models addressed by XPOP. Indeed, models with latent variables (i.e. non observed variables such as random effects in a mixed effects model), models with censored data (e.g. data below some limit of quantification) or models with dropout mechanism (e.g. when a subject in a clinical trial fails to continue in the study) can be seen as missing data models.



## TRIBE Project-Team

### 3. Research Program

#### 3.1. Research program

Following up on the effort initiated by the team members during the last few years and building on an approach combining protocol design, data analytics, and experimental research, we propose a research program organized around three closely related objectives that are briefly described in the following.

- **Technologies for accommodating low-end IoT devices:** The IoT is expected to gradually connect billions of low-end devices to the Internet, and thereby drastically increase communication without human source or destination. Low-end IoT devices differ starkly from high-end IoT devices in terms of resources such as energy, memory, and computational power. Projections show this divide will not fundamentally change in the future and that IoT should ultimately interconnect a dense population of devices as tiny as dust particles, feeding off ambient power sources (energy harvesting). These characteristics constrain the software and communication protocols running on low-end IoT devices: they are neither able to run a common software platform such as Linux (or its derivatives), nor the standard protocol stack based on TCP/IP. Solutions for low-end IoT devices require thus: **(i) optimized communication protocols** taking into account radio technology evolution and devices constrained requirements; **(ii) tailored software platforms** providing high level programming, modular software updates as well as advanced support for new security and energy concentration features; **(iii) unification of technologies** for low-end IoT, which is too fragmented at the moment, guaranteeing integration with core or other edge networks.
- **Technologies for leveraging high-end IoT devices' advents:** High-end IoT devices are one of the most important instances of the connected devices supporting a noteworthy shift towards mobile Internet access. As our lives become more dependent on pervasive connectivity, our social patterns (as human being in the Internet era) are nowadays being reflected from our real life onto the virtual binary world. This gives birth to two tendencies. From one side, edge networks can now be utilized as mirrors to reflect the inherent human dynamics, their context, and interests thanks to their well organized recording and almost ubiquitous coverage. From the other side, social norms and structure dictating human behavior (e.g., interactions, mobility, interest, cultural patterns) are now directly influencing the way individuals interact with the network services and demand resources or content. In particular, we observe the particularities present in human dynamics *shape the way (i.e., where, when, how, or what) resources, services, and infrastructures are used at the edge of the Internet*. Hence, we claim a need to digitally study high-end IoT devices' end-users behaviors and to leverage this understanding in networking solutions' design, so as to optimize network exploitation. This suggests the **integration of the heterogeneity and uncertainty of behaviors in designed networking solutions**. For this, *useful knowledge* allowing the understanding of behaviors and context of users has to be *extracted and delivered out* of large masses of data. Such knowledge has to be then *integrated in current design practices*. This brings the idea of a more *tactful networking design practice* where the network is assigned with the human like capability of observation, interpretation, and reaction to daily life features and entities involving high-end IoT devices. Research activities here include: **(i) the quest for meaningful data**, which includes the integration of data from different sources, the need for scaling up data analysis, the usage and analysis of fine-grained datasets, or still, the completion of sparse and coarse grained datasets; **(ii) expanding edge networks' usage understanding**, which concerns analysis on how and when contextual information impact network usage, fine-grained analysis of short-term mobility of individuals, or the identification of patterns of behavior and novelty-seeking of individuals; **(iii) human-driven prediction models**, extensible to context awareness and adapted to individuals preferences in terms of novelty, diversity, or routines.

- **Articulating the IoT edge with the core of the network:** The edge is the interface between the IoT devices and the core network: some of the challenges encountered by IoT devices have their continuity at the edge of the network inside the gateway (i.e., interoperability, heterogeneity and mobility support). Besides, the edge should be able to support intermediary functions between devices and the rest of the core (e.g., the cloud). This includes: **(i) proxying functionality**, facilitating connections between devices and the Internet; **(ii) machine learning enhanced IoT solutions**, designed to improve performance of advanced IoT networked systems (e.g., through methods such as supervised, unsupervised or reinforcement learning) at adapted levels of the protocol stack (e.g., for multiple access, coding, choices); **(iii) IoT data contextualization**, so collection of meaningful IoT data (i.e., right data collected at the right time) can be earlier determined closer to the data source; **(iv) intermediary computation** through fog or Mobile Edge Computing (MEC) models, where IoT devices can obtain computing, data storage, and communication means with lower latency in a decentralized way; or **(v) security of end-to-end IoT software supply-chain**, including remote management and over-the-air updates.

## AVIZ Project-Team

### 3. Research Program

#### 3.1. Scientific Foundations

The scientific foundations of Visual Analytics lie primarily in the domains of Visualization and Data Mining. Indirectly, it inherits from other established domains such as graphic design, Exploratory Data Analysis (EDA), statistics, Artificial Intelligence (AI), Human-Computer Interaction (HCI), and Psychology.

The use of graphic representation to understand abstract data is a goal Visual Analytics shares with Tukey's Exploratory Data Analysis (EDA) [67], graphic designers such as Bertin [54] and Tufte [66], and HCI researchers in the field of Information Visualization [53].

EDA is complementary to classical statistical analysis. Classical statistics starts from a *problem*, gathers *data*, designs a *model* and performs an *analysis* to reach a *conclusion* about whether the data follows the model. While EDA also starts with a problem and data, it is most useful *before* we have a model; rather, we perform visual analysis to discover what kind of model might apply to it. However, statistical validation is not always required with EDA; since often the results of visual analysis are sufficiently clear-cut that statistics are unnecessary.

Visual Analytics relies on a process similar to EDA, but expands its scope to include more sophisticated graphics and areas where considerable automated analysis is required before the visual analysis takes place. This richer data analysis has its roots in the domain of Data Mining, while the advanced graphics and interactive exploration techniques come from the scientific fields of Data Visualization and HCI, as well as the expertise of professions such as cartography and graphic designers who have long worked to create effective methods for graphically conveying information.

The books of the cartographer Bertin and the graphic designer Tufte are full of rules drawn from their experience about how the meaning of data can be best conveyed visually. Their purpose is to find effective visual representation that describe a data set but also (mainly for Bertin) to discover structure in the data by using the right mappings from abstract dimensions in the data to visual ones.

For the last 25 years, the field of Human-Computer Interaction (HCI) has also shown that interacting with visual representations of data in a tight perception-action loop improves the time and level of understanding of data sets. Information Visualization is the branch of HCI that has studied visual representations suitable to understanding and interaction methods suitable to navigating and drilling down on data. The scientific foundations of Information Visualization come from theories about perception, action and interaction.

Several theories of perception are related to information visualization such as the "Gestalt" principles, Gibson's theory of visual perception [59] and Triesman's "preattentive processing" theory [65]. We use them extensively but they only have a limited accuracy for predicting the effectiveness of novel visual representations in interactive settings.

Information Visualization emerged from HCI when researchers realized that interaction greatly enhanced the perception of visual representations.

To be effective, interaction should take place in an interactive loop faster than 100ms. For small data sets, it is not difficult to guarantee that analysis, visualization and interaction steps occur in this time, permitting smooth data analysis and navigation. For larger data sets, more computation should be performed to reduce the data size to a size that may be visualized effectively.

In 2002, we showed that the practical limit of InfoVis was on the order of 1 million items displayed on a screen [57]. Although screen technologies have improved rapidly since then, eventually we will be limited by the physiology of our vision system: about 20 millions receptor cells (rods and cones) on the retina. Another problem will be the limits of human visual attention, as suggested by our 2006 study on change blindness in large and multiple displays [55]. Therefore, visualization alone cannot let us understand very large data sets. Other techniques such as aggregation or sampling must be used to reduce the visual complexity of the data to the scale of human perception.

Abstracting data to reduce its size to what humans can understand is the goal of Data Mining research. It uses data analysis and machine learning techniques. The scientific foundations of these techniques revolve around the idea of finding a good model for the data. Unfortunately, the more sophisticated techniques for finding models are complex, and the algorithms can take a long time to run, making them unsuitable for an interactive environment. Furthermore, some models are too complex for humans to understand; so the results of data mining can be difficult or impossible to understand directly.

Unlike pure Data Mining systems, a Visual Analytics system provides analysis algorithms and processes compatible with human perception and understandable to human cognition. The analysis should provide understandable results quickly, even if they are not ideal. Instead of running to a predefined threshold, algorithms and programs should be designed to allow trading speed for quality and show the tradeoffs interactively. This is not a temporary requirement: it will be with us even when computers are much faster, because good quality algorithms are at least quadratic in time (e.g. hierarchical clustering methods). Visual Analytics systems need different algorithms for different phases of the work that can trade speed for quality in an understandable way.

Designing novel interaction and visualization techniques to explore huge data sets is an important goal and requires solving hard problems, but how can we assess whether or not our techniques and systems provide real improvements? Without this answer, we cannot know if we are heading in the right direction. This is why we have been actively involved in the design of evaluation methods for information visualization [63], [62], [60], [61], [58]. For more complex systems, other methods are required. For these we want to focus on longitudinal evaluation methods while still trying to improve controlled experiments.

## 3.2. Innovation

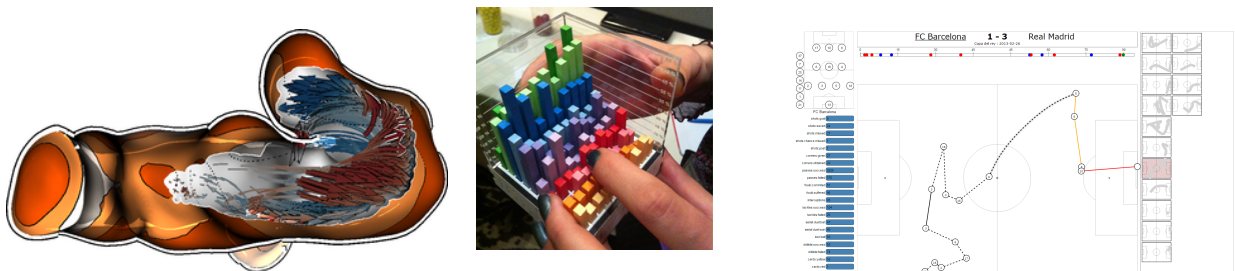


Figure 1. Example novel visualization techniques and tools developed by the team. Left: a non-photorealistic rendering technique that visualizes blood flow and vessel thickness. Middle: a physical visualization showing economic indicators for several countries, right: SoccerStories a tool for visualizing soccer games.

We design novel visualization and interaction techniques (see, for example, Figure 1 ). Many of these techniques are also evaluated throughout the course of their respective research projects. We cover application domains such as sports analysis, digital humanities, fluid simulations, and biology. A focus of Aviz' work is the improvement of graph visualization and interaction with graphs. We further develop individual techniques

for the design of tabular visualizations and different types of data charts. Another focus is the use of animation as a transition aid between different views of the data. We are also interested in applying techniques from illustrative visualization to visual representations and applications in information visualization as well as scientific visualization [8].

### 3.3. Evaluation Methods

Evaluation methods are required to assess the effectiveness and usability of visualization and analysis methods. Aviz typically uses traditional HCI evaluation methods, either quantitative (measuring speed and errors) or qualitative (understanding users tasks and activities). Moreover, Aviz is also contributing to the improvement of evaluation methods by reporting on the best practices in the field, by co-organizing workshops (BELIV 2010–2018) to exchange on novel evaluation methods, by improving our ways of reporting, interpreting and communicating statistical results, and by applying novel methodologies, for example to assess visualization literacy [3], [4].

### 3.4. Software Infrastructures

We want to understand the requirements that software and hardware architectures should provide to support exploratory analysis of large amounts of data. So far, “big data” has been focusing on issues related to storage management and predictive analysis: applying a well-known set of operations on large amounts of data. Visual Analytics is about exploration of data, with sometimes little knowledge of its structure or properties. Therefore, interactive exploration and analysis is needed to build knowledge and apply appropriate analyses; this knowledge and appropriateness is supported by visualizations. However, applying analytical operations on large data implies long-lasting computations, incompatible with interactions, and generates large amounts of results, impossible to visualize directly without aggregation or sampling. Visual Analytics has started to tackle these problems for specific applications but not in a general manner, leading to fragmentation of results and difficulties to reuse techniques from one application to the other. We are interested in abstracting-out the issues and finding general architectural models, patterns, and frameworks to address the Visual Analytics challenge in more generic ways.

### 3.5. Emerging Technologies



Figure 2. Example emerging technology solutions developed by the team for multi-display environments, wall displays, and token-based visualization.

We want to use different types of display media to empower humans to visually and interactively explore information, in order to better understand and exploit it. This includes novel display equipment and accompanying input techniques. The Aviz team specifically focuses on the exploration of the use of large displays in visualization contexts as well as emerging physical and tangible visualizations (e. g. [6], [5]). In terms of interaction modalities our work focuses on using touch and tangible interaction. Aviz participates to the Digiscope project that funds 11 wall-size displays at multiple places in the Paris area (see <http://www.digiscope.fr>),

connected by telepresence equipment and a Fablab for creating devices. Aviz is in charge of creating and managing the Fablab, uses it to create physical visualizations, and is also using the local wall-size display (called WILD) to explore visualization on large screens. The team also investigates the perceptual, motor and cognitive implications of using such technologies for visualization.

### **3.6. Psychology**

More cross-fertilization is needed between psychology and information visualization. The only key difference lies in their ultimate objective: understanding the human mind vs. helping to develop better tools. We focus on understanding and using findings from psychology to inform new tools for information visualization. In many cases, our work also extends previous work in psychology. Our approach to the psychology of information visualization is largely holistic and helps bridge gaps between perception, action and cognition in the context of information visualization. Our focus includes the perception of charts in general, perception in large display environments, collaboration, perception of animations, how action can support perception and cognition, and judgment under uncertainty (e. g. [9]).

## CEDAR Project-Team

### 3. Research Program

#### 3.1. Scalable Heterogeneous Stores

Big Data applications increasingly involve *diverse* data sources, such as: structured or unstructured documents, data graphs, relational databases etc. and it is often impractical to load (consolidate) diverse data sources in a single repository. Instead, interesting data sources need to be exploited “as they are”, with the added value of the data being realized especially through the ability to combine (join) together data from several sources. Systems capable of exploiting diverse Big Data in this fashion are usually termed *polystores*. A current limitation of polystores is that data stays captive of its original storage system, which may limit the data exploitation performance. We work to devise highly efficient storage systems for heterogeneous data across a variety of data stores.

#### 3.2. Semantic Query Answering

In the presence of data semantics, query evaluation techniques are insufficient as they only take into account the database, but do not provide the reasoning capabilities required in order to reflect the semantic knowledge. In contrast, (ontology-based) query answering takes into account both the data and the semantic knowledge in order to compute the full query answers, blending query evaluation and semantic reasoning.

We aim at designing efficient semantic query answering algorithms, both building on cost-based reformulation algorithms developed in the team and exploring new approaches mixing materialization and reformulation.

#### 3.3. Multi-Model Querying

As the world’s affairs get increasingly more digital, a large and varied set of data sources becomes available: they are either structured databases, such as government-gathered data (demographics, economics, taxes, elections, ...), legal records, stock quotes for specific companies, un-structured or semi-structured, including in particular graph data, sometimes endowed with semantics (see e.g. the Linked Open Data cloud). Modern data management applications, such as data journalism, are eager to combine in innovative ways both static and dynamic information coming from structured, semi-structured, and un-structured databases and social feeds. However, current content management tools for this task are not suited for the task, in particular when they require a lengthy rigid cycle of data integration and consolidation in a warehouse. Thus, we see a need for flexible tools allowing to interconnect various kinds of data sources and to query them together.

#### 3.4. Interactive Data Exploration at Scale

In the Big Data era we are faced with an increasing gap between the fast growth of data and the limited human ability to comprehend data. Consequently, there has been a growing demand of data management tools that can bridge this gap and help users retrieve high-value content from data more effectively. To respond to such user information needs, we aim to build interactive data exploration as a new database service, using an approach called “explore-by-example”.

#### 3.5. Exploratory Querying of Semantic Graphs

Semantic graphs including data and knowledge are hard to apprehend for users, due to the complexity of their structure and oftentimes to their large volumes. To help tame this complexity, in prior research (2014), we have presented a full framework for RDF data warehousing, specifically designed for heterogeneous and semantic-rich graphs. However, this framework still leaves to the users the burden of choosing the most interesting warehousing queries to ask. More user-friendly data management tools are needed, which help the user discover the interesting structure and information hidden within RDF graphs. This research has benefitted from the arrival in the team of Mirjana Mazuran, as well as from the start of the PhD thesis of Pawel Guzewicz, co-advised by Yanlei Diao and Ioana Manolescu.

### **3.6. An Unified Framework for Optimizing Data Analytics**

Data analytics in the cloud has become an integral part of enterprise businesses. Big data analytics systems, however, still lack the ability to take user performance goals and budgetary constraints for a task, collectively referred to as task objectives, and automatically configure an analytic job to achieve the objectives.

Our goal, is to come up with a data analytics optimizer that can automatically determine a cluster configuration with a suitable number of cores as well as other runtime system parameters that best meet the task objectives. To achieve this, we also need to design a multi-objective optimizer that constructs a Pareto optimal set of job configurations for task-specific objectives, and recommends new job configurations to best meet these objectives.



## EX-SITU Project-Team

### 3. Research Program

#### 3.1. Research Program

We characterize Extreme Situated Interaction as follows:

**Extreme users.** We study extreme users who make extreme demands on current technology. We know that human beings take advantage of the laws of physics to find creative new uses for physical objects. However, this level of adaptability is severely limited when manipulating digital objects. Even so, we find that creative professionals—artists, designers and scientists—often adapt interactive technology in novel and unexpected ways and find creative solutions. By studying these users, we hope to not only address the specific problems they face, but also to identify the underlying principles that will help us to reinvent virtual tools. We seek to shift the paradigm of interactive software, to establish the laws of interaction that significantly empower users and allow them to control their digital environment.

**Extreme situations.** We develop extreme environments that push the limits of today’s technology. We take as given that future developments will solve “practical” problems such as cost, reliability and performance and concentrate our efforts on interaction in and with such environments. This has been a successful strategy in the past: Personal computers only became prevalent after the invention of the desktop graphical user interface. Smartphones and tablets only became commercially successful after Apple cracked the problem of a usable touch-based interface for the iPhone and the iPad. Although wearable technologies, such as watches and glasses, are finally beginning to take off, we do not believe that they will create the major disruptions already caused by personal computers, smartphones and tablets. Instead, we believe that future disruptive technologies will include fully interactive paper and large interactive displays.

Our extensive experience with the Digiscope WILD and WILDER platforms places us in a unique position to understand the principles of distributed interaction that extreme environments call for. We expect to integrate, at a fundamental level, the collaborative capabilities that such environments afford. Indeed almost all of our activities in both the digital and the physical world take place within a complex web of human relationships. Current systems only support, at best, passive sharing of information, e.g., through the distribution of independent copies. Our goal is to support active collaboration, in which multiple users are actively engaged in the lifecycle of digital artifacts.

**Extreme design.** We explore novel approaches to the design of interactive systems, with particular emphasis on extreme users in extreme environments. Our goal is to empower creative professionals, allowing them to act as both designers and developers throughout the design process. Extreme design affects every stage, from requirements definition, to early prototyping and design exploration, to implementation, to adaptation and appropriation by end users. We hope to push the limits of participatory design to actively support creativity at all stages of the design lifecycle. Extreme design does not stop with purely digital artifacts. The advent of digital fabrication tools and FabLabs has significantly lowered the cost of making physical objects interactive. Creative professionals now create hybrid interactive objects that can be tuned to the user’s needs. Integrating the design of physical objects into the software design process raises new challenges, with new methods and skills to support this form of extreme prototyping.

Our overall approach is to identify a small number of specific projects, organized around four themes: *Creativity, Augmentation, Collaboration* and *Infrastructure*. Specific projects may address multiple themes, and different members of the group work together to advance these different topics.

## ILDA Project-Team

### 3. Research Program

#### 3.1. Introduction

Our ability to acquire or generate, store, process, interlink and query data has increased spectacularly over the last few years. The corresponding advances are commonly grouped under the umbrella of so called *Big Data*. Even if the latter has become a buzzword, these advances are real, and they are having a profound impact in domains as varied as scientific research, commerce, social media, industrial processes or e-government. Yet, looking ahead, emerging technologies related to what we now call the *Web of Data* (a.k.a the Semantic Web) have the potential to create an even larger revolution in data-driven activities, by making information accessible to machines as semistructured data [26] that eventually becomes actionable knowledge. Indeed, novel Web data models considerably ease the interlinking of semi-structured data originating from multiple independent sources. They make it possible to associate machine-processable semantics with the data. This in turn means that heterogeneous systems can exchange data, infer new data using reasoning engines, and that software agents can cross data sources, resolving ambiguities and conflicts between them [77]. Datasets are becoming very rich and very large. They are gradually being made even larger and more heterogeneous, but also much more useful, by interlinking them, as exemplified by the Linked Data initiative [49].

These advances raise research questions and technological challenges that span numerous fields of computer science research: databases, communication networks, security and trust, data mining, as well as human-computer interaction. Our research is based on the conviction that interactive systems play a central role in many data-driven activity domains. Indeed, no matter how elaborate the data acquisition, processing and storage pipelines are, data eventually get processed or consumed one way or another by users. The latter are faced with large, increasingly interlinked heterogeneous datasets (see, *e.g.*, Figure 1 ) that are organized according to complex structures, resulting in overwhelming amounts of both raw data and structured information. Users thus require effective tools to make sense of their data and manipulate them.

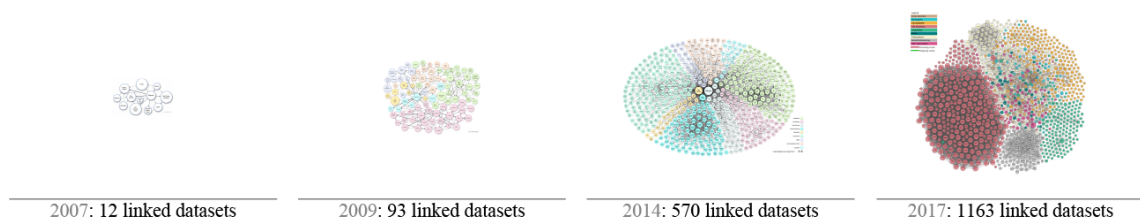


Figure 1. Linking Open Data cloud diagram from 2007 to 2017 – <http://lod-cloud.net>

We approach this problem from the perspective of the Human-Computer Interaction (HCI) field of research, whose goal is to study how humans interact with computers and inspire novel hardware and software designs aimed at optimizing properties such as efficiency, ease of use and learnability, in single-user or cooperative work contexts. More formally, HCI is about designing systems that lower the barrier between users' cognitive model of what they want to accomplish, and computers' understanding of this model. HCI is about the design, implementation and evaluation of computing systems that humans interact with [54], [79]. It is a highly multidisciplinary field, with experts from computer science, cognitive psychology, design, engineering, ethnography, human factors and sociology.

In this broad context, ILDA aims at designing interactive systems that display [35], [61], [87] the data and let users interact with them, aiming to help users better *navigate* and *comprehend* large webs of data represented visually [7], as well as *relate* and *manipulate* them.

Our research agenda consists of the three complementary axes detailed in the following subsections. Designing systems that consider interaction in close conjunction with data semantics is pivotal to all three axes. Those semantics will help drive navigation in, and manipulation of, the data, so as to optimize the communication bandwidth between users and data.

### 3.2. Semantics-driven Data Manipulation

**Participants:** Emmanuel Pietriga, Caroline Appert, Anastasia Bezerianos, Marie Destandau, Hugo Romat, Tong Xue, Léo Colombaro.

The Web of Data has been maturing for the last fifteen years and is starting to gain adoption across numerous application domains (Figure 1). Now that most foundational building blocks are in place, from knowledge representation, inference mechanisms and query languages [50], all the way up to the expression of data presentation knowledge [70] and to mechanisms like look-up services [86] or spreading activation [43], we need to pay significant attention to how human beings are going to interact with this new Web, if it is to “*reach its full potential*” [44].

Most efforts in terms of user interface design and development for the Web of data have essentially focused on tools for software developers or subject-matter experts who create ontologies and populate them [56], [41]. Tools more oriented towards end-users are starting to appear [32], [34], [51], [52], [55], [64], including the so-called *linked data browsers* [49]. However, those browsers are in most cases based on quite conventional point-and-click hypertext interfaces that present data to users in a very page-centric, web-of-documents manner that is ill-suited to navigating in, and manipulating, webs of data.

To be successful, interaction paradigms that let users navigate and manipulate data on the Web have to be tailored to the radically different way of browsing information enabled by it, where users directly interact with the data rather than with monolithic documents. The general research question addressed in this part of our research program is how to design novel interaction techniques that help users manipulate their data more efficiently. By data manipulation, we mean all low-level tasks related to manually creating new content, modifying and cleaning existing content, merging data from different sources, establishing connections between datasets, categorizing data, and eventually sharing the end results with other users; tasks that are currently considered quite tedious because of the sheer complexity of the concepts, data models and syntax, and the interplay between all of them.

Our approach is based on the conviction that there is a strong potential for cross-fertilization, as mentioned earlier: on the one hand, user interface design is essential to the management and understanding of webs of data; on the other hand, interlinked datasets enriched with even a small amount of semantics can help create more powerful user interfaces, that provide users with the right information at the right time.

We envision systems that focus on the data themselves, exploiting the underlying *semantics and structure* in the background rather than exposing them – which is what current user interfaces for the Web of Data often do. We envision interactive systems in which the semantics and structure are not exposed directly to users, but serve as input to the system to generate interactive representations that convey information relevant to the task at hand and best afford the possible manipulation actions.

Relevant publications by team members this year: [22], [15], [17] and major ones in recent years: [7].

### 3.3. Generalized Multi-scale Navigation

**Participants:** Caroline Appert, Anastasia Bezerianos, Olivier Chapuis, Emmanuel Pietriga, Vanessa Peña-Araya, Marie Destandau, Anna Gogolou, Hugo Romat, Dylan Lebout.

The foundational question addressed here is what to display when, where and how, so as to provide effective support to users in their data understanding and manipulation tasks. ILDA targets contexts in which workers have to interact with complementary views on the same data, or with views on different-but-related datasets, possibly at different levels of abstraction. Being able to combine or switch between representations of the data at different levels of detail and merge data from multiple sources in a single representation is central to many scenarios. This is especially true in both of the application domains we consider: mission-critical systems (e.g., natural disaster crisis management) and the exploratory analysis of scientific data (e.g., correlate theories and heterogeneous observational data for an analysis of a given celestial body in Astrophysics).

A significant part of our research over the last ten years has focused on multi-scale interfaces. We designed and evaluated novel interaction techniques, but also worked actively on the development of open-source UI toolkits for multi-scale interfaces (<http://zvtm.sf.net>). These interfaces let users navigate large but relatively homogeneous datasets at different levels of detail, on both workstations [73], [29], [69], [68], [67], [30], [72], [28], [74] and wall-sized displays [63], [58], [71], [62], [31], [37], [36]. This part of the ILDA research program is about extending multi-scale navigation in two directions: 1. Enabling the representation of multiple, spatially-registered but widely varying, multi-scale data layers in Geographical Information Systems (GIS); 2. Generalizing the multi-scale navigation paradigm to interconnected, heterogeneous datasets as found on the Web of Data.

The first research problem has been mainly investigated in collaboration with IGN in the context of ANR project MapMuxing, which stands for *multi-dimensional map multiplexing*, from 2014 to early 2019. Project MapMuxing aimed at going beyond the traditional pan & zoom and overview+detail interface schemes, and at designing and evaluating novel cartographic visualizations that rely on high-quality generalization, *i.e.*, the simplification of geographic data to make it legible at a given map scale [82], [83], and symbol specification. Beyond project MapMuxing, we are also investigating multi-scale multiplexing techniques for geo-localized data in the specific context of ultra-high-resolution wall-sized displays, where the combination of a very high pixel density and large physical surface enable us to explore designs that involve collaborative interaction and physical navigation in front of the workspace. This is work done in cooperation with team Massive Data at Inria Chile.

The second research problem is about the extension of multi-scale navigation to interconnected, heterogeneous datasets. Generalization has a rather straightforward definition in the specific domain of geographical information systems, where data items are geographical entities that naturally aggregate as scale increases. But it is unclear how generalization could work for representations of the more heterogeneous webs of data that we consider in the first axis of our research program. Those data form complex networks of resources with multiple and quite varied relationships between them, that cannot rely on a single, unified type of representation (a role played by maps in GIS applications).

Addressing the limits of current generalization processes is a longer-term, more exploratory endeavor. Here again, the machine-processable semantics and structure of the data give us an opportunity to rethink how users navigate interconnected heterogeneous datasets. Using these additional data, we investigate ways to generalize the multi-scale navigation paradigm to datasets whose layout and spatial relationships can be much richer and much more diverse than what can be encoded with static linear hierarchies as typically found today in interfaces for browsing maps or large imagery. Our goal is thus to design and develop highly dynamic and versatile multi-scale information spaces for heterogeneous data whose structure and semantics are not known in advance, but discovered incrementally.

Relevant publications by team members this year: [24], [20], [13], [14], [11], [19] and major ones in recent years: [10], [2].

### 3.4. Novel Forms of Input for Groups and Individuals

**Participants:** Caroline Appert, Anastasia Bezerianos, Olivier Chapuis, Emmanuel Pietriga, Eugénie Brasier, Emmanuel Courtoux, Raphaël James.

Analyzing and manipulating large datasets can involve multiple users working together in a coordinated manner in multi-display environments: workstations, handheld devices, wall-sized displays [31]. Those users work towards a common goal, navigating and manipulating data displayed on various hardware surfaces in a coordinated manner. Group awareness [48], [25] is central in these situations, as users, who may or may not be co-located in the same room, can have an optimal individual behavior only if they have a clear picture of what their collaborators have done and are currently doing in the global context. We work on the design and implementation of interactive systems that improve group awareness in co-located situations [57], making individual users able to figure out what other users are doing without breaking the flow of their own actions.

In addition, users need a rich interaction vocabulary to handle large, structured datasets in a flexible and powerful way, regardless of the context of work. Input devices such as mice and trackpads provide a limited number of input actions, thus requiring users to switch between modes to perform different types of data manipulation and navigation actions. The action semantics of these input devices are also often too much dependent on the display output. For instance, a mouse movement and click can only be interpreted according to the graphical controller (widget) above which it is moved. We focus on designing powerful input techniques based upon technologies such as tactile surfaces (supported by UI toolkits developed in-house), 3D motion tracking systems, or custom-built controllers [60] *to complement (rather than replace) traditional input devices* such as keyboards, that remain the best method so far for text entry, and indirect input devices such as mice or trackpads for pixel-precise pointing actions.

The input vocabularies we investigate enable users to navigate and manipulate large and structured datasets in environments that involve multiple users and displays that vary in their size, position and orientation [31], [45], each having their own characteristics and affordances: wall displays [63], [89], workstations, tabletops [66], [40], tablets [65], [84], smartphones [88], [38], [80], [81], and combinations thereof [39], [85], [62], [31].

We aim at designing rich interaction vocabularies that go far beyond what current touch interfaces offer, which rarely exceeds five gestures such as simple slides and pinches. Designing larger gesture vocabularies requires identifying discriminating dimensions (e.g., the presence or absence of anchor points and the distinction between internal and external frames of reference [65]) in order to structure a space of gestures that interface designers can use as a dictionary for choosing a coherent set of controls. These dimensions should be few and simple, so as to provide users with gestures that are easy to memorize and execute. Beyond gesture complexity, the scalability of vocabularies also depends on our ability to design robust gesture recognizers that will allow users to fluidly chain simple gestures that make it possible to interlace navigation and manipulation actions.

We also study how to further extend input vocabularies by combining touch [65], [88], [66] and mid-air gestures [63] with physical objects [53], [78], [60] and classical input devices such as keyboards to enable users to input commands to the system or to involve other users in their workflow (request for help, delegation, communication of personal findings, etc.) [33], [59]. Gestures and objects encode a lot of information in their shape, dynamics and direction, that can be directly interpreted in relation with the user, independently from the display output. Physical objects can also greatly improve coordination among actors for, e.g., handling priorities or assigning specific roles.

Relevant publications by team members this year: [9], [23], [16], [22] and major ones in recent years: [1], [10], [5], [3], [8].

## **PETRUS Project-Team**

### **3. Research Program**

#### **3.1. Research Program**

To tackle the challenge introduced above, we identify three main lines of research:

- (Axis 1) Personal cloud server architectures. Based on the intuition that user control, security and privacy are key properties in the definition of trusted personal cloud solutions, our objective is to propose new architectures (encompassing both software and hardware aspects) for secure personal cloud data management and formally prove important bricks of the architecture. We also focus in this axis on administration models and their enforcement in relation to the architecture of the system, so that the exclusive control of a non expert individual can be ensured.
- (Axis 2) Global query evaluation. The goal of this line of research is to provide capabilities for crossing data belonging to multiple individuals (e.g., performing statistical queries over personal data, computing queries on social graphs or organizing participatory data collection) in a fully decentralized setting while providing strong and personalized privacy guarantees. This means proposing new secure distributed database indexing models and query processing strategies. In addition, we concentrate on locally ensuring to each participant the good behaviour of the processing, such that no collective results can be produced if privacy conditions are not respected by other participants.
- (Axis 3) Economic, legal and societal issues. This research axis is more transverse and entails multidisciplinary research, addressing the links between economic, legal, societal and technological aspects. We will follow here a multi-disciplinary approach based on a 3-step methodology: i) identifying important common issues related to privacy and to the exploitation of personal data; ii) characterizing their dimensions in all relevant disciplines and jointly study their entanglement; iii) validating the proposed analysis, models and trade-offs thanks to in vivo experiments.

These contributions will also rely on tools (algorithms, protocols, proofs, etc.) from other communities, namely security (cryptography, secure multiparty computations, formal methods, differential privacy, etc.) and distributed systems (distributed hash tables, gossip protocols, etc.). Beyond the research actions, we structure our software activity around a single common platform (rather than isolated demonstrators), integrating our main research contributions, called PlugDB. This platform is the cornerstone to help validating our research results through accurate performance measurements on a real platform, a common practice in the DB community, and target the best conferences. It is also a strong vector to federate the team, simplify the bootstrapping of new PhD or master students, conduct multi-disciplinary research and open the way to industrial collaborations and technological transfers.