# Activity Report 2019

# Section Scientific Foundations

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

# 3. Research Program

## 3.1. Number fields, class groups and other invariants

**Participants:** Bill Allombert, Jared Guissmo Asuncion, Karim Belabas, Jean-Paul Cerri, Henri Cohen, Jean-Marc Couveignes, Andreas Enge, Fredrik Johansson, Aurel Page.

Modern number theory has been introduced in the second half of the 19th century by Dedekind, Kummer, Kronecker, Weber and others, motivated by Fermat's conjecture: There is no non-trivial solution in integers to the equation $x^n + y^n = z^n$ for $n \geqslant 3$. Kummer's idea for solving Fermat's problem was to rewrite the equation as $(x + y)(x + \zeta y)(x + \zeta^2 y) \cdots (x + \zeta^{n-1} y) = z^n$ for a primitive $n$-th root of unity $\zeta$, which seems to imply that each factor on the left hand side is an $n$-th power, from which a contradiction can be derived.

The solution requires to augment the integers by *algebraic numbers*, that are roots of polynomials in $\mathbb{Z}[X]$. For instance, $\zeta$ is a root of $X^n - 1$, $\sqrt[3]{2}$ is a root of $X^3 - 2$ and $\frac{\sqrt{3}}{5}$ is a root of $25X^2 - 3$. A *number field* consists of the rationals to which have been added finitely many algebraic numbers together with their sums, differences, products and quotients. It turns out that actually one generator suffices, and any number field $K$ is isomorphic to $\mathbb{Q}[X]/(f(X))$, where $f(X)$ is the minimal polynomial of the generator. Of special interest are *algebraic integers*, "numbers without denominators", that are roots of a monic polynomial. For instance, $\zeta$ and $\sqrt[3]{2}$ are integers, while $\frac{\sqrt{3}}{5}$ is not. The *ring of integers* of $K$ is denoted by $\mathcal{O}_K$; it plays the same role in $K$ as $\mathbb{Z}$ in $\mathbb{Q}$.

Unfortunately, elements in $\mathcal{O}_K$ may factor in different ways, which invalidates Kummer's argumentation. Unique factorisation may be recovered by switching to *ideals*, subsets of $\mathcal{O}_K$ that are closed under addition and under multiplication by elements of $\mathcal{O}_K$. In $\mathbb{Z}$, for instance, any ideal is *principal*, that is, generated by one element, so that ideals and numbers are essentially the same. In particular, the unique factorisation of ideals then implies the unique factorisation of numbers. In general, this is not the case, and the *class group* $\mathrm{Cl}_K$ of ideals of $\mathcal{O}_K$ modulo principal ideals and its *class number* $h_K = |\mathrm{Cl}_K|$ measure how far $\mathcal{O}_K$ is from behaving like $\mathbb{Z}$.

Using ideals introduces the additional difficulty of having to deal with *units*, the invertible elements of $\mathcal{O}_K$: Even when $h_K = 1$, a factorisation of ideals does not immediately yield a factorisation of numbers, since ideal generators are only defined up to units. For instance, the ideal factorisation $(6) = (2) \cdot (3)$ corresponds to the two factorisations $6 = 2 \cdot 3$ and $6 = (-2) \cdot (-3)$. While in $\mathbb{Z}$, the only units are $1$ and $-1$, the unit structure in general is that of a finitely generated $\mathbb{Z}$-module, whose generators are the *fundamental units*. The *regulator* $R_K$ measures the "size" of the fundamental units as the volume of an associated lattice.

One of the main concerns of algorithmic algebraic number theory is to explicitly compute these invariants ($\mathrm{Cl}_K$ and $h_K$, fundamental units and $R_K$), as well as to provide the data allowing to efficiently compute with numbers and ideals of $\mathcal{O}_K$; see [36] for a recent account.

The *analytic class number formula* links the invariants $h_K$ and $R_K$ (unfortunately, only their product) to the $\zeta$-function of $K$, $\zeta_K(s) := \prod_{\mathfrak{p} \text{ prime ideal of } \mathcal{O}_K} \left(1 - \mathrm{N}\,\mathfrak{p}^{-s}\right)^{-1}$, which is meaningful when $\Re(s) > 1$, but which may be extended to arbitrary complex $s \neq 1$. Introducing characters on the class group yields a generalisation of $\zeta$- to $L$-functions. The *generalised Riemann hypothesis (GRH)*, which remains unproved even over the rationals, states that any such $L$-function does not vanish in the right half-plane $\Re(s) > 1/2$. The validity of the GRH has a dramatic impact on the performance of number theoretic algorithms. For instance, under GRH, the class group admits a system of generators of polynomial size; without GRH, only exponential bounds are known. Consequently, an algorithm to compute $\mathrm{Cl}_K$ via generators and relations (currently the only viable practical approach) either has to assume that GRH is true or immediately becomes exponential.

When $h_K = 1$ the number field $K$ may be norm-Euclidean, endowing $\mathcal{O}_K$ with a Euclidean division algorithm. This question leads to the notions of the Euclidean minimum and spectrum of $K$, and another task in algorithmic number theory is to compute explicitly this minimum and the upper part of this spectrum, yielding for instance generalised Euclidean gcd algorithms.

## 3.2. Function fields, algebraic curves and cryptology

**Participants:** Karim Belabas, Guilhem Castagnos, Jean-Marc Couveignes, Andreas Enge, Damien Robert, Jean Kieffer, Razvan Barbulescu.

Algebraic curves over finite fields are used to build the currently most competitive public key cryptosystems. Such a curve is given by a bivariate equation $\mathcal{C}(X, Y) = 0$ with coefficients in a finite field $\mathbb{F}_q$. The main classes of curves that are interesting from a cryptographic perspective are *elliptic curves* of equation $\mathcal{C} = Y^2 - (X^3 + aX + b)$ and *hyperelliptic curves* of equation $\mathcal{C} = Y^2 - (X^{2g+1} + \cdots)$ with $g \geqslant 2$.

The cryptosystem is implemented in an associated finite abelian group, the *Jacobian* $\mathrm{Jac}_\mathcal{C}$. Using the language of function fields exhibits a close analogy to the number fields discussed in the previous section. Let $\mathbb{F}_q(X)$ (the analogue of $\mathbb{Q}$) be the *rational function field* with subring $\mathbb{F}_q[X]$ (which is principal just as $\mathbb{Z}$). The *function field* of $\mathcal{C}$ is $K_\mathcal{C} = \mathbb{F}_q(X)[Y]/(\mathcal{C})$; it contains the *coordinate ring* $\mathcal{O}_\mathcal{C} = \mathbb{F}_q[X, Y]/(\mathcal{C})$. Definitions and properties carry over from the number field case $K/\mathbb{Q}$ to the function field extension $K_\mathcal{C}/\mathbb{F}_q(X)$. The Jacobian $\mathrm{Jac}_\mathcal{C}$ is the divisor class group of $K_\mathcal{C}$, which is an extension of (and for the curves used in cryptography usually equals) the ideal class group of $\mathcal{O}_\mathcal{C}$.

The size of the Jacobian group, the main security parameter of the cryptosystem, is given by an $L$-function. The GRH for function fields, which has been proved by Weil, yields the Hasse–Weil bound $(\sqrt{q} - 1)^{2g} \leqslant |\mathrm{Jac}_\mathcal{C}| \leqslant (\sqrt{q} + 1)^{2g}$, or $|\mathrm{Jac}_\mathcal{C}| \approx q^g$, where the *genus* $g$ is an invariant of the curve that correlates with the degree of its equation. For instance, the genus of an elliptic curve is 1, that of a hyperelliptic one is $\frac{\deg_X \mathcal{C} - 1}{2}$. An important algorithmic question is to compute the exact cardinality of the Jacobian.

The security of the cryptosystem requires more precisely that the *discrete logarithm problem* (DLP) be difficult in the underlying group; that is, given elements $D_1$ and $D_2 = xD_1$ of $\mathrm{Jac}_\mathcal{C}$, it must be difficult to determine $x$. Computing $x$ corresponds in fact to computing $\mathrm{Jac}_\mathcal{C}$ explicitly with an isomorphism to an abstract product of finite cyclic groups; in this sense, the DLP amounts to computing the class group in the function field setting.

For any integer $n$, the *Weil pairing* $e_n$ on $\mathcal{C}$ is a function that takes as input two elements of order $n$ of $\mathrm{Jac}_\mathcal{C}$ and maps them into the multiplicative group of a finite field extension $\mathbb{F}_{q^k}$ with $k = k(n)$ depending on $n$. It is bilinear in both its arguments, which allows to transport the DLP from a curve into a finite field, where it is potentially easier to solve. The *Tate-Lichtenbaum pairing*, that is more difficult to define, but more efficient to implement, has similar properties. From a constructive point of view, the last few years have seen a wealth of cryptosystems with attractive novel properties relying on pairings.

For a random curve, the parameter $k$ usually becomes so big that the result of a pairing cannot even be output any more. One of the major algorithmic problems related to pairings is thus the construction of curves with a given, smallish $k$.

## 3.3. Complex multiplication

**Participants:** Jared Guissmo Asuncion, Karim Belabas, Henri Cohen, Jean-Marc Couveignes, Andreas Enge, Fredrik Johansson, Chloe Martindale, Damien Robert.

Complex multiplication provides a link between number fields and algebraic curves; for a concise introduction in the elliptic curve case, see [38], for more background material, [37]. In fact, for most curves $\mathcal{C}$ over a finite field, the endomorphism ring of $\mathrm{Jac}_\mathcal{C}$, which determines its $L$-function and thus its cardinality, is an order in a special kind of number field $K$, called *CM field*. The CM field of an elliptic curve is an imaginary-quadratic field $\mathbb{Q}(\sqrt{D})$ with $D < 0$, that of a hyperelliptic curve of genus $g$ is an imaginary-quadratic extension of a totally real number field of degree $g$. Deuring's lifting theorem ensures that $\mathcal{C}$ is the reduction modulo some prime of a curve with the same endomorphism ring, but defined over the *Hilbert class field* $H_K$ of $K$.

Algebraically, $H_K$ is defined as the maximal unramified abelian extension of $K$; the Galois group of $H_K/K$ is then precisely the class group $\mathrm{Cl}_K$. A number field extension $H/K$ is called *Galois* if $H \simeq K[X]/(f)$ and $H$ contains all complex roots of $f$. For instance, $\mathbb{Q}(\sqrt{2})$ is Galois since it contains not only $\sqrt{2}$, but also the second root $-\sqrt{2}$ of $X^2 - 2$, whereas $\mathbb{Q}(\sqrt[3]{2})$ is not Galois, since it does not contain the root $e^{2\pi i/3}\sqrt[3]{2}$ of $X^3 - 2$. The *Galois group* $\mathrm{Gal}_{H/K}$ is the group of automorphisms of $H$ that fix $K$; it permutes the roots of $f$. Finally, an *abelian* extension is a Galois extension with abelian Galois group.

Analytically, in the elliptic case $H_K$ may be obtained by adjoining to $K$ the *singular value* $j(\tau)$ for a complex valued, so-called *modular* function $j$ in some $\tau \in \mathcal{O}_K$; the correspondence between $\mathrm{Gal}_{H/K}$ and $\mathrm{Cl}_K$ allows to obtain the different roots of the minimal polynomial $f$ of $j(\tau)$ and finally $f$ itself. A similar, more involved construction can be used for hyperelliptic curves. This direct application of complex multiplication yields algebraic curves whose $L$-functions are known beforehand; in particular, it is the only possible way of obtaining ordinary curves for pairing-based cryptosystems.

The same theory can be used to develop algorithms that, given an arbitrary curve over a finite field, compute its $L$-function.

A generalisation is provided by *ray class fields*; these are still abelian, but allow for some well-controlled ramification. The tools for explicitly constructing such class fields are similar to those used for Hilbert class fields.

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

# 3. Research Program

## 3.1. The scientific context

### *3.1.1. Computational fluid mechanics: modeling or not before discretizing ?*

A typical continuous solution of the Navier-Stokes equations at sufficiently large values of the Reynolds number is governed by a wide spectrum of temporal and spatial scales closely connected with the turbulent nature of the flow. The term deterministic chaos employed by Frisch in his enlightening book [46] is certainly conveying most adequately the difficulty in analyzing and simulating this kind of flows. The broadness of the turbulence spectrum is directly controlled by the Reynolds number defined as the ratio between the inertial forces and the viscous forces. This number is not only useful to determine the transition from a laminar to a turbulent flow regime, it also indicates the range of scales of fluctuations that are present in the flow under consideration. Typically, for the velocity field and far from solid walls, the ratio between the largest scale (the integral length scale) and the smallest one (Kolmogorov scale) is proportional to $Re_t^{3/4}$ per dimension, where $Re_t^{3/4}$ is the turbulent Reynolds number, based on the length and velocity scales of the largest turbulent eddies. In addition, for internal flows, viscous effects near the solid walls yield a scaling proportional to $Re_\tau$ per dimension, where $Re_\tau$ is the friction Reynolds number. The smallest scales play a crucial role in the dynamics of the largest ones, which implies that an accurate framework for the computation of turbulent flows must take into account all the scales, which can lead to unrealistic computational costs in real-world applications. Thus, the usual practice to deal with turbulent flows is to choose between an a priori modeling (in most situations) or not (low Re number and rather simple configurations) before proceeding to the discretization step, followed by the simulation itself. If a modeling phase is on the agenda, then one has to choose again among the above-mentioned variety of approaches. The different simulation options and their date of availability for high-Reynolds-number applications are illustrated in Fig. <span style="color:red">1</span> : simulation of turbulent flows can be achieved either by directly solving the Navier-Stokes equations (DNS) or by first applying to the equations a statistical averaging (RANS), a spatial filtering (LES), or a combination of these two operators (hybrid RANS/LES). The new terms brought about by the operator have to be modeled. From a computational point of view, the RANS approach is the least demanding, which explains why historically it has been the workhorse in both the academic and the industrial sectors, and it remains the standard approach nowadays for industrial design, except for very specific applications. It has permitted quite a substantial progress in the understanding of various phenomena such as turbulent combustion or heat transfer. Its inherent inability to provide a time-dependent information has led to promote in the last decade the recourse to either LES or DNS to supplement if not replace RANS. By simulating the large scale structures while modeling the smallest ones, assumed more isotropic, LES proved to be quite a breakthrough to fully take advantage of the increasing power of computers to study complex flow configurations. At the same time, DNS was gradually applied to geometries of increasing complexity (channel flows with values of $Re_\tau$ multiplied by 45 during the last 30 years, jets, turbulent premixed flames, among many others), and proved to be a formidable tool to **(i)** improve our knowledge on turbulent flows and **(ii)** test (i.e., validate or invalidate) and improve the modeling hypotheses inherently associated to the RANS and LES approaches. From a numerical point of view, due to the steady nature of the RANS equations, numerical accuracy is generally not ensured via the use of high-order schemes, but rather on careful grid convergence studies. In contrast, the high computational cost of LES or DNS makes necessary the use of highly-accurate numerical schemes in order to optimize the use of computational resources.

To the noticeable exception of the hybrid RANS-LES modeling, which is not yet accepted as a reliable tool for industrial design, as mentioned in the preamble of the Go4hybrid European program [0], a turbulence model represents turbulent mechanisms in the same way in the whole flow. Thus, depending on its intrinsic strengths

---

[0]<span style="color:red">https://cordis.europa.eu/project/rcn/109107/factsheet/en</span>

and weaknesses, accuracy will be a rather volatile quantity, strongly dependent on the flow configuration. For instance, RANS is perfectly suited to attached boundary layers, but exhibits severe limitations in massively-separated flow regions. Therefore, the turbulence modeling and industrial design communities waver between the desire to continue to rely on the RANS approach, which is unrivaled in terms of computational cost, but is still not able to accurately represent all the complex phenomena; and the temptation to switch to LES, which outperforms RANS in many situations, but is prohibitively expensive in high-Reynolds number wall-bounded flows. In order to account for the limitations of the two approaches and to combine them for significantly improving the overall performance of the models, the hybrid RANS-LES approach has emerged during the last two decades as a viable, intermediate way, and we are definitely inscribing our project in this innovative field of research, with an original approach though, based on temporal filtering (Hybrid temporal LES, HTLES) rather than spatial filtering, and a systematic and progressive validation process against experimental data produced by the team.
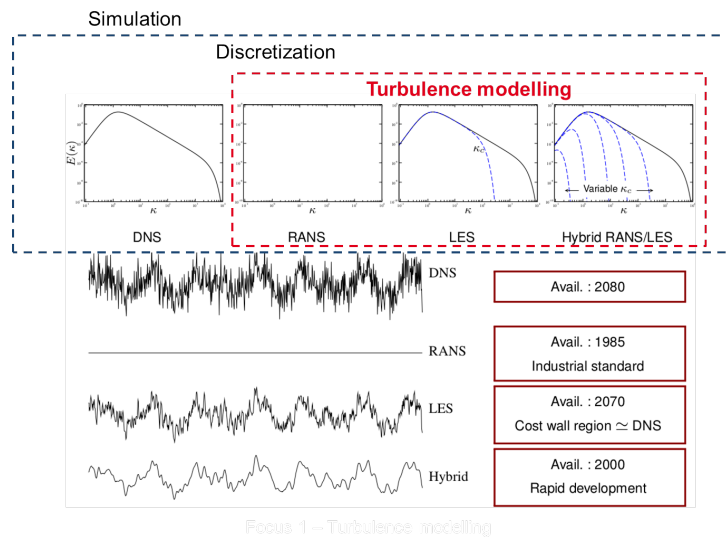


*Figure 1. Schematic view of the different nested steps for turbulent flow simulation: from DNS to hybrid RANS-LES. The approximate dates at which the different approaches are or will be routinely used in the industry are indicated in the boxes on the right (extrapolations based on the present rate of increase in computer performances).*

### 3.1.2. Computational fluid mechanics: high order discretization on unstructured meshes and efficient methods of solution

All the methods considered in the project are mesh-based methods: the computational domain is divided into cells, that have an elementary shape: triangles and quadrangles in two dimensions, and tetrahedra, hexahedra, pyramids, and prisms in three dimensions. If the cells are only regular hexahedra, the mesh is said to be structured. Otherwise, it is said to be unstructured. If the mesh is composed of more than one sort of elementary shape, the mesh is said to be hybrid. In the project, the numerical strategy is based on discontinuous Galerkin methods. These methods were introduced by Reed and Hill [57] and first studied by Lesaint and Raviart [53]. The extension to the Euler system with explicit time integration was mainly led by Shu, Cockburn and their collaborators. The steps of time integration and slope limiting were similar to high-order ENO schemes, whereas specific constraints given by the finite-element nature of the scheme were gradually solved for scalar conservation laws [41], [40], one dimensional systems [39], multidimensional scalar conservation laws [38], and multidimensional systems [42]. For the same system, we can also cite the work of [45], [50], which is

slightly different: the stabilization is made by adding a nonlinear term, and the time integration is implicit. In contrast to continuous Galerkin methods, the discretization of diffusive operators is not straightforward. This is due to the discontinuous approximation space, which does not fit well with the space function in which the diffusive system is well posed. A first stabilization was proposed by Arnold [31]. The first application of discontinuous Galerkin methods to Navier-Stokes equations was proposed in [36] by mean of a mixed formulation. Actually, this first attempt led to a non-compact computational stencil, and was later proved to be unstable. A compactness improvement was made in [37], which was later analyzed, and proved to be stable in a more unified framework [32]. The combination with the $k - \omega$ RANS model was made in [35]. As far as Navier-Stokes equations are concerned, we can also cite the work of [48], in which the stabilization is closer to the one of [32], the work of [54] on local time stepping, or the first use of discontinuous Galerkin methods for direct numerical simulation of a turbulent channel flow done in [43]. Discontinuous Galerkin methods became very popular because:

- They can be developed for any order of approximation.
- The computational stencil of one given cell is limited to the cells with which it has a common face. This stencil does not depend on the order of approximation. This is a pro, compared for example with high-order finite volumes, for which the number of neighbors required increases with the order of approximation.
- They can be developed for any kind of mesh, structured, unstructured, but also for aggregated grids [34]. This is a pro compared not only with finite-difference schemes, which can be developed only on structured meshes, but also compared with continuous finite-element methods, for which the definition of the approximation basis is not clear on aggregated elements.
- $p$-adaptivity is easier than with continuous finite elements, because neighboring elements having a different order are only weakly coupled.
- Upwinding is as natural as for finite volumes methods, which is a benefit for hyperbolic problems.
- As the formulation is weak, boundary conditions are naturally weakly formulated. This is a benefit compared with strong formulations, for example point centered formulation when a point is at the intersection of two kinds of boundary conditions.

For concluding this section, there already exists numerical schemes based on the discontinuous Galerkin method, which proved to be efficient for computing compressible viscous flows. Nevertheless, there remain many things to be improved, which include: efficient shock capturing methods for supersonic flows, high-order discretization of curved boundaries, low-Mach-number behavior of these schemes and combination with second-moment RANS closures. Another aspect that deserves attention is the computational cost of discontinuous Galerkin methods, due to the accurate representation of the solution, calling for a particular care of implementation for being efficient. We believe that this cost can be balanced by the strong memory locality of the method, which is an asset for porting on emerging many-core architectures.

### 3.1.3. *Experimental fluid mechanics: a relevant tool for physical modeling and simulation development*

With the considerable and constant development of computer performance, many people were thinking at the turn of the 21st century that in the short term, CFD would replace experiments, considered as too costly and not flexible enough. Simply flipping through scientific journals such as Journal of Fluid Mechanics, Combustion and Flame, Physics of Fluids or Journal of Computational Physics or through websites such that of Ercoftac [0] is sufficient to convince oneself that the recourse to experiments to provide either a quantitative description of complex phenomena or reference values for the assessment of the predictive capabilities of models and simulations is still necessary. The major change that can be noted though concerns the content of the interaction between experiments and CFD (understood in the broad sense). Indeed, LES or DNS assessment calls for the experimental determination of temporal and spatial turbulent scales, as well as time-resolved measurements and determination of single or multi-point statistical properties of the velocity field. Thus, the team methodology incorporates from the very beginning an experimental component that is operated in strong interaction with the modeling and simulation activities.

---

[0] http://www.ercoftac.org

## 3.2. Research directions

### 3.2.1. Boundary conditions

*3.2.1.1. Generating synthetic turbulence*

A crucial point for any multi-scale simulation able to locally switch (in space or time) from a coarse to a fine level of description of turbulence, is the enrichment of the solution by fluctuations as physically meaningful as possible. Basically, this issue is an extension of the problem of the generation of realistic inlet boundary conditions in DNS or LES of subsonic turbulent flows. In that respect, the method of anisotropic linear forcing (ALF) we have developed in collaboration with EDF proved very encouraging, by its efficiency, its generality and simplicity of implementation. So, it seems natural, on the one hand, to extend this approach to the compressible framework and to implement it in AeroSol. On the other hand, we shall concentrate (in cooperation with EDF R&D in Chatou in the framework of a the CIFRE PhD of V. Duffal) on the theoretical link between the local variations of the scale of description of turbulence (e.g. a sudden variations in the size of the time filter) and the intensity of the ALF forcing, transiently applied to promote the development of missing fluctuating scales.

*3.2.1.2. Stable and non reflecting boundary conditions*

In aerodynamics, and especially for subsonic computations, handling inlet and outlet boundary conditions is a difficult issue. A significant amount of work has already been performed for second-order schemes for Navier-Stokes equations, see [56], [59] and the huge number of papers citing it. On the one hand, we believe that decisive improvements are necessary for higher-order schemes: indeed, the less dissipative the scheme is, the worse impact have the spurious reflections. For this purpose, we will first concentrate on the linearized Navier-Stokes system, and analyze the way to impose boundary conditions in a discontinuous Galerkin framework with a similar approach as in [47]. We will also try to extend the work of [60], which deals with Euler equations, to the Navier-Stokes equations.

### 3.2.2. Turbulence models and model agility

*3.2.2.1. Extension of zero-Mach models to the compressible system*

We shall develop in parallel our multi-scale turbulence modeling and the related adaptive numerical methods of AeroSol. Without prejudice to methods that will be on the podium in the future, a first step in this direction will be to extend to a compressible framework the continuous temporal hybrid RANS/LES method we have developed up to now in a Mach zero context.

*3.2.2.2. Study of wall flows with and without mass or heat transfer at the wall: determination and validation of relevant criteria for hybrid turbulence models*

In the targeted application domains, turbulence/wall interactions and heat transfer at the fluid-solid interface are physical phenomena whose numerical prediction is at the heart of the concerns of our industrial partners. For instance, for a jet engine manufacturer, being able to properly design the configuration of the cooling of the walls of its engine combustion chamber in the presence of thermoacoustic instabilities is based on the proper identification and a thorough understanding of the major mechanisms that drive the dynamics of the parietal transfer. Our objective is to take advantage of our analysis, experimental and computational tools to actively participate in the improvement of the collective knowledge of such kind of transfer. The flow configurations dealt with from the beginning of the project are those of subsonic, single-phase impinging jets or JICF (jets in crossflow) with the possible presence of an interacting acoustic wave. The issue of conjugate heat transfer at the wall will be also gradually investigated. The existing switchover criteria of the hybrid RANS/LES models will be tested on these flow configurations in order to determine their domain of validity. In parallel, the hydrodynamic instability modes of the JICF will be studied experimentally and theoretically (in cooperation with the SIAME laboratory) in order to determine the possibility to drive a change of instability regime (e.g., from absolute to convective) and thus to propose challenging flow conditions that would be relevant for the setting-up of an hybrid LES/DNS approach aimed at supplementing the hybrid RANS/LES approach.

*3.2.2.3. Improvement of turbulence models*

The production and subsequent use of DNS (AeroSol library) and experimental (MAVERIC bench) databases dedicated to the improvement of the physical models is a significant part of our activity. In that respect, our present capability of producing in-situ experimental data for simulation validation and flow analysis is clearly a strongly differentiating mark of our project. The analysis of the DNS and experimental data produced make the improvement of the hybrid RANS/LES approach possible. Our hybrid temporal LES (HTLES) method has a decisive advantage over all other hybrid RANS/LES approaches since it relies on a well-defined time-filtering formalism. This feature greatly facilitates the proper extraction from the databases of the various terms appearing in transport equations obtained at the different scales involved (e.g. from RANS to LES). But we would not be comprehensive in that matter if we were not questioning the relevance of any simulation-experiment comparisons. In other words, a central issue is the following question: are we comparing the same quantities between simulations and experiment? From an experimental point of view, the questions to be raised will be, among others, the possible difference in resolution between the experiment and the simulations, the similar location of the measurement points and simulation points, the acceptable level of random error associated to the necessary finite number of samples. In that respect, the recourse to uncertainty quantification techniques will be advantageously considered.

### 3.2.3. Development of an efficient implicit high-order compressible solver scalable on new architectures

As the flows simulated are very computationally demanding, we will maintain our efforts in the development of AeroSol in the following directions:

- Efficient implementation of the discontinuous Galerkin method.
- Implicit methods based on Jacobian-Free-Newton-Krylov methods and multigrid.
- Porting on heterogeneous architectures.
- Implementation of models.

*3.2.3.1. Efficient implementation of the discontinuous Galerkin method*

In high-order discontinuous Galerkin methods, the unknown vector is composed of a concatenation of the unknowns in the cells of the mesh. An explicit residual computation is composed of three loops: an integration loop on the cells, for which computations in two different cells are independent, an integration loop on boundary faces, in which computations depend on data of one cell and on the boundary conditions, and an integration loop on the interior faces, in which computations depend on data of the two neighboring cells. Each of these loops is composed of three steps: the first step consists in interpolating data at the quadrature points; the second step in computing a nonlinear flux at the quadrature points (the physical flux for the cell loop, an upwind flux for interior faces or a flux adapted to the kind of boundary condition for boundary faces); and the third step in projecting the nonlinear flux on the degrees of freedom.

In this research direction, we propose to exploit the strong memory locality of the method (i.e., the fact that all the unknowns of a cell are stocked contiguously). This formulation can reduce the linear steps of the method (interpolation on the quadrature points and projection on the degrees of freedom) to simple matrix-matrix product which can be optimized. For the nonlinear steps, composed of the computation of the physical flux on the cells and of the numerical flux on the faces, we will try to exploit vectorization.

*3.2.3.2. Implicit methods based on Jacobian-Free-Newton-Krylov methods and multigrid*

For our computations of the IMPACT-AE project, we have used explicit time stepping. The time stepping is limited by the CFL condition, and in our flow, the time step is limited by the acoustic wave velocity. As the Mach number of the flow we simulated in IMPACT-AE was low, the acoustic time restriction is much lower than the turbulent time scale, which is driven by the velocity of the flow. We hope to have a better efficiency by using time implicit methods, for using a time step driven by the velocity of the flow.

Using implicit time stepping in compressible flows in particularly difficult, because the system is fully nonlinear, such that the nonlinear solving theoretically requires to build many times the Jacobian. Our experience in implicit methods is that the building of a Jacobian is very costly, especially in three dimensions and in a high-order framework, because the optimization of the memory usage is very difficult. That is why we propose to use a Jacobian-free implementation, based on [52]. This method consists in solving the linear steps of the Newton method by a Krylov method, which requires Jacobian-vector product. The smart idea of this method is to replace this product by an approximation based on a difference of residual, therefore avoiding any Jacobian computation. Nevertheless, Krylov methods are known to converge slowly, especially for the compressible system when the Mach number is low, because the system is ill-conditioned. In order to precondition, we propose to use an aggregation-based multigrid method, which consists in using the same numerical method on coarser meshes obtained by aggregation of the initial mesh. This choice is driven by the fact that multigrid methods are the only one to scale linearly [61], [62] with the number of unknowns in term of number of operations, and that this preconditioning does not require any Jacobian computation.

Beyond the technical aspects of the multigrid approach, which is challenging to implement, we are also interested in the design of an efficient aggregation. This often means to perform an aggregation based on criteria (anisotropy of the problem, for example) [55]. To this aim, we propose to extend the scalar analysis of [63] to a linearized version of the Euler and Navier-Stokes equations, and try to deduce an optimal strategy for anisotropic aggregation, based on the local characteristics of the flow. Note that discontinuous Galerkin methods are particularly well suited to h-p aggregation, as this kind of methods can be defined on any shape [34].

### 3.2.3.3. Porting on heterogeneous architectures

Until the beginning of the 2000s, the computing capacities have been improved by interconnecting an increasing number of more and more powerful computing nodes. The computing capacity of each node was increased by improving the clock speed, the number of cores per processor, the introduction of a separate and dedicated memory bus per processor, but also the instruction level parallelism, and the size of the memory cache. Even if the number of transistors kept on growing up, the clock speed improvement has flattened since the mid 2000s [58]. Already in 2003, [49] pointed out the difficulties for efficiently using the biggest clusters: "While these super-clusters have theoretical peak performance in the Teraflops range, sustained performance with real applications is far from the peak. Salinas, one of the 2002 Gordon Bell Awards was able to sustain 1.16 Tflops on ASCI White (less than 10% of peak)." From the current multi-core architectures, the trend is now to use many-core accelerators. The idea behind many-core is to use an accelerator composed of a lot of relatively slow and simplified cores for executing the most simple parts of the algorithm. The larger the part of the code executed on the accelerator, the faster the code may become. Therefore, it is necessary to work on the heterogeneous aspects of computations. These heterogeneities are intrinsic to our computations and have two sources. The first one is the use of hybrid meshes, which are necessary for using a locally-structured mesh in a boundary layer. As the different cell shapes (pyramids, hexahedra, prisms and tetrahedra) do not have the same number of degrees of freedom, nor the same number of quadrature points, the execution time on one face or one cell depends on its shape. The second source of heterogeneity are the boundary conditions. Depending on the kind of boundary conditions, user-defined boundary values might be needed, which induces a different computational cost. Heterogeneities are typically what may decrease efficiency in parallel if the workload is not well balanced between the cores. Note that heterogeneities were not dealt with in what we consider as one of the most advanced work on discontinuous Galerkin on GPU [51], as only straight simplicial cell shapes were addressed. For managing at best our heterogeneous computations on heterogeneous architectures, we propose to use the execution runtime StarPU [33]. For this, the discontinuous Galerkin algorithm will be reformulated in terms of a graph of tasks. The previous tasks on the memory management will be useful for that. The linear steps of the discontinuous Galerkin methods require also memory transfers, and one issue consists in determining the optimal task granularity for this step, i.e. the number of cells or face integrations to be sent in parallel on the accelerator. On top of that, the question of which device is the most appropriate to tackle such kind of tasks is to be discussed.

Last, we point out that the combination of shared-memory and distributed-memory parallel programming models is better suited than only the distributed-memory one for multigrid, because in a hybrid version, a wider part of the mesh shares the same memory, therefore making a coarser aggregation possible.

These aspects will benefit from a particularly stimulating environment in the Inria Bordeaux Sud Ouest center around high-performance computing, which is one of the strategic axes of the center.

*3.2.3.4. Implementation of turbulence models in AeroSol and validation*

We will gradually insert models developed in research direction 3.2.2.1  in the AeroSol library in which we develop methods for the DNS of compressible turbulent flows at low Mach number. Indeed, due to its formalism based on temporal filtering, the HTLES approach offers a consistent theoretical framework characterized by a continuous transition from RANS to DNS, even for complex flow configurations (e.g. without directions of spatial homogeneity). As for the discontinuous Galerkin method available presently in AeroSol, it is the best suited and versatile method able to meet the requirements of accuracy, stability and cost related to the local (varying) level of resolution of the turbulent flow at hand, regardless of its complexity. The first step in this direction was taken in 2017 during the internship of Axelle Perraud, who has implemented a turbulence model ($k$-$\omega$-SST) in the Aerosol library.

### 3.2.4. Validation of the simulations: test flow configurations

To supplement whenever necessary the test flow configuration of MAVERIC and apart from configurations that could emerge in the course of the project, the following configurations for which either experimental data, simulation data or both have been published will be used whenever relevant for benchmarking the quality of our agile computations:

- The impinging turbulent jet (simulations).
- The ORACLES two-channel dump combustor developed in the European projects LES4LPP and MOLECULES.
- The non reactive single-phase PRECCINSTA burner (monophasic swirler), a configuration that has been extensively calculated in particular with the AVBP and Yales2 codes.
- The LEMCOTEC configuration (monophasic swirler + effusion cooling).
- The ONERA MERCATO two-phase injector configuration provided the question of confidentiality of the data is not an obstacle.
- Rotating turbulent flows with wall interaction and heat transfer.
- Turbulent flows with buoyancy.

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

# 3. Research Program

## 3.1. Variational discrete asymptotic modelling

In many of the applications we consider, intermediate fidelity models are or can be derived using an asymptotic expansion for the relevant scale resolving PDEs, and eventually considering some averaged for of the resulting continuous equations. The resulting systems of PDEs are often very complex and their characterization, e.g. in terms of stability, unclear, or poor, or too complex to allow to obtain discrete analogy of the continuous properties. This makes the numerical approximation of these PDE systems a real challenge. Moreover, most of these models are often based on asymptotic expansions involving small geometrical scales. This is true for many applications considered here involving flows in/of thin layers (free surface waves, liquid films on wings generating ice layers, oxide flows in material cracks, etc). This asymptotic expansion is nothing else than a discretization (some sort of Taylor expansion) in terms of the small parameter. The actual discretization of the PDE system is another expansion in space involving as a small parameter the mesh size. What is the interaction between these two expansions ? Could we use the spatial discretization (truncation error) as means of filtering undesired small scales instead of having to explicitly derive PDEs for the large scales ? We will investigate in depth the relations between asymptotics and discretization by :

- comparing the asymptotic limits of discretized forms of the relevant scale resolving equations with the discretization of the analogous continuous asymptotic PDEs. Can we discretize a well understood system of PDEs instead of a less understood and more complex one ? ;

- study the asymptotic behaviour of error terms generated by coarse one-dimensional discretization in the direction of the "small scale". What is the influence of the number of cells along the vertical direction, and of their clustering ? ;

- derive equivalent continuous equations (modified equations) for anisotropic discretizations in which the direction is direction of the "small scale" is approximated with a small number of cells. What is the relation with known asymptotic PDE systems ?

Our objective is to gain sufficient control of the interaction between discretization and asymptotics to be able to replace the coupling of several complex PDE systems by adaptive strongly anisotrotropic finite elemient approximations of relevant and well understood PDEs. Here the anisotropy is intended in the sense of having a specific direction in which a much poorer (and possibly variable with the flow conditions) polynomial approximation (expansion) is used. The final goal is, profiting from the availability of faster and cheaper computational platforms, to be able to automatically control numerical *and* physical accuracy of the model with the same techniques. This activity will be used to improve our modelling in coastal engineering as well as for de-anti icing systems, wave energy converters, composite materials (cf. next sections).

In parallel to these developments, we will make an effort in to gain a better understanding of continuous asymptotic PDE models. We will in particular work on improving, and possibly, simplifying their numerical approximation. An effort will be done in trying to embed in these more complex nonlinear PDE models discrete analogs of operator identities necessary for stability (see e.g. the recent work of [70], [72] and references therein).

## 3.2. High order discretizations on moving adaptive meshes

We will work on both the improvement of high order mesh generation and adaptation techniques, and the construction of more efficient, adaptive high order discretisation methods.

Concerning curved mesh generation, we will focus on two points. First propose a robust and automatic method to generate curved simplicial meshes for realistic geometries. The untangling algorithm we plan to develop is a hybrid technique that gathers a local mesh optimization applied on the surface of the domain and a linear elasticity analogy applied in its volume. Second we plan to extend the method proposed in [26] to hybrid meshes (prism/tetra).

For time dependent adaptation we will try to exploit as much as possible the use of $r-$adaptation techniques based on the solution of some PDE system for the mesh. We will work on enhancing the results of [29] by developing more robust nonlinear variants allowing to embed rapidly moving objects. For this the use of non-linear mesh PDEs (cf e.g. [80], [85], [38]), combined with Bezier type approximations for the mesh displacements to accommodate high order curved meshes [26], and with improved algorithms to discretize accurately and fast the elliptic equations involved. For this we will explore different type of relaxation methods, including those proposed in [71], [75], [74] allowing to re-use high order discretizations techniques already used for the flow variables. All these modelling approaches for the mesh movement are based on some minimization argument, and do not allow easily to take into account explicitly properties such as e.g. the positivity of nodal volumes. An effort will be made to try to embed these properties, as well as to improve the control on the local mesh sizes obtained. Developments made in numerical methods for Lagrangian hydrodynamics and compressible materials may be a possible path for these objectives (see e.g. [49], [91], [90] and references therein). We will stretch the use of these techniques as much as we can, and couple them with remeshing algorithms based on local modifications plus conservative, high order, and monotone ALE (or other) remaps (cf. [27], [58], [92], [47] and references therein).

The development of high order schemes for the discretization of the PDE will be a major part of our activity. We will work from the start in an Arbitrary Lagrangian Eulerian setting, so that mesh movement will be easily accommodated, and investigate the following main points:

- the ALE formulation is well adapted both to handle moving meshes, and to provide conservative, high order, and monotone remaps between different meshes. We want to address the issue of cost-accuracy of adaptive mesh computations by exploring different degrees of coupling between the flow and the mesh PDEs. Initial experience has indicated that a clever coupling may lead to a considerable CPU time reduction for a given resolution [29]. This balance is certainly dependent on the nature of the PDEs, on the accuracy level sought, on the cost of the scheme, and on the time stepping technique. All these elements will be taken into account to try to provide the most efficient formulation ;

- the conservation of volume, and the subsequent preservation of constant mass-momentum-energy states on deforming domains is one of the most primordial elements of Arbitrary Lagrangian-Eulerian formulations. For complex PDEs as the ones considered here, of especially for some applications, there may be a competition between the conservation of e.g. mass, an the conservation of other constant states, as important as mass. This is typically the case for free surface flows, in which mass preservation is in competitions with the preservation of constant free surface levels [29]. Similar problems may arise in other applications. Possible solutions to this competition may come from super-approximation (use of higher order polynomials) of some of the data allowing to reduce (e.g. bathymetry) the error in the preservation of one of the competing quantities. This is similar to what is done in super-parametric approximations of the boundaries of an object immersed in the flow, except that in our case the data may enter the PDE explicitly and not only through the boundary conditions. Several efficient solutions for this issue will be investigated to obtain fully conservative moving mesh approaches:

- an issue related to the previous one is the accurate treatment of wall boundaries. It is known that even for standard lower order (second) methods, a higher order, curved, approximation of the boundaries may be beneficial. This, however, may become difficult when considering moving objects, as in the case e.g. of the study of the impact of ice debris in the flow. To alleviate this issue, we plan to follow on with our initial work on the combined use of immersed boundaries techniques with high order, anisotropic (curved) mesh adaptation. In particular, we will develop combined approaches involving high order hybrid meshes on fixed boundaries with the use of penalization techniques and immersed

boundaries for moving objects. We plan to study the accuracy obtainable across discontinuous functions with $r-$adaptive techniques, and otherwise use whenever necessary anisotropic meshes to be able to provide a simplified high order description of the wall boundary (cf. [69]). The use of penalization will also provide a natural setting to compute immediate approximations of the forces on the immersed body [73], [76]. An effort will be also made on improving the accuracy of these techniques using e.g. higher order approaches, either based on generalizations of classical splitting methods [59], or on some iterative Defect Correction method (see e.g. [40]) ;

- the proper treatment of different physics may be addressed by using mixed/hybrid schemes in which different variables/equations are approximated using a different polynomial expansion. A typical example is our work on the discretization of highly non-linear wave models [54] in which we have shown how to use a standard continuous Galerkin method for the elliptic equation/variable representative of the dispersive effects, while the underlying hyperbolic system is evolved using a (discontinuous) third order finite volume method. This technique will be generalized to other classes of discontinuous methods, and similar ideas will be used in other context to provide a flexible approximation. Such mathods have clear advantages in multiphase flows but not only. A typical example where such mixed methods are beneficial are flows involving different species and tracer equations, which are typically better treated with a discontinuous approximation. Another example is the use of this mixed approximation to describe the topography with a high order continuous polynomial even in discontinuous method. This allows to greatly simplify the numerical treatment of the bathymetric source terms ;

- the enhancement of stabilized methods based on some continuous finite element approximation will remain a main topic. We will further pursue the study on the construction of simplified stabilization operators which do not involve any contributions to the mass matrix. We will in particular generalize our initial results to higher order spatial approximations using cubature points, or Bezier polynomials, or also hierarchical approximations. This will also be combined with time dependent variants of the reconstruction techniques initially proposed by D. Caraeni [39], allowing to have a more flexible approach similar to the so-called $P^nP^m$ method [52], [84]. How to localize these enhancements, and to efficiently perform local reconstructions/enrichment, as well as $p-$adaptation, and handling hanging nodes will also be a main line of work. A clever combination of hierarchical enrichment of the polynomials, with a constrained approximation will be investigated. All these developments will be combined with the shock capturing/positivity preserving construction we developed in the past. Other discontinuity resolving techniques will be investigated as well, such as face limiting techniques as those partially studied in [56] ;

- time stepping is an important issue, especially in presence of local mesh adaptation. The techniques we use will force us to investigate local and multilevel techniques. We will study the possibility constructing semi-implicit methods combining extrapolation techniques with space-time variational approaches. Other techniques will be considered, as multi-stage type methods obtained using Defect-Correction, Multi-step Runge-Kutta methods [36], as well as spatial partitioning techniques [65]. A major challenge will be to be able to guarantee sufficient locality to the time integration method to allow to efficiently treat highly refined meshes, especially for viscous reactive flows. Another challenge will be to embed these methods in the stabilized methods we will develop.

## 3.3. Coupled approximation/adaptation in parameter and physical space

As already remarked, classical methods for uncertainty quantification are affected by the so-called Curse-of-Dimensionality. Adaptive approaches proposed so far, are limited in terms of efficiency, or of accuracy. Our aim here is to develop methods and algorithms permitting a very high-fidelity simulation in the physical and in the stochastic space at the same time. We will focus on both non-intrusive and intrusive approaches.

Simple non-intrusive techniques to reduce the overall cost of simulations under uncertainty will be based on adaptive quadrature in stochastic space with mesh adaptation in physical space using error monitors related to the variance of to the sensitivities obtained e.g. by an ANOVA decomposition. For steady state problems,

remeshing using metric techniques is enough. For time dependent problems both mesh deformation and re-meshing techniques will be used. This approach may be easily used in multiple space dimensions to minimize the overall cost of model evaluations by using high order moments of the properly chosen output functional for the adaptation (as in optimization). Also, for high order curved meshes, the use of high order moments and sensitivities issued from the UQ method or optimization provides a viable solution to the lack of error estimators for high order schemes.

Despite the coupling between stochastic and physical space, this approach can be made massively parallel by means of extrapolation/interpolation techniques for the high order moments, in time and on a reference mesh, guaranteeing the complete independence of deterministic simulations. This approach has the additional advantage of being feasible for several different application codes due to its non-intrusive character.

To improve on the accuracy of the above methods, intrusive approaches will also be studied. To propagate uncertainties in stochastic differential equations, we will use Harten's multiresolution framework, following [25]. This framework allows a reduction of the dimensionality of the discrete space of function representation, defined in a proper stochastic space. This reduction allows a reduction of the number of explicit evaluations required to represent the function, and thus a gain in efficiency. Moreover, multiresolution analysis offers a natural tool to investigate the local regularity of a function and can be employed to build an efficient refinement strategy, and also provides a procedure to refine/coarsen the stochastic space for unsteady problems. This strategy should allow to capture and follow all types of flow structures, and, as proposed in [25], allows to formulate a non-linear scheme in terms of compression capabilities, which should allow to handle non-smooth problems. The potential of the method also relies on its moderate intrusive behaviour, compared to e.g. spectral Galerkin projection, where a theoretical manipulation of the original system is needed.

Several activities are planned to generalize our initial work, and to apply it to complex flows in multiple (space) dimensions and with many uncertain parameters.

The first is the improvement of the efficiency. This may be achieved by means of anisotropic mesh refinement, and by experimenting with a strong parallelization of the method. Concerning the first point, we will investigate several anisotropic refinement criteria existing in literature (also in the UQ framework), starting with those already used in the team to adapt the physical grid. Concerning the implementation, the scheme formulated in [25] is conceived to be highly parallel due to the external cycle on the number of dimensions in the space of uncertain parameters. In principle, a number of parallel threads equal to the number of spatial cells could be employed. The scheme should be developed and tested for treating unsteady and discontinuous probability density function, and correlated random variables. Both the compression capabilities and the accuracy of the scheme (in the stochastic space) should be enhanced with a high-order multidimensional conservative and non-oscillatory polynomial reconstruction (ENO/WENO).

Another main objective is related to the use of multiresolution in both physical and stochastic space. This requires a careful handling of data and an updated definition of the wavelet. Until now, only a weak coupling has been performed, since the number of points in the stochastic space varies according to the physical space, but the number of points in the physical space remains unchanged. Several works exist on the multiresolution approach for image compression, but this could be the first time i in which this kind of approach would be applied at the same time in the two spaces with an unsteady procedure for refinement (and coarsening). The experimental code developed using these technologies will have to fully exploit the processing capabilities of modern massively parallel architectures, since there is a unique mesh to handle in the coupled physical/stochastic space.

## 3.4. Robust multi-fidelity modelling for optimization and certification

Due to the computational cost, it is of prominent importance to consider multi-fidelity approaches gathering high-fidelity and low-fidelity computations. Note that low-fidelity solutions can be given by both the use of surrogate models in the stochastic space, and/or eventually some simplified choices of physical models of some element of the system. Procedures which deal with optimization considering uncertainties for complex problems may require the evaluation of costly objective and constraint functions hundreds or even thousands of times. The associated costs are usually prohibitive. For these reason, the robustness of the optimal

solution should be assessed, thus requiring the formulation of efficient methods for coupling optimization and stochastic spaces. Different approaches will be explored. Work will be developed along three axes:

1. a robust strategy using the statistics evaluation will be applied separately, *i.e.* using only low or high-fidelity evaluations. Some classical optimization algorithms will be used in this case. Influence of high-order statistics and model reduction in the robust design optimization will be explored, also by further developing some low-cost methods for robust design optimization working on the so-called Simplex$^2$ method [45] ;

2. a multi-fidelity strategy by using in an efficient way low fidelity and high-fidelity estimators both in physical and stochastic space will be conceived, by using a Bayesian framework for taking into account model discrepancy and a PC expansion model for building a surrogate model ;

3. develop advanced methods for robust optimization. In particular, the Simplex$^2$ method will be modified for introducing a hierarchical refinement with the aim to reduce the number of stochastic samples according to a given design in an adaptive way.

This work is related to the activities foreseen in the EU contract MIDWEST, in the ANR LabCom project VIPER (currently under evaluation), in a joint project with DGA and VKI, in two projects under way with AIRBUS and SAFRAN-HERAKLES.

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

# 3. Research Program

## 3.1. Introduction

The scientific objectives of the team are to provide mathematical tools for modeling and optimization of complex systems. These systems require mathematical representations which are in essence dynamic, multi-model and stochastic. This increasing complexity poses genuine scientific challenges in the domain of modeling and optimization. More precisely, our research activities are focused on stochastic optimization and (parametric, semi-parametric, multidimensional) statistics which are complementary and interlinked topics. It is essential to develop simultaneously statistical methods for the estimation and control methods for the optimization of the models.

## 3.2. Main research topics

**Stochastic modeling**: Markov chain, Piecewise Deterministic Markov Processes (PDMP), Markov Decision Processes (MDP).

The mathematical representation of complex systems is a preliminary step to our final goal corresponding to the optimization of its performance. The team CQFD focuses on two complementary types of approaches. The first approach is based on mathematical representations built upon physical models where the dynamic of the real system is described by *stochastic processes*. The second one consists in studying the modeling issue in an abstract framework where the real system is considered as black-box. In this context, the outputs of the system are related to its inputs through a *statistical model*. Regarding stochastic processes, the team studies Piecewise Deterministic Markov Processes (PDMPs) and Markov Decision Processes (MDPs). These two classes of Markov processes form general families of controlled stochastic models suitable for the design of sequential decision-making problems. They appear in many fields such as biology, engineering, computer science, economics, operations research and provide powerful classes of processes for the modeling of complex systems. Our contribution to this topic consists in expressing real-life industrial problems into these mathematical frameworks. Regarding statistical methods, the team works on dimension reduction models. They provide a way to understand and visualize the structure of complex data sets. Furthermore, they are important tools in several different areas such as data analysis and machine learning, and appear in many applications such as biology, genetics, environment and recommendation systems. Our contribution to this topic consists in studying semiparametric modeling which combines the advantages of parametric and nonparametric models.

**Estimation methods**: estimation for PDMP; estimation in non- and semi- parametric regression modeling.

To the best of our knowledge, there does not exist any general theory for the problems of estimating parameters of PDMPs although there already exist a large number of tools for sub-classes of PDMPs such as point processes and marked point processes. To fill the gap between these specific models and the general class of PDMPs, new theoretical and mathematical developments will be on the agenda of the whole team. In the framework of non-parametric regression or quantile regression, we focus on kernel estimators or kernel local linear estimators for complete data or censored data. New strategies for estimating semi-parametric models via recursive estimation procedures have also received an increasing interest recently. The advantage of the recursive estimation approach is to take into account the successive arrivals of the information and to refine, step after step, the implemented estimation algorithms. These recursive methods do require restarting calculation of parameter estimation from scratch when new data are added to the base. The idea is to use only the previous estimations and the new data to refresh the estimation. The gain in time could be very interesting and there are many applications of such approaches.

**Dimension reduction**: dimension-reduction via SIR and related methods, dimension-reduction via multidimensional and classification methods.

Most of the dimension reduction approaches seek for lower dimensional subspaces minimizing the loss of some statistical information. This can be achieved in modeling framework or in exploratory data analysis context.

In modeling framework we focus our attention on semi-parametric models in order to conjugate the advantages of parametric and nonparametric modeling. On the one hand, the parametric part of the model allows a suitable interpretation for the user. On the other hand, the functional part of the model offers a lot of flexibility. In this project, we are especially interested in the semi-parametric regression model $Y = f(X'\theta) + \varepsilon$, the unknown parameter $\theta$ belongs to $\mathbb{R}^p$ for a single index model, or is such that $\theta = [\theta_1, \cdots, \theta_d]$ (where each $\theta_k$ belongs to $\mathbb{R}^p$ and $d \leq p$ for a multiple indices model), the noise $\varepsilon$ is a random error with unknown distribution, and the link function $f$ is an unknown real valued function. Another way to see this model is the following: the variables $X$ and $Y$ are independent given $X'\theta$. In our semi-parametric framework, the main objectives are to estimate the parametric part $\theta$ as well as the nonparametric part which can be the link function $f$, the conditional distribution function of $Y$ given $X$ or the conditional quantile $q_\alpha$. In order to estimate the dimension reduction parameter $\theta$ we focus on the Sliced Inverse Regression (SIR) method which has been introduced by Li [37] and Duan and Li [35].

Methods of dimension reduction are also important tools in the field of data analysis, data mining and machine learning.They provide a way to understand and visualize the structure of complex data sets.Traditional methods among others are principal component analysis for quantitative variables or multiple component analysis for qualitative variables. New techniques have also been proposed to address these challenging tasks involving many irrelevant and redundant variables and often comparably few observation units. In this context, we focus on the problem of synthetic variables construction, whose goals include increasing the predictor performance and building more compact variables subsets. Clustering of variables is used for feature construction. The idea is to replace a group of ''similar'' variables by a cluster centroid, which becomes a feature. The most popular algorithms include K-means and hierarchical clustering. For a review, see, e.g., the textbook of Duda [36].

**Stochastic control**: optimal stopping, impulse control, continuous control, linear programming.

The main objective is to develop *approximation techniques* to provide quasi-optimal feasible solutions and to derive *optimality results* for control problems related to MDPs and PDMPs:

- *Approximation techniques*. The analysis and the resolution of such decision models mainly rely on the maximum principle and/or the dynamic/linear programming techniques together with their various extensions such as the value iteration (VIA) and the policy iteration (PIA) algorithm. However, it is well known that these approaches are hardly applicable in practice and suffer from the so-called *curse of dimensionality*. Hence, solving numerically a PDMP or an MDP is a difficult and important challenge. Our goal is to obtain results which are both consistent from a theoretical point of view and computationally tractable and accurate from an application standpoint. It is important to emphasize that these research objectives were not planned in our initial 2009 program.

  Our objective is to propose approximation techniques to efficiently compute the optimal value function and to get quasi-optimal controls for different classes of constrained and unconstrained MDPs with general state/action spaces, and possibly unbounded cost function. Our approach is based on combining the linear programming formulation of an MDP with probabilistic approximation techniques related to quantization techniques and the theory of empirical processes. An other aim is to apply our methods to specific industrial applications in collaboration with industrial partners such as Airbus Defence & Space, Naval Group and Thales.

  Asymptotic approximations are also developed in the context of queueing networks, a class of models where the decision policy of the underlying MDP is in some sense fixed a priori, and our main goal is to study the transient or stationary behavior of the induced Markov process. Even though the decision policy is fixed, these models usually remain intractable to solve. Given this complexity, the team has developed analyses in some limiting regime of practical interest, i.e., queueing models in the large-network, heavy-traffic, fluid or mean-field limit. This approach is helpful to obtain a simpler mathematical description of the system under investigation, which is often given in terms of ordinary differential equations or convex optimization problems.

- *Optimality results.* Our aim is to investigate new important classes of optimal stochastic control problems including constraints and combining continuous and impulse actions for MDPs and PDMPs. In this framework, our objective is to obtain different types of optimality results. For example, we intend to provide conditions to guarantee the existence and uniqueness of the optimality equation for the problem under consideration and to ensure existence of an optimal (and $\epsilon$-optimal) control strategy. We also plan to analyze the structural properties of the optimal strategies as well as to study the associated infinite dimensional linear programming problem. These results can be seen as a first step toward the development of numerical approximation techniques in the sense described above.

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

# 3. Research Program

## 3.1. General methodology

**Fully Developed Turbulence (FDT)**   Turbulence at very high Reynolds numbers; systems in FDT are beyond deterministic chaos, and symmetries are restored in a statistical sense only, and multi-scale correlated structures are landmarks. Generalizing to more random uncorrelated multi-scale structured turbulent fields.

**Compact Representation**   Reduced representation of a complex signal (dimensionality reduction) from which the whole signal can be reconstructed. The reduced representation can correspond to points randomly chosen, such as in Compressive Sensing, or to geometric localization related to statistical information content (framework of reconstructible systems).

**Sparse representation**   The representation of a signal as a linear combination of elements taken in a dictionary (frame or Hilbertian basis), with the aim of finding as less as possible non-zero coefficients for a large class of signals.

**Universality class**   In theoretical physics, the observation of the coincidence of the critical exponents (behaviour near a second order phase transition) in different phenomena and systems is called universality. Universality is explained by the theory of the renormalization group, allowing for the determination of the changes followed by structured fluctuations under rescaling, a physical system is the stage of. The notion is applicable with caution and some differences to generalized out-of-equilibrium or disordered systems. Non-universal exponents (without definite classes) exist in some universal slowing dynamical phenomena like the glass transition and kindred. As a consequence, different macroscopic phenomena displaying multiscale structures (and their acquisition in the form of complex signals) may be grouped into different sets of generalized classes.

Every signal conveys, as a measure experiment, information on the physical system whose signal is an acquisition of. As a consequence, it seems natural that signal analysis or compression should make use of physical modelling of phenomena: the goal is to find new methodologies in signal processing that goes beyond the simple problem of interpretation. Physics of disordered systems, and specifically physics of (spin) glasses is putting forward new algorithmic resolution methods in various domains such as optimization, compressive sensing etc. with significant success notably for NP hard problem heuristics. Similarly, physics of turbulence introduces phenomenological approaches involving multifractality. Energy cascades are indeed closely related to geometrical manifolds defined through random processes. At these structures' scales, information in the process is lost by dissipation (close to the lower bound of inertial range). However, all the cascade is encoded in the geometric manifolds, through long or short distance correlations depending on cases. How do these geometrical manifold structures organize in space and time, in other words, how does the scale entropy cascades itself ? To unify these two notions, a description in term of free energy of a generic physical model is sometimes possible, such as an elastic interface model in a random nonlinear energy landscape : This is for instance the correspondence between compressible stochastic Burgers equation and directed polymers in a disordered medium. Thus, trying to unlock the fingerprints of cascade-like structures in acquired natural signals becomes a fundamental problem, from both theoretical and applicative viewpoints.

To illustrate the general methodology undertaken, let us focus on an example conducted in the study of physiological time series: the analysis of signals recorded from the electrical activity of the heart in the general setting of Atrial Fibrillation (AF). AF is a cardiac arrhythmia characterized by rapid and irregular atrial electrical activity with a high clinical impact on stroke incidence. Best available therapeutic strategies combine pharmacological and surgical means. But when successful, they do not always prevent long-term relapses.

Initial success becomes all the more tricky to achieve as the arrhythmia maintains itself and the pathology evolves into sustained or chronic AF. This raises the open crucial issue of deciphering the mechanisms that govern the onset of AF as well as its perpetuation. We have developed a wavelet-based multi-scale strategy to analyze the electrical activity of human hearts recorded by catheter electrodes, positioned in the coronary sinus (CS), during episodes of chronic AF. We have computed the so-called multifractal spectra using two variants of the wavelet transform modulus maxima method, the moment (partition function) method and the magnitude cumulant method (checking confidence intervals with surrogate data). Application of these methods to long time series recorded in a patient with chronic AF provides quantitative evidence of the multifractal intermittent nature of the electric energy of passing cardiac impulses at low frequencies, *i.e.* for times ($> \sim 0.5$ s) longer than the mean interbeat ($\simeq 10^{-1}$s). We have also reported the results of a two-point magnitude correlation analysis which infers the absence of a multiplicative time-scale structure underlying multifractal scaling. The electric energy dynamics looks like a "multifractal white noise" with quadratic (log-normal) multifractal spectra. *These observations challenge concepts of functional reentrant circuits in mechanistic theories of AF.* A transition is observed in the computed multifractal spectra which group according to two distinct areas, consistently with the anatomical substrate binding to the CS, namely the left atrial posterior wall, and the ligament of Marshall which is innervated by the ANS. These negative results challenge also the existing models, which by principle cannot explain such results. As a consequence, we go beyond the existing models and propose a mathematical model of a denervated heart where the kinetics of gap junction conductance alone induces a desynchronization of the myocardial excitable cells, accounting for the multifractal spectra found experimentally in the left atrial posterior wall area (devoid of ANS influence).

## 3.2. Turbulence in insterstellar clouds and Earth observation data

The research described in this section is a collaboration effort of GEOSTAT, CNRS LEGOS (Toulouse), CNRS LAM (Marseille Laboratory for Astrophysics), MERCATOR (Toulouse), IIT Roorkee, Moroccan Royal Center for Teledetection (CRST), Moroccan Center for Science CNRST, Rabat University, University of Heidelberg. Researchers involved:

- GEOSTAT: H. Yahia, N. Brodu, K. Daoudi, A. El Aouni, A. Tamim
- CNRS LAB: S. Bontemps, N. Schneider
- CNRS LEGOS: V. Garçon, I. Hernandez-Carrasco, J. Sudre, B. Dewitte
- CNRS LAM: T. Fusco
- CNRST, CRTS, Rabat University: D. Aboutajdine, A. Atillah, K. Minaoui
- Universiy of Heidelberg: C. Garbe

The analysis and modeling of natural phenomena, specially those observed in geophysical sciences and in astronomy, are influenced by statistical and multiscale phenomenological descriptions of turbulence; indeed these descriptions are able to explain the partition of energy within a certain range of scales. A particularly important aspect of the statistical theory of turbulence lies in the discovery that the support of the energy transfer is spatially highly non uniform, in other terms it is *intermittent* [70]. Because of the absence of localization of the Fourier transform, linear methods are not successful to unlock the multiscale structures and cascading properties of variables which are of primary importance as stated by the physics of the phenomena. This is the reason why new approaches, such as DFA (Detrented Fluctuation Analysis), Time-frequency analysis, variations on curvelets [66] etc. have appeared during the last decades. Recent advances in dimensionality reduction, and notably in Compressive Sensing, go beyond the Nyquist rate in sampling theory using nonlinear reconstruction, but data reduction occur at random places, independently of geometric localization of information content, which can be very useful for acquisition purposes, but of lower impact in signal analysis. We are successfully making use of a microcanonical formulation of the multifractal theory, based on predictability and reconstruction, to study the turbulent nature of interstellar molecular or atomic clouds. Another important result obtained in GEOSTAT is the effective use of multiresolution analysis associated to optimal inference along the scales of a complex system. The multiresolution analysis is performed on dimensionless quantities given by the *singularity exponents* which encode properly the geometrical structures

associated to multiscale organization. This is applied successfully in the derivation of high resolution ocean dynamics, or the high resolution mapping of gaseous exchanges between the ocean and the atmosphere; the latter is of primary importance for a quantitative evaluation of global warming. Understanding the dynamics of complex systems is recognized as a new discipline, which makes use of theoretical and methodological foundations coming from nonlinear physics, the study of dynamical systems and many aspects of computer science. One of the challenges is related to the question of *emergence* in complex systems: large-scale effects measurable macroscopically from a system made of huge numbers of interactive agents [31], [61]. Some quantities related to nonlinearity, such as Lyapunov exponents, Kolmogorov-Sinai entropy etc. can be computed at least in the phase space [32]. Consequently, knowledge from acquisitions of complex systems (which include *complex signals*) could be obtained from information about the phase space. A result from F. Takens [67] about strange attractors in turbulence has motivated the theoretical determination of nonlinear characteristics associated to complex acquisitions. Emergence phenomena can also be traced inside complex signals themselves, by trying to localize information content geometrically. Fundamentally, in the nonlinear analysis of complex signals there are broadly two approaches: characterization by attractors (embedding and bifurcation) and time-frequency, multiscale/multiresolution approaches. In real situations, the phase space associated to the acquisition of a complex phenomenon is unknown. It is however possible to relate, inside the signal's domain, local predictability to local reconstruction [13] and to deduce relevant information associated to multiscale geophysical signals [14]. A multiscale organization is a fundamental feature of a complex system, it can be for example related to the cascading properties in turbulent systems. We make use of this kind of description when analyzing turbulent signals: intermittency is observed within the inertial range and is related to the fact that, in the case of FDT (fully developed turbulence), symmetry is restored only in a statistical sense, a fact that has consequences on the quality of any nonlinear signal representation by frames or dictionaries.

The example of FDT as a standard "template" for developing general methods that apply to a vast class of complex systems and signals is of fundamental interest because, in FDT, the existence of a multiscale hierarchy $\mathcal{F}_h$ which is of multifractal nature and geometrically localized can be derived from physical considerations. This geometric hierarchy of sets is responsible for the shape of the computed singularity spectra, which in turn is related to the statistical organization of information content in a signal. It explains scale invariance, a characteristic feature of complex signals. The analogy from statistical physics comes from the fact that singularity exponents are direct generalizations of *critical exponents* which explain the macroscopic properties of a system around critical points, and the quantitative characterization of *universality classes*, which allow the definition of methods and algorithms that apply to general complex signals and systems, and not only turbulent signals: signals which belong to a same universality class share common statistical organization. During the past decades, canonical approaches permitted the development of a well-established analogy taken from thermodynamics in the analysis of complex signals: if $\mathcal{F}$ is the free energy, $\mathcal{T}$ the temperature measured in energy units, $\mathcal{U}$ the internal energy per volume unit $\mathcal{S}$ the entropy and $\widehat{\beta} = 1/\mathcal{T}$, then the scaling exponents associated to moments of intensive variables $p \to \tau_p$ corresponds to $\widehat{\beta}\mathcal{F}$, $\mathcal{U}(\widehat{\beta})$ corresponds to the singularity exponents values, and $\mathcal{S}(\mathcal{U})$ to the singularity spectrum [27]. The research goal is to be able to determine universality classes associated to acquired signals, independently of microscopic properties in the phase space of various complex systems, and beyond the particular case of turbulent data [53].

We show in figure 1  the result of the computation of singularity exponents on an *Herschel* astronomical observation map (the Musca galactic cloud) which has been edge-aware filtered using sparse $L^1$ filtering to eliminate the cosmic infrared background (or CIB), a type of noise that can modify the singularity spectrum of a signal.

## 3.3. Causal modeling

The team is working on a new class of models for modeling physical systems, starting from measured data and accounting for their dynamics [40]. The idea is to statistically describe the evolution of a system in terms of causally-equivalent states; states that lead to the same predictions [33]. Transitions between these states can be reconstructed from data, leading to a theoretically-optimal predictive model [63]. In practice, however, no algorithm is currently able to reconstruct these models from data in a reasonable time and without substantial
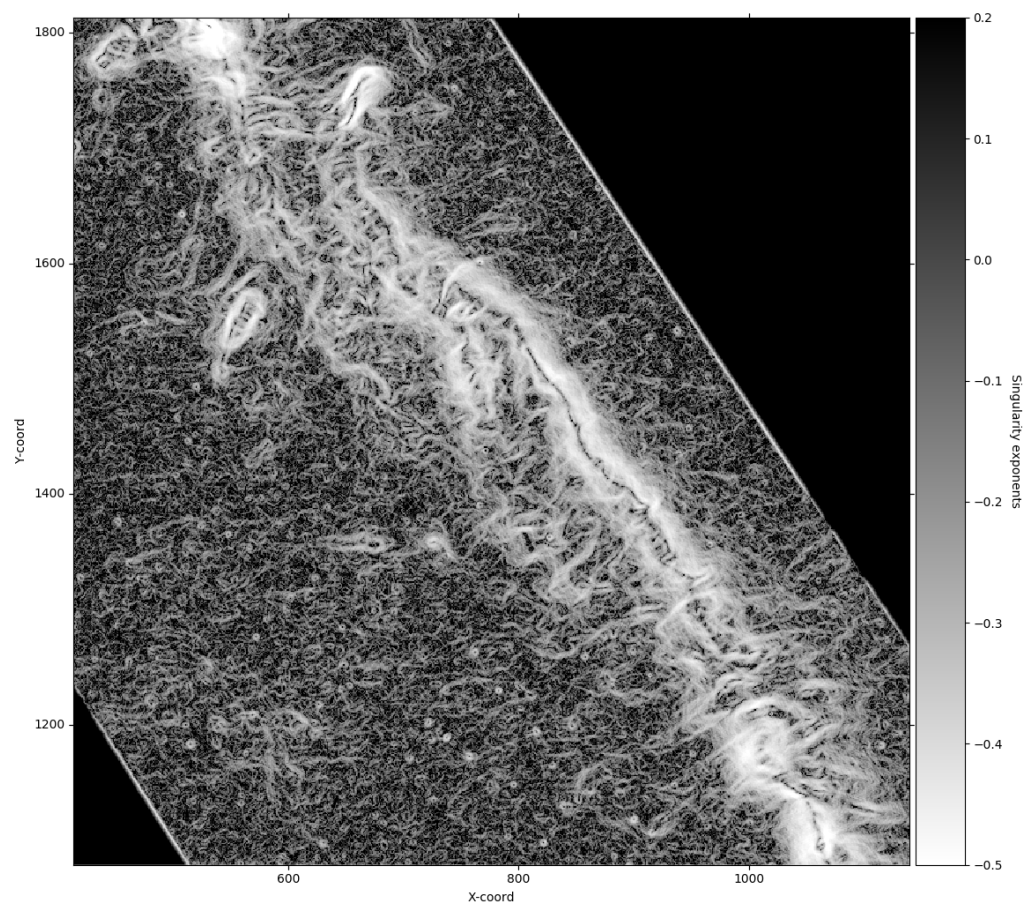
*Figure 1. Visualization of the singularity exponents, computed on a edge-aware filtered Musca Herschel observation map .*

discrete approximations. Recent progress now allows a continuous formulation of predictive causal models. Within this framework, more efficient algorithms may be found. The broadened class of predictive models promises a new perspective on structural complexity in many applications.

## 3.4. Speech analysis

Phonetic and sub-phonetic analysis: We developed a novel algorithm for automatic detection of Glottal Closure Instants (GCI) from speech signals using the Microcanonical Multiscale Formalism (MMF). This state of the art algorithm is considered as a reference in this field. We made a Matlab code implementing it available to the community (link). Our approach is based on the Microcanonical Multiscale Formalism. We showed that in the case of clean speech, our algorithm performs almost as well as a recent state-of-the-art method. In presence of different types of noises, we showed that our method is considerably more accurate (particularly for very low SNRs). Moreover, our method has lower computational times does not rely on an estimate of pitch period nor any critical choice of parameters. Using the same MMF, we also developed a method for phonetic segmentation of speech signal. We showed that this method outperforms state of the art ones in term of accuracy and efficiency.

Pathological speech analysis and classification: we made a critical analysis of some widely used methodologies in pathological speech classification. We then introduced some novel methods for extracting some common features used in pathological speech analysis and proposed more robust techniques for classification.

Speech analysis of patients with Parkinsonism: with our collaborators from the Czech Republic, we started preliminary studies of some machine learning issues in the field essentially due the small amount of training data.

## 3.5. Excitable systems: analysis of physiological time series

The research described in this section is a collaboration effort of GEOSTAT, CNRS LOMA (Laboratoire Ondes et Matière d'Aquitaine) and Laboratory of Physical Foundation of Strength, Institute of Continuous Media Mechanics (Perm, Russia Federation).

AF is an arrhythmia originating in the rapid and irregular electrical activity of the atria (the heart's two upper chambers) that causes their pump function to fail, increasing up to fivefold the risk of embolic stroke. The prevailing electrophysiological concepts describing tachy-arrhythmias are more than a century old. They involve abnormal automaticity and conduction [52]. Initiation and maintenance are thought to arise from a vulnerable substrate prone to the emergence of multiple self-perpetuating reentry circuits, also called "multiple wavelets" [59], [60]. Reentries may be driven structurally, for instance because of locally high fibrous tissue content which badly conducts, or functionally because of high spatial dispersion of decreased refractoriness and APD [58]. The latter is coined the leading circle concept with the clinically more relevant notion of a critical "wavelength" (in fact the length) of the cardiac impulse [26], [65], [62], [30]. The related concept of vulnerability was originally introduced to uncover a physiological substrate evolving from normality to pathology. It was found in vulnerable patients that high rate frequency would invariably lead to functional disorder as cardiac cells would no longer properly adapt their refractoriness [29]. Mathematical models have managed to exhibit likewise phenomena, with the generation of breaking spiral waves in various conditions [49], [54]. The triggering role of abnormal ectopic activity of the pulmonary veins has been demonstrated on patients with paroxysmal AF resistant to drug therapy [48], but its origin still remains poorly understood. This region is highly innervated with sympathetic and parasympathetic stimulation from the ANS [68], [69], [28]. In particular, Coumel et al. [39], [38] have revealed the pathophysiological role of the vagal tone on a vulnerable substrate. It is frequently observed that rapid tachycardia of ectopic origin transits to AF. This is known to result from electrical remodeling. As described for the first time by Allessie et al. [25], remodeling is a transient and reversible process by which the impulse properties such as its refractory period are altered during the course of the arrhythmia, promoting its perpetuation: "AF begets AF" [71]. Under substantial beating rate increase, cells may undergo remodeling to overcome the toxicity of their excessive intercellular calcium loading, by a rapid down regulation (a few minutes) of their L-type calcium membrane current. Moreover, other ionic channel

functions are also modified such as the potassium channel function, inducing a change in the conduction properties including the conduction velocity. The intercellular coupling at the gap junction level shows also alterations of their connexin expression and dispersion.

Wavelet-based methods (WTMM, log-cumulants, two point scale correlations), and confidence statistical methodology, have been applied to catheter recordings in the coronary sinus vein right next to the left atria of a small sample of patients with various conditions, and exhibit clear multifractal scaling without cross-scale correlation, which are coined "multifractal white noise", and that can be grouped according to two anatomical regions. One of our main result was to show that this is incompatible with the common lore for atrial fibrillation based on so-called circuit reentries. We used two declinations of a wavelet-based multiscale method, the moment (partition function) method and the magnitude cumulant method, as originally introduced in the field of fully developed turbulence. In the context of cardiac physiology, this methodology was shown to be valuable in assessing congestive heart failure from the monitoring of sinus heart rate variability [50]. We develop a model such that the substrate function is modulated by the kinetics of conduction. A simple reversible mechanism of short term remodeling under rapid pacing is demonstrated, by which ionic overload acts locally (dynamical feedback) on the kinetics of gap junction conductance. The whole process may propagate and pervade the myocardium via electronic currents, becoming desynchronized. In a new description, we propose that circuit reentries may well exist before the onset of fibrillation, favoring onset but not contributing directly to the onset and perpetuation. By contrast, cell-to-cell coupling is considered fundamentally dynamical. The rationale stems from the observation that multifractal scaling necessitates a high number of degrees of freedom (tending to infinity with system size), which can originate in excitable systems in hyperbolic spatial coupling.

## 3.6. Data-based identification of characteristic scales and automated modeling

Data are often acquired at the highest possible resolution, but that scale is not necessarily the best for modeling and understanding the system from which data was measured. The intrinsic properties of natural processes do not depend on the arbitrary scale at which data is acquired; yet, usual analysis techniques operate at the acquisition resolution. When several processes interact at different scales, the identification of their characteristic scales from empirical data becomes a necessary condition for properly modeling the system. A classical method for identifying characteristic scales is to look at the work done by the physical processes, the energy they dissipate over time. The assumption is that this work matches the most important action of each process on the studied natural system, which is usually a reasonable assumption. In the framework of time-frequency analysis [45], the power of the signal can be easily computed in each frequency band, itself matching a temporal scale.

However, in open and dissipating systems, energy dissipation is a prerequisite and thus not necessarily the most useful metric to investigate. In fact, most natural, physical and industrial systems we deal with fall in this category, while balanced quasi-static assumptions are practical approximation only for scales well below the characteristic scale of the involved processes. Open and dissipative systems are not locally constrained by the inevitable rise in entropy, thus allowing the maintaining through time of mesoscopic ordered structures. And, according to information theory [47], more order and less entropy means that these structures have a higher information content than the rest of the system, which usually gives them a high functional role.

We propose to identify characteristic scales not only with energy dissipation, as usual in signal processing analysis, but most importantly with information content. Information theory can be extended to look at which scales are most informative (e.g. multi-scale entropy [37], $\varepsilon$-entropy [36]). Complexity measures quantify the presence of structures in the signal (e.g. statistical complexity [42], MPR [56] and others [44]). With these notions, it is already possible to discriminate between random fluctuations and hidden order, such as in chaotic systems [41], [56]. The theory of how information and structures can be defined through scales is not complete yet, but the state of art is promising [43]. Current research in the team focuses on how informative scales can be found using collections of random paths, assumed to capture local structures as they reach out [35].

Building on these notions, it should also possible to fully automate the modeling of a natural system. Once characteristic scales are found, causal relationships can be established empirically. They are then clustered together in internal states of a special kind of Markov models called $\epsilon$-machines [42]. These are known to be the

optimal predictors of a system, with the drawback that it is currently quite complicated to build them properly, except for small system [64]. Recent extensions with advanced clustering techniques [34], [46], coupled with the physics of the studied system (e.g. fluid dynamics), have proved that $\epsilon$-machines are applicable to large systems, such as global wind patterns in the atmosphere [51]. Current research in the team focuses on the use of reproducing kernels, coupled possibly with sparse operators, in order to design better algorithms for $\epsilon$-machines reconstruction. In order to help with this long-term project, a collaboration is ongoing with J. Crutchfield lab at UC Davis.

<p align="center"><span style="color:red"><strong>MEMPHIS Project-Team</strong></span></p>

# 3. Research Program

## 3.1. Reduced-order models

Massive parallelization and rethinking of numerical schemes will allow the use of mathematical models for a broader class of physical problems. For industrial applications, there is an increasing need for rapid and reliable numerical simulators to tackle design and control tasks. To provide a concrete example, in the design process of an aircraft, the flight conditions and manoeuvres, which provide the largest aircraft loads, are not known *a priori*. Therefore, the aerodynamic and inertial forces are calculated for a large number of conditions to give an estimate of the maximum loads, and hence stresses, that the structure of the detailed aircraft design might experience in service. As a result, the number of simulations required for a realistic design problem could easily be in the order of tens of millions. Even with simplistic models of the aircraft behavior this is an unfeasible number of separate simulations. However, engineering experience is used to identify the most likely critical load conditions, meaning that approximately hundreds of thousands simulations are required for conventional aircraft configurations. Furthermore, these analyses have to be repeated every time that there is an update in the aircraft structure.

Compared to existing approaches for ROMs  [35], our interest will be focused on two axes. On the one hand, we start from the consideration that small, highly nonlinear scales are typically concentrated in limited spatial regions of the full simulation domain. So for example, in the flow past a wing, the highly non-linear phenomena take place in the proximity of the walls at the scale of a millimeter, for computational domains that are of the order of hundreds of meters. Based on these considerations, we propose in [31] a multi-scale model where the large scales are described by far-field models based on ROMs and the small scales are simulated by high-fidelity models. The whole point for this approach is to optimally decouple the far field from the near field.

A second characterizing feature of our ROM approach is non-linear interpolation. We start from the consideration that dynamical models derived from the projection of the PDE model in the reduced space are neither stable to numerical integration nor robust to parameter variation when hard non-linear multi-scale phenomena are considered.

However, thanks to Proper Orthogonal Decomposition (POD) [41], [47], [30] we can accurately approximate large solution databases using a low-dimensional base. Recent techniques to investigate the temporal evolution of the POD modes (Koopman modes  [42], [28], Dynamic Mode Decomposition  [45]) and allow a dynamic discrimination of the role played by each of them. This in turn can be exploited to interpolate between modes in parameter space, thanks to ideas relying on optimal transportation  [50], [32] that we have started developing in the FP7 project FFAST and H2020 AEROGUST.

## 3.2. Hierarchical Cartesian schemes

We intend to conceive schemes that will simplify the numerical approximation of problems involving complex unsteady objects together with multi-scale physical phenomena. Rather than using extremely optimized but non-scalable algorithms, we adopt robust alternatives that bypass the difficulties linked to grid generation. Even if the mesh problem can be tackled today thanks to powerful mesh generators, it still represents a severe difficulty, in particular when highly complex unsteady geometries need to be dealt with. Industrial experience and common practice shows that mesh generation accounts for about 20% of overall analysis time, whereas creation of a simulation-specific geometry requires about 60%, and only 20% of overall time is actually devoted to analysis. The methods that we develop bypass the generation of tedious geometrical models by automatic implicit geometry representation and hierarchical Cartesian schemes.

The approach that we plan to develop combines accurate enforcement of unfitted boundary conditions with adaptive octree and overset grids. The core idea is to use an octree/overset mesh for the approximation of the solution fields, while the geometry is captured by level set functions  [46], [40] and boundary conditions are imposed using appropriate interpolation methods  [27], [49], [44]. This eliminates the need for boundary-conforming meshes that require time-consuming and error-prone mesh generation procedures, and opens the door for simulation of very complex geometries. In particular, it will be possible to easily import the industrial geometry and to build the associated level set function used for simulation.

Hierarchical octree grids offer several considerable advantages over classical adaptive mesh refinement for body-fitted meshes, in terms of data management, memory footprint and parallel HPC performance. Typically, when refining unstructured grids, like for example tetrahedral grids, it is necessary to store the whole data tree corresponding to successive subdivisions of the elements and eventually recompute the full connectivity graph. In the linear octree case that we develop, only the tree leaves are stored in a linear array, with a considerable memory advantage. The mapping between the tree leaves and the linear array as well as the connectivity graph is efficiently computed thanks to an appropriate space-filling curve. Concerning parallelization, linear octrees guarantee a natural load balancing thanks to the linear data structure, whereas classical unstructured meshes require sophisticated (and moreover time consuming) tools to achieve proper load distribution (SCOTCH, METIS etc.). Of course, using unfitted hierarchical meshes requires further development and analysis of methods to handle the refinement at level jumps in a consistent and conservative way, accuracy analysis for new finite-volume or finite-difference schemes, efficient reconstructions at the boundaries to recover appropriate accuracy and robustness. These subjects, that are currently virtually absent at Inria, are among the main scientific challenges of our team.

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

# 3. Research Program

## 3.1. Introduction

integer programming, graph theory, decomposition approaches, polyhedral approaches, quadratic programming approaches, constraint programming.

*Combinatorial optimization* is the field of discrete optimization problems. In many applications, the most important decisions (control variables) are binary (on/off decisions) or integer (indivisible quantities). Extra variables can represent continuous adjustments or amounts. This results in models known as *mixed integer programs* (MIP), where the relationships between variables and input parameters are expressed as linear constraints and the goal is defined as a linear objective function. MIPs are notoriously difficult to solve: good quality estimations of the optimal value (bounds) are required to prune enumeration-based global-optimization algorithms whose complexity is exponential. In the standard approach to solving an MIP is so-called *branch-and-bound algorithm* : $(i)$ one solves the linear programming (LP) relaxation using the simplex method; $(ii)$ if the LP solution is not integer, one adds a disjunctive constraint on a factional component (rounding it up or down) that defines two sub-problems; $(iii)$ one applies this procedure recursively, thus defining a binary enumeration tree that can be pruned by comparing the local LP bound to the best known integer solution. Commercial MIP solvers are essentially based on branch-and-bound (such IBM-CPLEX, FICO-Xpress-mp, or GUROBI). They have made tremendous progress over the last decade (with a speedup by a factor of 60). But extending their capabilities remains a continuous challenge; given the combinatorial explosion inherent to enumerative solution techniques, they remain quickly overwhelmed beyond a certain problem size or complexity.

Progress can be expected from the development of tighter formulations. Central to our field is the characterization of polyhedra defining or approximating the solution set and combinatorial algorithms to identify "efficiently" a minimum cost solution or separate an unfeasible point. With properly chosen formulations, exact optimization tools can be competitive with other methods (such as meta-heuristics) in constructing good approximate solutions within limited computational time, and of course has the important advantage of being able to provide a performance guarantee through the relaxation bounds. Decomposition techniques are implicitly leading to better problem formulation as well, while constraint propagation are tools from artificial intelligence to further improve formulation through intensive preprocessing. A new trend is robust optimization where recent progress have been made: the aim is to produce optimized solutions that remain of good quality even if the problem data has stochastic variations. In all cases, the study of specific models and challenging industrial applications is quite relevant because developments made into a specific context can become generic tools over time and see their way into commercial software.

Our project brings together researchers with expertise in mathematical programming (polyhedral approaches, decomposition and reformulation techniques in mixed integer programing, robust and stochastic programming, and dynamic programming), graph theory (characterization of graph properties, combinatorial algorithms) and constraint programming in the aim of producing better quality formulations and developing new methods to exploit these formulations. These new results are then applied to find high quality solutions for practical combinatorial problems such as routing, network design, planning, scheduling, cutting and packing problems, High Performance and Cloud Computing.

## 3.2. Polyhedral approaches for MIP

Adding valid inequalities to the polyhedral description of an MIP allows one to improve the resulting LP bound and hence to better prune the enumeration tree. In a cutting plane procedure, one attempt to identify valid inequalities that are violated by the LP solution of the current formulation and adds them to the formulation. This can be done at each node of the branch-and-bound tree giving rise to a so-called

*branch-and-cut algorithm* [65]. The goal is to reduce the resolution of an integer program to that of a linear program by deriving a linear description of the convex hull of the feasible solutions. Polyhedral theory tells us that if $X$ is a mixed integer program: $X = P \cap \mathbb{Z}^n \times \mathbb{R}^p$ where $P = \{x \in \mathbb{R}^{n+p} : Ax \leq b\}$ with matrix $(A, b) \in \mathbb{Q}^{m \times (n+p+1)}$, then $conv(X)$ is a polyhedron that can be described in terms of linear constraints, i.e. it writes as $conv(X) = \{x \in \mathbb{R}^{n+p} : C\,x \leq d\}$ for some matrix $(C, d) \in \mathbb{Q}^{m' \times (n+p+1)}$ although the dimension $m'$ is typically quite large. A fundamental result in this field is the equivalence of complexity between solving the combinatorial optimization problem $\min\{cx : x \in X\}$ and solving the *separation problem* over the associated polyhedron $conv(X)$: if $\widetilde{x} \notin conv(X)$, find a linear inequality $\pi\,x \geq \pi_0$ satisfied by all points in $conv(X)$ but violated by $\widetilde{x}$. Hence, for NP-hard problems, one can not hope to get a compact description of $conv(X)$ nor a polynomial time exact separation routine. Polyhedral studies focus on identifying some of the inequalities that are involved in the polyhedral description of $conv(X)$ and derive efficient *separation procedures* (cutting plane generation). Only a subset of the inequalities $C\,x \leq d$ can offer a good approximation, that combined with a branch-and-bound enumeration techniques permits to solve the problem. Using *cutting plane algorithm* at each node of the branch-and-bound tree, gives rise to the algorithm called *branch-and-cut*.

## 3.3. Decomposition-and-reformulation-approaches

An hierarchical approach to tackle complex combinatorial problems consists in considering separately different substructures (subproblems). If one is able to implement relatively efficient optimization on the substructures, this can be exploited to reformulate the global problem as a selection of specific subproblem solutions that together form a global solution. If the subproblems correspond to subset of constraints in the MIP formulation, this leads to Dantzig-Wolfe decomposition. If it corresponds to isolating a subset of decision variables, this leads to Bender's decomposition. Both lead to extended formulations of the problem with either a huge number of variables or constraints. Dantzig-Wolfe approach requires specific algorithmic approaches to generate subproblem solutions and associated global decision variables dynamically in the course of the optimization. This procedure is known as *column generation*, while its combination with branch-and-bound enumeration is called *branch-and-price*. Alternatively, in Bender's approach, when dealing with exponentially many constraints in the reformulation, the *cutting plane procedures* that we defined in the previous section are well-suited tools. When optimization on a substructure is (relatively) easy, there often exists a tight reformulation of this substructure typically in an extended variable space. This gives rise powerful reformulation of the global problem, although it might be impractical given its size (typically pseudo-polynomial). It can be possible to project (part of) the extended formulation in a smaller dimensional space if not the original variable space to bring polyhedral insight (cuts derived through polyhedral studies can often be recovered through such projections).

## 3.4. Integration of Artificial Intelligence Techniques in Integer Programming

When one deals with combinatorial problems with a large number of integer variables, or tightly constrained problems, mixed integer programming (MIP) alone may not be able to find solutions in a reasonable amount of time. In this case, techniques from artificial intelligence can be used to improve these methods. In particular, we use variable fixing techniques, primal heuristics and constraint programming.

Primal heuristics are useful to find feasible solutions in a small amount of time. We focus on heuristics that are either based on integer programming (rounding, diving, relaxation induced neighborhood search, feasibility pump), or that are used inside our exact methods (heuristics for separation or pricing subproblem, heuristic constraint propagation, ...). Such methods are likely to produce good quality solutions only if the integer programming formulation is of top quality, i.e., if its LP relaxation provides a good approximation of the IP solution.

In the same line, variable fixing techniques, that are essential in reducing the size of large scale problems, rely on good quality approximations: either tight formulations or tight relaxation solvers (as a dynamic program combined with state space relaxation). Then if the dual bound derives when the variable is fixed to one exceeds the incubent solution value, the variable can be fixed to zero and hence removed from the problem. The process can be apply sequentially by refining the degree of relaxation.

Constraint Programming (CP) focuses on iteratively reducing the variable domains (sets of feasible values) by applying logical and problem-specific operators. The latter propagates on selected variables the restrictions that are implied by the other variable domains through the relations between variables that are defined by the constraints of the problem. Combined with enumeration, it gives rise to exact optimization algorithms. A CP approach is particularly effective for tightly constrained problems, feasibility problems and min-max problems. Mixed Integer Programming (MIP), on the other hand, is known to be effective for loosely constrained problems and for problems with an objective function defined as the weighted sum of variables. Many problems belong to the intersection of these two classes. For such problems, it is reasonable to use algorithms that exploit complementary strengths of Constraint Programming and Mixed Integer Programming.

## 3.5. Robust Optimization

Decision makers are usually facing several sources of uncertainty, such as the variability in time or estimation errors. A simplistic way to handle these uncertainties is to overestimate the unknown parameters. However, this results in over-conservatism and a significant waste in resource consumption. A better approach is to account for the uncertainty directly into the decision aid model by considering mixed integer programs that involve uncertain parameters. Stochastic optimization account for the expected realization of random data and optimize an expected value representing the average situation. Robust optimization on the other hand entails protecting against the worst-case behavior of unknown data. There is an analogy to game theory where one considers an oblivious adversary choosing the realization that harms the solution the most. A full worst case protection against uncertainty is too conservative and induces very high over-cost. Instead, the realization of random data are bound to belong to a restricted feasibility set, the so-called uncertainty set. Stochastic and robust optimization rely on very large scale programs where probabilistic scenarios are enumerated. There is hope of a tractable solution for realistic size problems, provided one develops very efficient ad-hoc algorithms. The techniques for dynamically handling variables and constraints (column-and-row generation and Bender's projection tools) that are at the core of our team methodological work are specially well-suited to this context.

## 3.6. Approximation Algorithms

In some contexts, obtaining an exact solution to an optimization problem is not feasible: when instances are too large, or when decisions need to be taken rapidly. Since most of the combinatorial optimization problems are NP-hard, another direction to obtain good quality solutions in reasonable time is to focus on **approximation algorithms**. The definition of approximation algorithms is based on the notion of input set $\mathcal{I}$ and each $I \in \mathcal{I}$ defines a solution space $\mathcal{S}_I$. For a minimization problem $\min_{x \in \mathcal{S}_I} f(x)$, an algorithm $\mathcal{A}$ is an $\alpha$-approximation algorithm if it provides a solution within $\alpha$ of the optimal solution for all instances in the input set:

$$\forall I \in \mathcal{I}, \quad f(\mathcal{A}(I)) \leq \alpha \min_{x \in \mathcal{S}_I} f(x) = f^*(I)$$

The objective is to search for polynomial algorithms, with approximation ratios as close to 1 as possible. Such algorithms are called *worst-case* approximation algorithms, because the performance guarantee is expressed over all possible inputs of the problem. The design of these algorithms have strong links with the enumeration techniques described above: since computing $f^*(I)$ is an NP-hard problem, it is often required to derive **strong a priori bounds** on the optimal solution value which can afterward be compared to estimations of the value of the solution produced. In many cases, it is also possible to build $\alpha$-approximate solutions by a careful rounding of a solution obtained from the linear relaxation of an integer formulation of the problem. Members of the team have expertise in designing and evaluating approximation algorithms for resource allocation in computer systems, using a variety of techniques, such as dual approximation (where a guess of the optimal value $f^*$ is provided, and $\mathcal{A}$ either provides a solution within $\alpha f^*$, or guarantees that no solution of value $f^*$ or less exists), or resource augmentation (where an approximation is obtained by relaxing some of the constraints of the problem).

## 3.7. Polyhedral Combinatorics and Graph Theory

Many fundamental combinatorial optimization problems can be modeled as the search for a specific structure in a graph. For example, ensuring connectivity in a network amounts to building a *tree* that spans all the nodes. Inquiring about its resistance to failure amounts to searching for a minimum cardinality *cut* that partitions the graph. Selecting disjoint pairs of objects is represented by a so-called *matching*. Disjunctive choices can be modeled by edges in a so-called *conflict graph* where one searches for *stable sets* – a set of nodes that are not incident to one another. Polyhedral combinatorics is the study of combinatorial algorithms involving polyhedral considerations. Not only it leads to efficient algorithms, but also, conversely, efficient algorithms often imply polyhedral characterizations and related min-max relations. Developments of polyhedral properties of a fundamental problem will typically provide us with more interesting inequalities well suited for a branch-and-cut algorithm to more general problems. Furthermore, one can use the fundamental problems as new building bricks to decompose the more general problem at hand. For problem that let themselves easily be formulated in a graph setting, the graph theory and in particular graph decomposition theorem might help.

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

# 3. Research Program

## 3.1. Complex models for the propagation of cardiac action potentials

The contraction of the heart is coordinated by a complex electrical activation process which relies on about a million ion channels, pumps, and exchangers of various kinds in the membrane of each cardiac cell. Their interaction results in a periodic change in transmembrane potential called an action potential. Action potentials in the cardiac muscle propagate rapidly from cell to cell, synchronizing the contraction of the entire muscle to achieve an efficient pump function. The spatio-temporal pattern of this propagation is related both to the function of the cellular membrane and to the structural organization of the cells into tissues. Cardiac arrythmias originate from malfunctions in this process. The field of cardiac electrophysiology studies the multiscale organization of the cardiac activation process from the subcellular scale up to the scale of the body. It relates the molecular processes in the cell membranes to the propagation process and to measurable signals in the heart and to the electrocardiogram, an electrical signal on the torso surface.

Several improvements of current models of the propagation of the action potential are being developed in the Carmen team, based on previous work [56] and on the data available at IHU LIRYC:

- Enrichment of the current monodomain and bidomain models [56], [67] by accounting for structural heterogeneities of the tissue at an intermediate scale. Here we focus on multiscale analysis techniques applied to the various high-resolution structural data available at the LIRYC.

- Coupling of the tissues from the different cardiac compartments and conduction systems. Here, we develop models that couple 1D, 2D and 3D phenomena described by reaction-diffusion PDEs.

These models are essential to improve our in-depth understanding of cardiac electrical dysfunction. To this aim, we use high-performance computing techniques in order to numerically explore the complexity of these models.

We use these model codes for applied studies in two important areas of cardiac electrophysiology: atrial fibrillation [60] and sudden-cardiac-death (SCD) syndromes [7], [6] [64]. This work is performed in collaboration with several physiologists and clinicians both at IHU Liryc and abroad.

## 3.2. Simplified models and inverse problems

The medical and clinical exploration of the cardiac electric signals is based on accurate reconstruction of the patterns of propagation of the action potential. The correct detection of these complex patterns by non-invasive electrical imaging techniques has to be developed. This problem involves solving inverse problems that cannot be addressed with the more compex models. We want both to develop simple and fast models of the propagation of cardiac action potentials and improve the solutions to the inverse problems found in cardiac electrical imaging techniques.

The cardiac inverse problem consists in finding the cardiac activation maps or, more generally, the whole cardiac electrical activity, from high-density body surface electrocardiograms. It is a new and a powerful diagnosis technique, which success would be considered as a breakthrough. Although widely studied recently, it remains a challenge for the scientific community. In many cases the quality of reconstructed electrical potential is not adequate. The methods used consist in solving the Laplace equation on the volume delimited by the body surface and the epicardial surface. Our aim is to

- study in depth the dependance of this inverse problem on inhomogeneities in the torso, conductivity values, the geometry, electrode positions, etc., and

- improve the solution to the inverse problem by using new regularization strategies, factorization of boundary value problems, and the theory of optimal control.

Of course we will use our models as a basis to regularize these inverse problems. We will consider the following strategies:

- using complete propagation models in the inverse problem, like the bidomain equations, for instance in order to localize electrical sources;
- constructing families of reduced-order models using e.g. statistical learning techniques, which would accurately represent some families of well-identified pathologies; and
- constructing simple models of the propagation of the activation front, based on eikonal or level-set equations, but which would incorporate the representation of complex activation patterns.

Additionaly, we will need to develop numerical techniques dedicated to our simplified eikonal/level-set equations.

## 3.3. Numerical techniques

We want our numerical simulations to be efficient, accurate, and reliable with respect to the needs of the medical community. Based on previous work on solving the monodomain and bidomain equations [4], [5], [8], [1], we will focus on

- High-order numerical techniques with respect to the variables with physiological meaning, like velocity, AP duration and restitution properties.
- Efficient, dedicated preconditioning techniques coupled with parallel computing.

Existing simulation tools used in our team rely, among others, on mixtures of explicit and implicit integration methods for ODEs, hybrid MPI-OpenMP parallellization, algebraic multigrid preconditioning, and Krylov solvers. New developments include high-order explicit integration methods and task-based dynamic parallellism.

## 3.4. Cardiac Electrophysiology at the Microscopic Scale

Numerical models of whole-heart physiology are based on the approximation of a perfect muscle using homogenisation methods. However, due to aging and cardiomyopathies, the cellular structure of the tissue changes. These modifications can give rise to life-threatening arrhythmias. For our research on this subject and with cardiologists of the IHU LIRYC Bordeaux, we aim to design and implement models that describe the strong heterogeneity of the tissue at the cellular level and to numerically explore the mechanisms of these diseases.

The literature on this type of model is still very limited [74]. Existing models are two-dimensional [65] or limited to idealized geometries, and use a linear (purely resistive) behaviour of the gap-juction channels that connect the cells. We propose a three-dimensional approach using realistic cellular geometry (figure 1 ), nonlinear gap-junction behaviour, and a numerical approach that can scale to hundreds of cells while maintaining a sub-micrometer spatial resolution (10 to 100 times smaller than the size of a cardiomyocyte) [52], [51], [49]. P-E. Bécue defended his PhD thesis on this topic in December 2018.

A                                              B                                              C

*Figure 1. **A:** The cardiac muscle consists of a branching network of elongated muscle cells, interspersed with other structures. Sheets of connective tissue (blue) can grow between the muscle cells and become pathogenic. **B:** Current models can only represent such alterations in a coarse way by replacing model elements with different types; each cube in this illustration would represent hundreds of cells. **C:** This hand-crafted example illustrates the type of geometric model we are experimenting with. Each cell is here represented by hundreds of elements.*

# 3. Research Program

## 3.1. Introduction

Probing the invisible is a quest that is shared by a wide variety of scientists such as archaeologists, geologists, astrophysicists, physicists, etc... Magique-3D is mainly involved in Geophysical imaging which aims at understanding the internal structure of the Earth from the propagation of waves. Both qualitative and quantitative information are required and two geophysical techniques can be used: **seismic reflection** and **seismic inversion**. Seismic reflection provides a qualitative description of the subsurface from reflected seismic waves by indicating the position of the reflectors while seismic inversion transforms seismic reflection data into a quantitative description of the subsurface. Both techniques are inverse problems based upon the numerical solution of wave equations. Oil and Gas explorations have been pioneering application domains for seismic reflection and inversion and even if numerical seismic imaging is computationally intensive, oil companies clearly promote the use of numerical simulations to provide synthetic maps of the subsurface. This is due to the tremendous progresses of scientific computing which have pushed the limits of existing numerical methods and it is now conceivable to tackle realistic 3D problems. However, mathematical wave modeling has to be well-adapted to the region of interest and the numerical schemes which are employed to solve wave equations have to be both accurate and scalable enough to take full advantage of parallel computing. Today, geophysical imaging tackles more and more realistic problems and we contribute to this task by improving the modeling and by deriving advanced numerical methods for solving wave problems.

MAGIQUE-3D research program is divided into four axes that are: (1) Imaging the Earth; (2) Exploring the Sun; (3) Detecting defaults in complex media; (4) Designing objects with a variety of shapes. Those applications stand out from the collaborations that we have established with interested end-user groups. It is worth noting that they share basic common methodologies which imparts consistency to our program despite they may appear quite distant. MAGIQUE-3D keep modeling and simulating geophysical phenomena for understanding the Earth interior and developing its resources sustainably, and our xperience with numerical geophysics may help us to address other challenging applications. We mainly used DG finite elements and spectral elements and both have demonstrated very good performance. However, in particular for reducing the computational costs and/or for better capturing the propagation characteristics, we are working on the development of hybrid solvers based on the coupling of different finite element methods. Other open questions deserve attention like the problem of numerical pollution or the poor scalability of decomposition domain techniques which are both significantly hampering computations in very large domains. For those purposes, we focus ourselves on the development of Trefftz-like approximations that are based on a particular computation of the fluxes by making a judicious use of an auxiliary numerical method (e.g. boundary integral equations, spectral elements, etc...). Those problems cannot be ignored and can be found in all of our research axes. In addressing those issues, we participate in the construction of new numerical schemes and for that purpose, we continuously need to improve our understanding of the underlying physics. By this way, we make our mathematical models evolve to more realistic representations of the wave propagation phenomenon. This motivated us to introduce experimental studies in our activities and to collaborate with geophysicists of the UPPA who own experimental devices adapted to our concerns. Moreover, we have hired recently Yder Masson who is an experienced researcher developing modeling and imaging methods to investigate the Earth's internal structure. This creates all the conditions for improving our mathematical representation of waves in complex media. It is worth noting that modeling is a concern for both geophysicists and mathematicians. Indeed, the Physics must be reproduced accurately and the underlying mathematical properties should be clarified. By this way, we can develop a numerical scheme respecting the main properties of the continuous problem of interest (energy conservation or attenuation, stability, well-posedness, etc...). Magique-3D proposes to define its research program from in-house accurate solution methodologies for simulating wave propagation in realistic scenarios to various applications involving trans-disciplinary efforts. The development of high-order

numerical methods for wave simulations is serving as a basis for our contributions regarding applications. In particular, we pursue and strengthen our collaboration with HPC teams, in order to improve the scalability of our codes and to run them on very large heterogeneous architectures (using task based programming libraries as StarPU developed by Inria project-team Storm, improving the I/O by collaborating with UFRGS at Porto Alegre, using the metaprogramming framework Boast developed by Inria project-team Corse to produce portable and efficient computing kernels). We are also continuing our collaboration with Inria project team Hiepacs on the use of hybrid linear solvers, by considering the multiple Right-Hand Sides feature and by integrating appropriate transmission conditions between the various domains. During 2019, we have worked a lot on: (a) High-order numerical methods for modeling wave propagation in porous media: development and implementation; (b) Understanding the interior of the Earth and the Sun by solving inverse problems; (c) Full waveform inversion for the optimal design of wind musical instruments.

## 3.2. High-order numerical methods for modeling wave propagation in porous media: development and implementation

We aim at achieving the characterization of conducting porous media which are media favoring the conversion of seismic waves into electromagnetic waves.. This project is identified as a "New scientific challenge" which is a set of research projects funded by the E2S project of UPPA. The shape and form of porous media can vary depending on the size of the pore and the structure of the solid skeleton. Porous media are found in the nature (sandstone, volcanic rocks, etc) or can be manufactured (concrete, polyurethane foam, etc) as depicted in [68]. Instead of modeling such media as strongly heterogeneous, homogenization is used to describe the material on a macroscopic scale. Biot's theory describes the solid skeleton according to linear elasticity and adds to this the Navier-Stokes equation for a viscous fluid and Darcy's law governing the motion of the fluid  [63], [61]. For simplified linear elasticity, there are one equation of motion and one constitutive law, with the unknowns being the displacement field in the solid and the solid stress. In poro-elasticity, the added unknowns are the fluid displacement relative to the solid and the fluid pressure. There are two equations of motion, coupled with two constitutive laws. By plane wave analysis, one obtains three types of waves: S wave, fast P wave and slow P wave (Biot's wave). While the first two types are similar to those existing in elastic solid, the existence of a third wave with drastically smaller speed adds to the complications already encountered in elasticity. This is obviously even more challenging for conducting poroelastic media where the three poroelastic waves are coupled with an electric field. In this case, it is not realistic to use a unique scheme for all the waves. Standard finite element methods coupled with time schemes have indeed difficulties to deliver accurate solutions because there is a need of adapting the mesh size to the smallest wave velocity and the time discretization to the largest wave velocity. It is then tricky to numerically reproduce the Biot's wave while approximating correctly the regular elastic waves P and S. Moreover, there is a challenging question about the boundary condition to be used for limiting the computational domain. We have launched a Ph.D project (Rose-Cloé Meyer) aiming at developing a new piece of software for the simulation of time-harmonic waves in conducting porous media. This project is developed in collaboration with Steve Pride from the Lawrence Berkeley National Laboratory who has elaborated the corresponding physical theory [83], [87], [73], [74]. Next, once a new numerical method is developed, it is validated by comparing the numerical solution to an analytical one. This is a key step to us for assessing the accuracy of our simulations. Nevertheless, analytical solutions are not available for realistic media such as poroelastic or viscoelastic media represented by heterogeneous parameters. Engineers still argue that simulations may be inaccurate and could lead to wrong conclusions. Fortunately, it is possible to produce experimentally quite complex configurations where multi-physics measurements are used to monitor the wave propagation. There is thus a possibility of moving further on the validation of the numerical methods by comparing simulations and experiments. What is very exciting is that experiments are used to validate numerical methods which have the objective of simulating new phenomena that are not possible to reproduce in a lab. We have launched two Ph.D thesis (Chengyi Shen and Victor Martins Gomes) in collaboration with Daniel Brito (LFCR-UPPA) on the comparison of simulations with experiments. This topic is connected to another project that we have with Total on the use of waves for characterizing carbonates.

## 3.3. Understanding the interior of the Earth and the Sun by solving inverse problems

Even if the Earth and the Sun are actually very different media, their imaging is based on the same solution methodology [64]. However, our knowledge on Earth inversion is far more developed than for the Sun. Earth inversion is in the continuation of previous MAGIQUE-3D achievements while Sun inversion requires developing new technologies based on modeling, numerical analysis and implementation of a piece of software which is able to ask for new developments. For instance, we would like to develop a HDG software package for solving Galbrun and Linearized Euler equations. To the best of our knowledge, this has never be done and would be a major milestone for tackling vectorial equations. Regarding the modeling, we are pursuing our collaboration with the Max Planck Institute for Solar System Research (Göttingen, Germany) in the framework of the associate team ANTS. This partnership is essential to us for understanding a complex (and new to us) physics including gravity waves that we have never considered in the past. Even if we dispose of advanced solvers dealing with elasticity, the development of fast and accurate solvers for reproducing waves travelling in large 3D domains is still one of the positive developments towards realistic simulations. In particular, the techniques for the forward discretization and linear system solver must evolve accordingly to resolve large scale time-harmonic problems. For instance, we have elaborated a space-time Trefftz-DG formulation of the elasto-acoustic problem [58], which performs very well regarding the number of dofs and the order of convergence. We have also coupled spectral and DG elements to take advantage of both methods and we have performed some simulations which are very promising [57]. The formulation of FWI is in progress in the framework of Pierre Jacquet thesis launched in November 2017. Finally, we have also initiated research on seismology at the planetary scale, with the arrival of Yder Masson on the subject and new collaborators (such as Berkeley lab). This will further help widen our expertise on inverse wave problems and will feed all the four research axes of the future team-project. Regarding industrial partnerships, we have collaboration with Total and the SME RealtimeSeismic (Pau, France). We also continue to work with the UPV, the BCAM and the BSC, namely in the framework of Mathrocks project.

## 3.4. Hybrid time discretizations of high-order

Most of the meshes we consider are composed of cells greatly varying in size. This can be due to the physical characteristics (propagation speed, topography, ...) which may require to refine the mesh locally, very unstructured meshes can also be the result of dysfunction of the mesher. For practical reasons which are essentially guided by the aim of reducing the number of matrix inversions, explicit schemes are generally privileged. However, they work under a stability condition, the so-called Courant Friedrichs Lewy (CFL) condition which forces the time step being proportional to the size of the smallest cell. Then, it is necessary to perform a huge number of iterations in time and in most of the cases because of a very few number of small cells. This implies to apply a very small time step on grids mainly composed of coarse cells and thus, there is a risk of creating numerical dispersion that should not exist. However, this drawback can be avoided by using low degree polynomial basis in space in the small meshes and high degree polynomials in the coarse meshes. By this way, it is possible to relax the CFL condition and in the same time, the dispersion effects are limited. Unfortunately, the cell-size variations are so important that this strategy is not sufficient. One solution could be to apply implicit and unconditionally stable schemes, which would obviously free us from the CFL constraint. Unfortunately, these schemes require inverting a linear system at each iteration and thus needs huge computational burden that can be prohibitive in 3D. Moreover, numerical dispersion may be increased. Then, as second solution is the use of local time stepping strategies for matching the time step to the different sizes of the mesh. There are several attempts [65], [60], [82], [79], [71] and Magique 3D has proposed a new time stepping method which allows us to adapt both the time step and the order of time approximation to the size of the cells. Nevertheless, despite a very good performance assessment in academic configurations, we have observed to our detriment that its implementation inside industrial codes is not obvious and in practice, improvements of the computational costs are disappointing, especially in a HPC framework. Indeed, the local time stepping algorithm may strongly affect the scalability of the code. Moreover, the complexity of the algorithm is increased when dealing with lossy media [76].

Recently, Dolean *et al*  [70] have considered a novel approach consisting in applying hybrid schemes combining second order implicit schemes in the thin cells and second order explicit discretization in the coarse mesh. Their numerical results indicate that this method could be a good alternative but the numerical dispersion is still present. It would then be interesting to implement this idea with high-order time schemes to reduce the numerical dispersion. The recent arrival in the team of J. Chabassier should help us to address this problem since she has the expertise in constructing high-order implicit time scheme based on energy preserving Newmark schemes  [62]. We propose that our work be organized around the two following tasks. The first one is the extension of these schemes to the case of lossy media because applying existing schemes when there is attenuation is not straightforward. This is a key issue because there is artificial attenuation when absorbing boundary conditions are introduced and if not, there are cases with natural attenuation like in visco-elastic media. The second one is the coupling of high-order implicit schemes with high-order explicit schemes. These two tasks can be first completed independently, but the ultimate goal is obviously to couple the schemes for lossy media. We will consider two strategies for the coupling. The first one will be based on the method proposed by Dolean *et al*, the second one will consist in using Lagrange multiplier on the interface between the coarse and fine grids and write a novel coupling condition that ensures the high order consistency of the global scheme. Besides these theoretical aspects, we will have to implement the method in industrial codes and our discretization methodology is very suitable for parallel computing since it involves Lagrange multipliers. We propose to organize this task as follows. There is first the crucial issue of a systematic distribution of the cells in the coarse/explicit and in the fine/implicit part. Based on our experience on local time stepping, we claim that it is necessary to define a criterion which discriminates thin cells from coarse ones. Indeed, we intend to develop codes which will be used by practitioners, in particular engineers working in the production department of Total. It implies that the code will be used by people who are not necessarily experts in scientific computing. Considering real-world problems means that the mesh will most probably be composed of a more or less high number of subsets arbitrarily distributed and containing thin or coarse cells. Moreover, in the prospect of solving inverse problems, it is difficult to assess which cells are thin or not in a mesh which varies at each iteration.

Another important issue is the load balancing that we can not avoid with parallel computing. In particular, we will have to choose one of these two alternatives: dedicate one part of processors to the implicit computations and the other one to explicit calculus or distribute the resolution with both schemes on all processors. A collaboration with experts in HPC is then mandatory since we are not expert in parallel computing. We will thus continue to collaborate with the team-projects Hiepacs and Runtime with whom we have a long-term experience of collaborations. The load-balancing leads then to the issue of mesh partitioning. Main mesh partitioners are very efficient for the coupling of different discretizations in space but to the best of our knowledge, the case of non-uniform time discretization has never been addressed. The study of meshes being out of the scopes of Magique-3D, we will collaborate with experts on mesh partitioning. We get already on to François Pellegrini who is the principal investigator of Scotch (http://www.labri.fr/perso/pelegrin/scotch) and permanent member of the team project Bacchus (Inria Bordeaux Sud Ouest Research Center).

In the future, we aim at enlarging the application range of implicit schemes. The idea will be to use the degrees of freedom offered by the implicit discretization in order to tackle specific difficulties that may appear in some systems. For instance, in systems involving several waves (as P and S waves in porous elastic media, or coupled wave problems as previously mentioned) the implicit parameter could be adapted to each wave and optimized in order to reduce the computational cost. More generally, we aim at reducing numeric bottlenecks by adapting the implicit discretization to specific cases.

## 3.5. Full waveform inversion for the optimal design of wind musical instruments

Makers have improved wind musical instruments (as flutes, trumpets, clarinets, bassoons, ...) in the past by a "trial and error" procedure, where the final sound and ease of the instrument in playing conditions are the main criteria. Although the playing context should still be the final reference, we can consider intermediate measurements of the pipe entry impedance [75], [69], which quantifies the Dirichlet-to-Neumann map of the

wave propagation in the pipe, and relies on mathematical simulations based on accurate and concise models of the pipe [84], [59] and the embouchure [77], [59], [55], [56], [86] in order to foresee the behavior of a given instrument, and therefore optimize it. A strong interaction with makers and players is necessary for defining both operable criteria quantified as a cost function and a design parameters space. We aim at building efficient musical instrument via handcrafted techniques but also modern tools as additive synthesis (3D printers). We plan to implement state-of-the-art numerical methods (finite elements, full waveform inversion, neuronal networks fed by numerical simulations, diverse optimization techniques...) that are versatile (in terms of models, formulations, couplings...) in order to solve the optimization problem, after a proper modeling of the linear and nonlinear coupled phenomena. We wish to take advantage of the fact that sound waves in musical instruments satisfy the laws of acoustics in pipes (PDEs), which leads to use FWI technique, in harmonic or temporal regime. We propose to implement an iterative process between instrument making and optimal design in order to build instruments that optimize tone quality and playability. We are currently collaborating with musical acoustics teams who have a strong experimental background on this question [66], [67] [DCY12, DF07, GPP98], we wish to strenghten the links we have with other teams [78], [72], [80], [81], [85], we will participate to professional clusters [ITE], and we are currently collaborating with makers and museums directly : Augustin Humeau (Dordogne) for the bassoon, Luc Gallois (Oise) and the Museum of Cité de la Musique - Philharmonie de Paris for the brass instruments. This research axis is surely the most exploratory of our research program and follows the successful "Exploratory Research Program" Inria grant obtained in 2017. It could pave the way for significant progresses in inverse problem solving. Indeed, the problem depends on a few number of parameters unlike geophysical or astrophysical problems. We can thus use it to test different methods like neuronal networks, statistical methods, coupling with nonlinear phenomena, and decide if it could be applied to large scale applications.

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

# 3. Research Program

## 3.1. Integrative and Cognitive Neuroscience

The human brain is often considered as the most complex system dedicated to information processing. This multi-scale complexity, described from the metabolic to the network level, is particularly studied in integrative neuroscience, the goal of which is to explain how cognitive functions (ranging from sensorimotor coordination to executive functions) emerge from (are the result of the interaction of) distributed and adaptive computations of processing units, displayed along neural structures and information flows. Indeed, beyond the astounding complexity reported in physiological studies, integrative neuroscience aims at extracting, in simplifying models, regularities at various levels of description. From a mesoscopic point of view, most neuronal structures (and particularly some of primary importance like the cortex, cerebellum, striatum, hippocampus) can be described through a regular organization of information flows and homogenous learning rules, whatever the nature of the processed information. From a macroscopic point of view, the arrangement in space of neuronal structures within the cerebral architecture also obeys a functional logic, the sketch of which is captured in models describing the main information flows in the brain, the corresponding loops built in interaction with the external and internal (bodily and hormonal) world and the developmental steps leading to the acquisition of elementary sensorimotor skills up to the most complex executive functions.

In summary, integrative neuroscience builds, on an overwhelming quantity of data, a simplifying and interpretative grid suggesting homogenous local computations and a structured and logical plan for the development of cognitive functions. They arise from interactions and information exchange between neuronal structures and the external and internal world and also within the network of structures.

This domain is today very active and stimulating because it proposes, of course at the price of simplifications, global views of cerebral functioning and more local hypotheses on the role of subsets of neuronal structures in cognition. In the global approaches, the integration of data from experimental psychology and clinical studies leads to an overview of the brain as a set of interacting memories, each devoted to a specific kind of information processing [53]. It results also in longstanding and very ambitious studies for the design of cognitive architectures aiming at embracing the whole cognition. With the notable exception of works initiated by [50], most of these frameworks (e.g. Soar, ACT-R), though sometimes justified on biological grounds, do not go up to a *connectionist* neuronal implementation. Furthermore, because of the complexity of the resulting frameworks, they are restricted to simple symbolic interfaces with the internal and external world and to (relatively) small-sized internal structures. Our main research objective is undoubtly to build such a general purpose cognitive architecture (to model the brain *as a whole* in a systemic way), using a connectionist implementation and able to cope with a realistic environment.

## 3.2. Computational Neuroscience

From a general point of view, computational neuroscience can be defined as the development of methods from computer science and applied mathematics, to explore more technically and theoretically the relations between structures and functions in the brain [55], [44]. During the recent years this domain has gained an increasing interest in neuroscience and has become an essential tool for scientific developments in most fields in neuroscience, from the molecule to the system. In this view, all the objectives of our team can be described as possible progresses in computational neuroscience. Accordingly, it can be underlined that the systemic view that we promote can offer original contributions in the sense that, whereas most classical models in computational neuroscience focus on the better understanding of the structure/function relationship for isolated specific structures, we aim at exploring synergies between structures. Consequently, we target interfaces and interplay between heterogenous modes of computing, which is rarely addressed in classical computational neuroscience.

We also insist on another aspect of computational neuroscience which is, in our opinion, at the core of the involvement of computer scientists and mathematicians in the domain and on which we think we could particularly contribute. Indeed, we think that our primary abilities in numerical sciences imply that our developments are characterized above all by the effectiveness of the corresponding computations: We provide biologically inspired architectures with effective computational properties, such as robustness to noise, self-organization, on-line learning. We more generally underline the requirement that our models must also mimick biology through its most general law of homeostasis and self-adaptability in an unknown and changing environment. This means that we propose to numerically experiment such models and thus provide effective methods to falsify them.

Here, computational neuroscience means mimicking original computations made by the neuronal substratum and mastering their corresponding properties: computations are distributed and adaptive; they are performed without an homonculus or any central clock. Numerical schemes developed for distributed dynamical systems and algorithms elaborated for distributed computations are of central interest here [41], [49] and were the basis for several contributions in our group [54], [51], [56]. Ensuring such a rigor in the computations associated to our systemic and large scale approach is of central importance.

Equally important is the choice for the formalism of computation, extensively discussed in the connectionist domain. Spiking neurons are today widely recognized of central interest to study synchronization mechanisms and neuronal coupling at the microscopic level [42]; the associated formalism [47] can be possibly considered for local studies or for relating our results with this important domain in connectionism. Nevertheless, we remain mainly at the mesoscopic level of modeling, the level of the neuronal population, and consequently interested in the formalism developed for dynamic neural fields [39], that demonstrated a richness of behavior [43] adapted to the kind of phenomena we wish to manipulate at this level of description. Our group has a long experience in the study and adaptation of the properties of neural fields [51], [52] and their use for observing the emergence of typical cortical properties [46]. In the envisioned development of more complex architectures and interplay between structures, the exploration of mathematical properties such as stability and boundedness and the observation of emerging phenomena is one important objective. This objective is also associated with that of capitalizing our experience and promoting good practices in our software production. In summary, we think that this systemic approach also brings to computational neuroscience new case studies where heterogenous and adaptive models with various time scales and parameters have to be considered jointly to obtain a mastered substratum of computation. This is particularly critical for large scale deployments.

## 3.3. Machine Learning

The adaptive properties of the nervous system are certainly among its most fascinating characteristics, with a high impact on our cognitive functions. Accordingly, machine learning is a domain [48] that aims at giving such characteristics to artificial systems, using a mathematical framework (probabilities, statistics, data analysis, etc.). Some of its most famous algorithms are directly inspired from neuroscience, at different levels. Connectionist learning algorithms implement, in various neuronal architectures, weight update rules, generally derived from the hebbian rule, performing non supervised (e.g. Kohonen self-organizing maps), supervised (e.g. layered perceptrons) or associative (e.g. Hopfield recurrent network) learning. Other algorithms, not necessarily connectionist, perform other kinds of learning, like reinforcement learning. Machine learning is a very mature domain today and all these algorithms have been extensively studied, at both the theoretical and practical levels, with much success. They have also been related to many functions (in the living and artificial domains) like discrimination, categorisation, sensorimotor coordination, planning, etc. and several neuronal structures have been proposed as the substratum for these kinds of learning [45], [38]. Nevertheless, we believe that, as for previous models, machine learning algorithms remain isolated tools, whereas our systemic approach can bring original views on these problems.

At the cognitive level, most of the problems we face do not rely on only one kind of learning and require instead skills that have to be learned in preliminary steps. That is the reason why cognitive architectures are often referred to as systems of memory, communicating and sharing information for problem solving. Instead of the classical view in machine learning of a flat architecture, a more complex network of modules must be

considered here, as it is the case in the domain of deep learning. In addition, our systemic approach brings the question of incrementally building such a system, with a clear inspiration from developmental sciences. In this perspective, modules can generate internal signals corresponding to internal goals, predictions, error signals, able to supervise the learning of other modules (possibly endowed with a different learning rule), supposed to become autonomous after an instructing period. A typical example is that of episodic learning (in the hippocampus), storing declarative memory about a collection of past episods and supervising the training of a procedural memory in the cortex.

At the behavioral level, as mentionned above, our systemic approach underlines the fundamental links between the adaptive system and the internal and external world. The internal world includes proprioception and interoception, giving information about the body and its needs for integrity and other fundamental programs. The external world includes physical laws that have to be learned and possibly intelligent agents for more complex interactions. Both involve sensors and actuators that are the interfaces with these worlds and close the loops. Within this rich picture, machine learning generally selects one situation that defines useful sensors and actuators and a corpus with properly segmented data and time, and builds a specific architecture and its corresponding criteria to be satisfied. In our approach however, the first question to be raised is to discover what is the goal, where attention must be focused on and which previous skills must be exploited, with the help of a dynamic architecture and possibly other partners. In this domain, the behavioral and the developmental sciences, observing how and along which stages an agent learns, are of great help to bring some structure to this high dimensional problem.

At the implementation level, this analysis opens many fundamental challenges, hardly considered in machine learning : stability must be preserved despite on-line continuous learning; criteria to be satisfied often refer to behavioral and global measurements but they must be translated to control the local circuit level; in an incremental or developmental approach, how will the development of new functions preserve the integrity and stability of others? In addition, this continous re-arrangement is supposed to involve several kinds of learning, at different time scales (from msec to years in humans) and to interfer with other phenomena like variability and meta-plasticity.

In summary, our main objective in machine learning is to propose on-line learning systems, where several modes of learning have to collaborate and where the protocoles of training are realistic. We promote here a *really autonomous* learning, where the agent must select by itself internal resources (and build them if not available) to evolve at the best in an unknown world, without the help of any *deus-ex-machina* to define parameters, build corpus and define training sessions, as it is generally the case in machine learning. To that end, autonomous robotics (*cf.* § 3.4 ) is a perfect testbed.

## 3.4. Autonomous Robotics

Autonomous robots are not only convenient platforms to implement our algorithms; the choice of such platforms is also motivated by theories in cognitive science and neuroscience indicating that cognition emerges from interactions of the body in direct loops with the world (*embodiment of cognition* [40]). In addition to real robotic platforms, software implementations of autonomous robotic systems including components dedicated to their body and their environment will be also possibly exploited, considering that they are also a tool for studying conditions for a real autonomous learning.

A real autonomy can be obtained only if the robot is able to define its goal by itself, without the specification of any high level and abstract cost function or rewarding state. To ensure such a capability, we propose to endow the robot with an artificial physiology, corresponding to perceive some kind of pain and pleasure. It may consequently discriminate internal and external goals (or situations to be avoided). This will mimick circuits related to fundamental needs (e.g. hunger and thirst) and to the preservation of bodily integrity. An important objective is to show that more abstract planning capabilities can arise from these basic goals.

A real autonomy with an on-line continuous learning as described in  § 3.3  will be made possible by the elaboration of protocols of learning, as it is the case, in animal conditioning, for experimental studies where performance on a task can be obtained only after a shaping in increasingly complex tasks. Similarly,

developmental sciences can teach us about the ordered elaboration of skills and their association in more complex schemes. An important challenge here is to translate these hints at the level of the cerebral architecture.

As a whole, autonomous robotics permits to assess the consistency of our models in realistic condition of use and offers to our colleagues in behavioral sciences an object of study and comparison, regarding behavioral dynamics emerging from interactions with the environment, also observable at the neuronal level.

In summary, our main contribution in autonomous robotics is to make autonomy possible, by various means corresponding to endow robots with an artificial physiology, to give instructions in a natural and incremental way and to prioritize the synergy between reactive and robust schemes over complex planning structures.

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

# 3. Research Program

## 3.1. Introduction

We are working in the context of data-driven medicine against cancer. We aim at coupling mathematical models with data to address relevant challenges for biologists and clinicians in order for instance to improve our understanding in cancer biology and pharmacology, assist the development of novel therapeutic approaches or develop personalized decision-helping tools for monitoring the disease and evaluating therapies.

More precisely, our research on mathematical oncology is three-fold:

- Axis 1: Tumor modeling for patient-specific simulations: *Clinical monitoring. Numerical markers from imaging data. Radiomics.*

- Axis 2: Bio-physical modeling for personalized therapies: *Electroporation from cells to tissue. Radiotherapy.*

- Axis 3: Quantitative cancer modeling for biological, clinical and preclinical studies: *Biological mechanisms. Metastatic dissemination. Pharmacometrics.*

In the first axis, we aim at producing patient-specific simulations of the growth of a tumor or its response to treatment starting from a series of images. We hope to be able to offer a valuable insight on the disease to the clinicians in order to improve the decision process. This would be particularly useful in the cases of relapses or for metastatic diseases.

The second axis aims at modeling biophysical therapies like electroporation, but also radiotherapy, thermo-ablations, radio-frequency ablations or electroporation that play a crucial role for a local treatment of the disease if possible limiting the metastatic dissemination, which is precisely the clinical context where the techniques of axis 1 will be applied.

The third axis is essential since it is a way to better understand and model the biological reality of cancer growth and the (possibly complex) effects of therapeutic intervention. Modeling in this case also helps to interpret the experimental results and improve the accuracy of the models used in Axis 1. Technically speaking, some of the computing tools are similar to those of Axis 1.

## 3.2. Axis 1: Tumor modeling for patient-specific simulations

The gold standard treatment for most cancers is surgery. In the case where total resection of the tumor is possible, the patient often benefits from an adjuvant therapy (radiotherapy, chemotherapy, targeted therapy or a combination of them) in order to eliminate the potentially remaining cells that may not be visible. In this case personalized modeling of tumor growth is useless and statistical modeling will be able to quantify the risk of relapse, the mean progression-free survival time...However if total resection is not possible or if metastases emerge from distant sites, clinicians will try to control the disease for as long as possible. A wide set of tools are available. Clinicians may treat the disease by physical interventions (radiofrequency ablation, cryoablation, radiotherapy, electroporation, focalized ultrasound,...) or chemical agents (chemotherapies, targeted therapies, antiangiogenic drugs, immunotherapies, hormonotherapies). One can also decide to monitor the patient without any treatment (this is the case for slowly growing tumors like some metastases to the lung, some lymphomas or for some low grade glioma). A reliable patient-specific model of tumor evolution with or without therapy may have different uses:

- Case without treatment: the evaluation of the growth of the tumor would offer a useful indication for the time at which the tumor may reach a critical size. For example, radiofrequency ablation of pulmonary lesion is very efficient as long as the diameter of the lesion is smaller than 3 cm. Thus, the prediction can help the clinician plan the intervention. For slowly growing tumors, quantitative

modeling can also help to decide at what time interval the patient has to undergo a CT-scan. CT-scans are irradiative exams and there is a challenge for decreasing their occurrence for each patient. It has also an economical impact. And if the disease evolution starts to differ from the prediction, this might mean that some events have occurred at the biological level. For instance, it could be the rise of an aggressive phenotype or cells that leave a dormancy state. This kind of events cannot be predicted, but some mismatch with respect to the prediction can be an indirect proof of their existence. It could be an indication for the clinician to start a treatment.

- Case with treatment: a model can help to understand and to quantify the final outcome of a treatment using the early response. It can help for a redefinition of the treatment planning. Modeling can also help to anticipate the relapse by analyzing some functional aspects of the tumor. Again, a deviation with respect to reference curves can mean a lack of efficiency of the therapy or a relapse. Moreover, for a long time, the response to a treatment has been quantified by the RECIST criteria which consists in (roughly speaking) measuring the diameters of the largest tumor of the patient, as it is seen on a CT-scan. This criteria is still widely used and was quite efficient for chemotherapies and radiotherapies that induce a decrease of the size of the lesion. However, with the systematic use of targeted therapies and anti-angiogenic drugs that modify the physiology of the tumor, the size may remain unchanged even if the drug is efficient and deeply modifies the tumor behavior. One better way to estimate this effect could be to use functional imaging (Pet-scan, perfusion or diffusion MRI, ...), a model can then be used to exploit the data and to understand in what extent the therapy is efficient.

- Optimization: currently, we do not believe that we can optimize a particular treatment in terms of distribution of doses, number, planning with the model that we will develop in a medium term perspective.

The scientific challenge is therefore as follows: given the history of the patient, the nature of the primitive tumor, its histopathology, knowing the treatments that patients have undergone, some biological facts on the tumor and having a sequence of images (CT-scan, MRI, PET or a mix of them), are we able to provide a numerical simulation of the extension of the tumor and of its metabolism that fits as best as possible with the data (CT-scans or functional data) and that is predictive in order to address the clinical cases described above?

Our approach relies on the elaboration of PDE models and their parametrization with images by coupling deterministic and stochastic methods. The PDE models rely on the description of the dynamics of cell populations. The number of populations depends on the pathology. For example, for glioblastoma, one needs to use proliferative cells, invasive cells, quiescent cells as well as necrotic tissues to be able to reproduce realistic behaviors of the disease. In order to describe the relapse for hepatic metastases of gastro-intestinal stromal tumor (gist), one needs three cell populations: proliferative cells, healthy tissue and necrotic tissue.

The law of proliferation is often coupled with a model for the angiogenesis. However such models of angiogenesis involve too many non measurable parameters to be used with real clinical data and therefore one has to use simplified or even simplistic versions. The law of proliferation often mimics the existence of an hypoxia threshold, it consists of an ODE. or a PDE that describes the evolution of the growth rate as a combination of sigmoid functions of nutrients or roughly speaking oxygen concentration. Usually, several laws are available for a given pathology since at this level, there are no quantitative argument to choose a particular one.

The velocity of the tumor growth differs depending on the nature of the tumor. For metastases, we will derive the velocity thanks to Darcy's law in order to express that the extension of the tumor is basically due to the increase of volume. This gives a sharp interface between the metastasis and the surrounding healthy tissues, as observed by anatomopathologists. For primitive tumors like glioma or lung cancer, we use reaction-diffusion equations in order to describe the invasive aspects of such primitive tumors.

The modeling of the drugs depends on the nature of the drug: for chemotherapies, a death term can be added into the equations of the population of cells, while antiangiogenic drugs have to be introduced in a angiogenic model. Resistance to treatment can be described either by several populations of cells or with non-constant growth or death rates. As said before, it is still currently difficult to model the changes of phenotype

or mutations, we therefore propose to investigate this kind of phenomena by looking at deviations of the numerical simulations compared to the medical observations.

The calibration of the model is achieved by using a series (at least 2) of images of the same patient and by minimizing a cost function. The cost function contains at least the difference between the volume of the tumor that is measured on the images with the computed one. It also contains elements on the geometry, on the necrosis and any information that can be obtained through the medical images. We will pay special attention to functional imaging (PET, perfusion and diffusion MRI). The inverse problem is solved using a gradient method coupled with some Monte-Carlo type algorithm. If a large number of similar cases is available, one can imagine to use statistical algorithms like random forests to use some non quantitative data like the gender, the age, the origin of the primitive tumor...for example for choosing the model for the growth rate for a patient using this population knowledge (and then to fully adapt the model to the patient by calibrating this particular model on patient data) or for having a better initial estimation of the modeling parameters. We have obtained several preliminary results concerning lung metastases including treatments and for metastases to the liver.
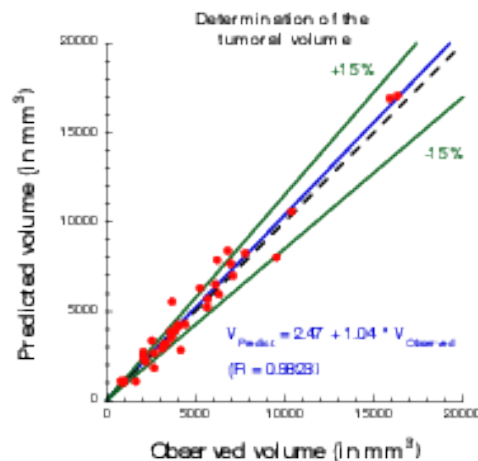


*Figure 4. Plot showing the accuracy of our prediction on meningioma volume. Each point corresponds to a patient whose two first exams were used to calibrate our model. A patient-specific prediction was made with this calibrated model and compared with the actual volume as measured on a third time by clinicians. A perfect prediction would be on the black dashed line. Medical data was obtained from Prof. Loiseau, CHU Pellegrin.*

## 3.3. Axis 2: Bio-physical modeling for personalized therapies

In this axis, we investigate locoregional therapies such as radiotherapy, irreversible electroporation. Electroporation consists in increasing the membrane permeability of cells by the delivery of high voltage pulses. This non-thermal phenomenon can be transient (reversible) or irreversible (IRE). IRE or electro-chemotherapy – which is a combination of reversible electroporation with a cytotoxic drug – are essential tools for the treatment of a metastatic disease. Numerical modeling of these therapies is a clear scientific challenge. Clinical applications of the modeling are the main target, which thus drives the scientific approach, even though theoretical studies in order to improve the knowledge of the biological phenomena, in particular for electroporation, should also be addressed. However, this subject is quite wide and we focus on two particular approaches: some aspects of radiotherapies and electro-chemotherapy. This choice is motivated partly by pragmatic reasons: we

already have collaborations with physicians on these therapies. Other treatments could be probably treated with the same approach, but we do not plan to work on this subject on a medium term.

- Radiotherapy (RT) is a common therapy for cancer. Typically, using a CT scan of the patient with the structures of interest (tumor, organs at risk) delineated, the clinicians optimize the dose delivery to treat the tumor while preserving healthy tissues. The RT is then delivered every day using low resolution scans (CBCT) to position the beams. Under treatment the patient may lose weight and the tumor shrinks. These changes may affect the propagation of the beams and subsequently change the dose that is effectively delivered. It could be harmful for the patient especially if sensitive organs are concerned. In such cases, a replanification of the RT could be done to adjust the therapeutical protocol. Unfortunately, this process takes too much time to be performed routinely. The challenges faced by clinicians are numerous, we focus on two of them:

  – *Detecting the need of replanification:* we are using the positioning scans to evaluate the movement and deformation of the various structures of interest. Thus we can detect whether or not a structure has moved out of the safe margins (fixed by clinicians) and thus if a replanification may be necessary. In a retrospective study, our work can also be used to determine RT margins when there are no standard ones. A collaboration with the RT department of Institut Bergonié is underway on the treatment of retroperitoneal sarcoma and ENT tumors (head and neck cancers). A retrospective study was performed on 11 patients with retro-peritoneal sarcoma. The results have shown that the safety margins (on the RT) that clinicians are currently using are probably not large enough. The tool used in this study was developed by an engineer funded by Inria (Cynthia Périer, ADT Sesar). We used well validated methods from a level-set approach and segmentation / registration methods. The originality and difficulty lie in the fact that we are dealing with real data in a clinical setup. Clinicians have currently no way to perform complex measurements with their clinical tools. This prevents them from investigating the replanification. Our work and the tools developed pave the way for easier studies on evaluation of RT plans in collaboration with Institut Bergonié. *There was no modeling involved in this work that arose during discussions with our collaborators.* The main purpose of the team is to have meaningful outcomes of our research for clinicians, sometimes it implies leaving a bit our area of expertise.

  – *Evaluating RT efficacy and finding correlation between the radiological responses and the clinical outcome:* our goal is to help doctors to identify correlation between the response to RT (as seen on images) and the longer term clinical outcome of the patient. Typically, we aim at helping them to decide when to plan the next exam after the RT. For patients whose response has been linked to worse prognosis, this exam would have to be planned earlier. This is the subject of collaborations with Institut Bergonié and CHU Bordeaux on different cancers (head and neck, pancreas). The response is evaluated from image markers (*e.g.* using texture information) or with a mathematical model developed in Axis 1. The other challenges are either out of reach or not in the domain of expertise of the team. Yet our works may tackle some important issues for adaptive radiotherapy.

- Both IRE and electrochemotherapy are anticancerous treatments based on the same phenomenon: the electroporation of cell membranes. This phenomenon is known for a few decades but it is still not well understood, therefore our interest is two fold:

  1. We want to use mathematical models in order to better understand the biological behavior and the effect of the treatment. We work in tight collaboration with biologists and bioeletromagneticians to derive precise models of cell and tissue electroporation, in the continuity of the research program of the Inria team-project MC2. These studies lead to complex non-linear mathematical models involving some parameters (as less as possible). Numerical methods to compute precisely such models and the calibration of the parameters with the experimental data are then addressed. Tight collaborations with the Vectorology and Anticancerous Therapies (VAT) of IGR at Villejuif, Laboratoire Ampère of Ecole

Centrale Lyon and the Karlsruhe Institute of technology will continue, and we aim at developing new collaborations with Institute of Pharmacology and Structural Biology (IPBS) of Toulouse and the Laboratory of Molecular Pathology and Experimental Oncology (LM-PEO) at CNR Rome, in order to understand differences of the electroporation of healthy cells and cancer cells in spheroids and tissues.

2. This basic research aims at providing new understanding of electroporation, however it is necessary to address, particular questions raised by radio-oncologists that apply such treatments. One crucial question is "What pulse or what train of pulses should I apply to electroporate the tumor if the electrodes are located as given by the medical images"? Even if the real-time optimization of the placement of the electrodes for deep tumors may seem quite utopian since the clinicians face too many medical constraints that cannot be taken into account (like the position of some organs, arteries, nerves...), one can expect to produce real-time information of the validity of the placement done by the clinician. Indeed, once the placement is performed by the radiologists, medical images are usually used to visualize the localization of the electrodes. Using these medical data, a crucial goal is to provide a tool in order to compute in real-time and visualize the electric field and the electroporated region directly on theses medical images, to give the doctors a precise knowledge of the region affected by the electric field. In the long run, this research will benefit from the knowledge of the theoretical electroporation modeling, but it seems important to use the current knowledge of tissue electroporation – even quite rough –, in order to rapidly address the specific difficulty of such a goal (real-time computing of non-linear model, image segmentation and visualization). Tight collaborations with CHU Pellegrin at Bordeaux, and CHU J. Verdier at Bondy are crucial.

- <u>Radiofrequency ablation.</u> In a collaboration with Hopital Haut Leveque, CHU Bordeaux we are trying to determine the efficacy and risk of relapse of hepatocellular carcinoma treated by radiofrequency ablation. For this matter we are using geometrical measurements on images (margins of the RFA, distance to the boundary of the organ) as well as texture information to statistically evaluate the clinical outcome of patients.

- <u>Intensity focused ultrasound.</u> In collaboration with Utrecht Medical center, we aim at tackling several challenges in clinical applications of IFU: target tracking, dose delivery...

## 3.4. Axis 3: Quantitative cancer modeling for biological and preclinical studies

With the emergence and improvement of a plethora of experimental techniques, the molecular, cellular and tissue biology has operated a shift toward a more quantitative science, in particular in the domain of cancer biology. These quantitative assays generate a large amount of data that call for theoretical formalism in order to better understand and predict the complex phenomena involved. Indeed, due to the huge complexity underlying the development of a cancer disease that involves multiple scales (from the genetic, intra-cellular scale to the scale of the whole organism), and a large number of interacting physiological processes (see the so-called "hallmarks of cancer"), several questions are not fully understood. Among these, we want to focus on the most clinically relevant ones, such as the general laws governing tumor growth and the development of metastases (secondary tumors, responsible of 90% of the deaths from a solid cancer). In this context, it is thus challenging to exploit the diversity of the data available in experimental settings (such as *in vitro* tumor spheroids or *in vivo* mice experiments) in order to improve our understanding of the disease and its dynamics, which in turn lead to validation, refinement and better tuning of the macroscopic models used in the axes 1 and 2 for clinical applications.

In recent years, several new findings challenged the classical vision of the metastatic development biology, in particular by the discovery of organism-scale phenomena that are amenable to a dynamical description in terms of mathematical models based on differential equations. These include the angiogenesis-mediated distant inhibition of secondary tumors by a primary tumor the pre-metastatic niche or the self-seeding phenomenon Building a general, cancer type specific, comprehensive theory that would integrate these dynamical processes

remains an open challenge. On the therapeutic side, recent studies demonstrated that some drugs (such as the Sunitinib), while having a positive effect on the primary tumor (reduction of the growth), could *accelerate* the growth of the metastases. Moreover, this effect was found to be scheduling-dependent. Designing better ways to use this drug in order to control these phenomena is another challenge. In the context of combination therapies, the question of the *sequence* of administration between the two drugs is also particularly relevant.

One of the technical challenge that we need to overcome when dealing with biological data is the presence of potentially very large inter-animal (or inter-individual) variability.

Starting from the available multi-modal data and relevant biological or therapeutic questions, our purpose is to develop adapted mathematical models (*i.e.* identifiable from the data) that recapitulate the existing knowledge and reduce it to its more fundamental components, with two main purposes:

1. to generate quantitative and empirically testable predictions that allow to assess biological hypotheses or

2. to investigate the therapeutic management of the disease and assist preclinical studies of anti-cancerous drug development.

We believe that the feedback loop between theoretical modeling and experimental studies can help to generate new knowledge and improve our predictive abilities for clinical diagnosis, prognosis, and therapeutic decision. Let us note that the first point is in direct link with the axes 1 and 2 of the team since it allows us to experimentally validate the models at the biological scale (*in vitro* and *in vivo* experiments) for further clinical applications.

More precisely, we first base ourselves on a thorough exploration of the biological literature of the biological phenomena we want to model: growth of tumor spheroids, *in vivo* tumor growth in mice, initiation and development of the metastases, effect of anti-cancerous drugs. Then we investigate, using basic statistical tools, the data we dispose, which can range from: spatial distribution of heterogeneous cell population within tumor spheroids, expression of cell markers (such as green fluorescent protein for cancer cells or specific antibodies for other cell types), bioluminescence, direct volume measurement or even intra-vital images obtained with specific imaging devices. According to the data type, we further build dedicated mathematical models that are based either on PDEs (when spatial data is available, or when time evolution of a structured density can be inferred from the data, for instance for a population of tumors) or ODEs (for scalar longitudinal data). These models are confronted to the data by two principal means:

1. when possible, experimental assays can give a direct measurement of some parameters (such as the proliferation rate or the migration speed) or

2. statistical tools to infer the parameters from observables of the model.

This last point is of particular relevance to tackle the problem of the large inter-animal variability and we use adapted statistical tools such as the mixed-effects modeling framework.

Once the models are shown able to describe the data and are properly calibrated, we use them to test or simulate biological hypotheses. Based on our simulations, we then aim at proposing to our biological collaborators new experiments to confirm or infirm newly generated hypotheses, or to test different administration protocols of the drugs. For instance, in a collaboration with the team of the professor Andreas Bikfalvi (Laboratoire de l'Angiogénèse et du Micro-environnement des Cancers, Inserm, Bordeaux), based on confrontation of a mathematical model to multi-modal biological data (total number of cells in the primary and distant sites and MRI), we could demonstrate that the classical view of metastatic dissemination and development (one metastasis is born from one cell) was probably inaccurate, in mice grafted with metastatic kidney tumors. We then proposed that metastatic germs could merge or attract circulating cells. Experiments involving cells tagged with two different colors are currently performed in order to confirm or infirm this hypothesis.

Eventually, we use the large amount of temporal data generated in preclinical experiments for the effect of anti-cancerous drugs in order to design and validate mathematical formalisms translating the biological mechanisms of action of these drugs for application to clinical cases, in direct connection with the axis 1. We have a special focus on targeted therapies (designed to specifically attack the cancer cells while sparing the

healthy tissue) such as the Sunitinib. This drug is indeed indicated as a first line treatment for metastatic renal cancer and we plan to conduct a translational study coupled between A. Bikfalvi's laboratory and medical doctors, F. Cornelis (radiologist) and A. Ravaud (head of the medical oncology department).

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

# 3. Research Program

## 3.1. A Geometric View of Diversity

Diversity may be studied as a set of dissimilarities between objects. The underlying mathematical construction is the notion of distance. Knowing a set of objects, it is possible, after computation of pairwise distances, or sometimes dissimilarities, to build a Euclidean image of it as a point cloud in a space of relevant dimension. Then, diversity can be associated with the shape of the point cloud. The human eye is often far better than an algorithm at recognizing a pattern or shape. One objective of our project is to narrow the gap between the story that a human eye can tell, and that an algorithm can tell. Several directions will be explored. First, this requires mastering classical tools in dimension reduction, mainly algebraic tools (PCA, NGS, Isomap, eigenmaps, etc ...). Second, neighborhoods in point clouds naturally lead to graphs describing the neighborhood networks. There is a natural link between modular structures in distance arrays and communities on graphs. Third, points (representing, say, DNA sequences) are samples of diversity. Dimension reduction may show that they live on a given manifold. This leads to geometry (differential or Riemannian geometry). It is expected that some properties of the manifold can tell something of the constraints on the space where measured individuals live. The connection between Riemannian geometry and graphs, where weighted graphs are seen as mesh embedded in a manifold, is currently an active field of research [28], [27]. See as well [30] for a link between geometric structure, linear and nonlinear dimensionality reduction.

Biodiversity and high-performance computing: Most methods and tools for characterizing diversity have been designed for datasets that can be analyzed on a laptop, but NGS datasets produced for metabarcoding are far too large. Data analysis algorithms and tools must be revisited and scaled up. We will mobilize both distributed algorithms like the Arnoldi method and new algorithms, like random projection or column selection methods, to build point clouds in Euclidean spaces from massive data sets, and thus to overcome the cubic complexity of computation of eigenvectors and eigenvalues of very large dense matrices. We will also link distance geometry [22] with convex optimization procedures through matrix completion [15], [17].

Intercalibration: There is a considerable difference between supervised and unsupervised clustering: in supervised clustering, the result for an item $i$ is independent from the result for an item $j \neq i$, whereas in unsupervised clustering, the result for an item $i$ (e.g. the cluster it belongs to, and its composition) depends on nearby items $j \neq i$. Which means that the result may change if some items are added to or subtracted from the sample. This raises the more global problem of how to merge two studies to yield a more comprehensive view of biodiversity?

## 3.2. Knowledge Management for Biology

The heterogenous data generated in computational molecular biology and ecology are distinguished not only by their volume, but by the richness of the many levels of interpretation that biologists create. The same nucleic acid sequence can be seen as a molecule with a structure, a sequence of base pairs, a collection of genes, an allele, or a molecular fingerprint. To extract the maximum benefit from this treasure trove we must organize the knowledge in ways that facilitate extraction, analysis, and inference. Our focus has been on the efficient representation of relations between biological objects and operations on those representations, in particular heuristic analyses and logical inference.

PLEIADE will develop applications in comparative genomics of related organisms, using new mathematical tools for representing compactly, at different scales of difference, comparisons between related genomes. New methods based on distance geometry will refine these comparisons. Compact representations can be stored, exchanged, and combined. They will form the basis of new simultaneous genome annotation methods, linked directly to abductive inference methods for building functional models of the organisms and their communities.
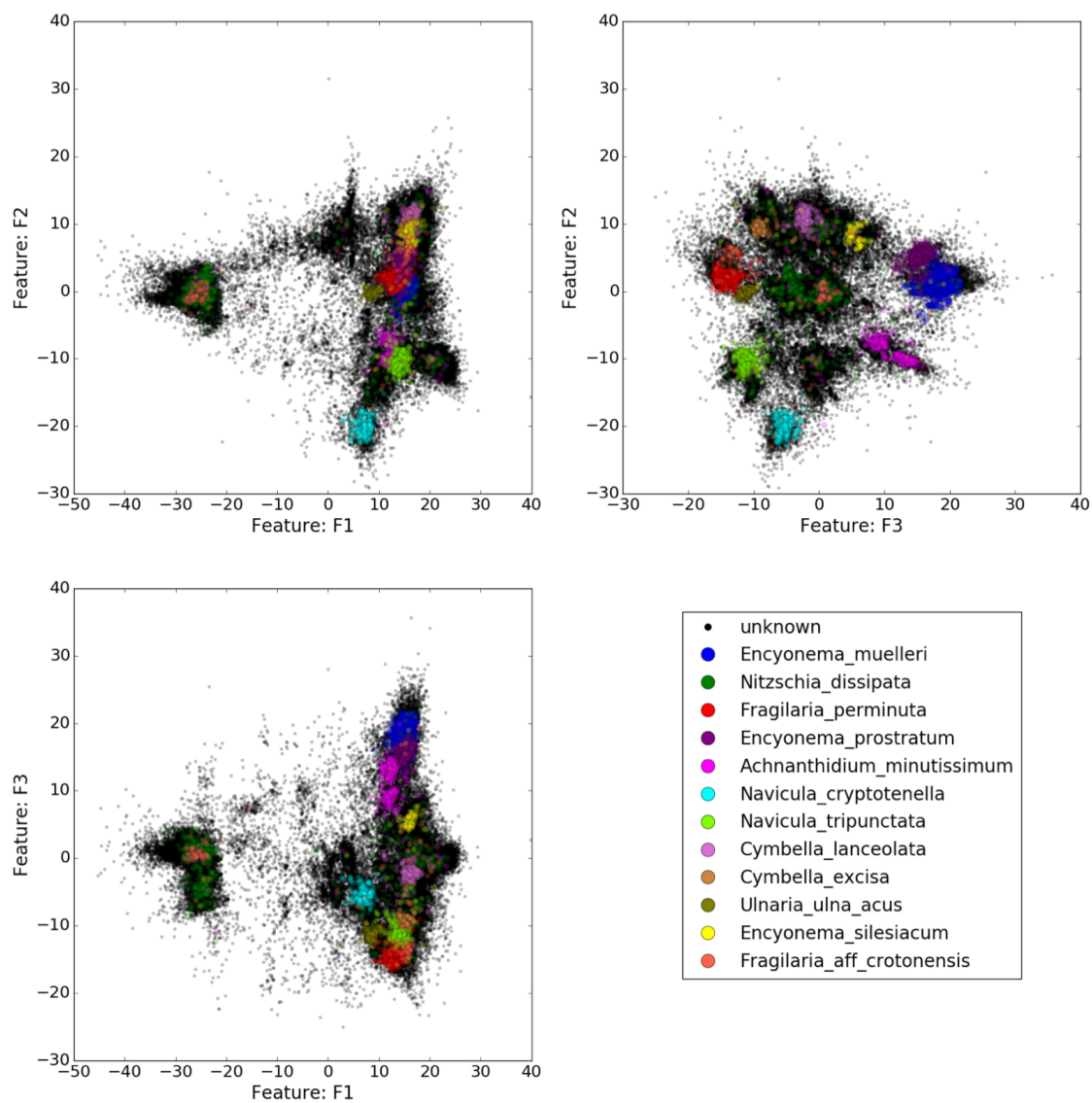
*Figure 3. Validation of high density islands using supervised classification. Metagenomic reads from diatoms in Lake Geneva [26] were analyzed by the method from [16] and colored by species according to a reference database.*

Since a goal of PLEIADE is to integrate diversity throughout the analysis process, it is necessary to incorporate **diversity as a form of knowledge** that can be stored in a knowledge base. Diversity can be represented using various compact representations, such as trees and quotient graphs storing nested sets of relations. Extracting structured representations and logical relations from integrated knowledge bases (Figure 2 ) will require domain-specific query methods that can express forms of diversity.

## 3.3. Modeling by successive refinement

Describing the links between diversity in traits and diversity in function will require comprehensive models, assembled from and refining existing models. A recurring difficulty in building comprehensive models of biological systems is that accurate models for subsystems are built using different formalisms and simulation techniques, and hand-tuned models tend to be so focused in scope that it is difficult to repurpose them [13]. Our belief is that a sustainable effort in building efficient behavioral models must proceed incrementally, rather than by modeling individual processes *de novo*. *Hierarchical modeling* [10] is one way of combining specific models into networks. Effective use of hierarchical models requires both formal definition of the semantics of such composition, and efficient simulation tools for exploring the large space of complex behaviors. We have previously shown that this approach can be effective for certains kinds of systems in biotechnology [2], [14] and medicine [12]. Our challenge is to adapt incremental, hierarchical refinement to modeling organisms and communities in metagenomic and comparative genomic applications.

<p align="center" style="color:red"><strong>SISTM Project-Team</strong></p>

# 3. Research Program

## 3.1. Mechanistic learning

When studying the dynamics of a given marker, say the HIV concentration in the blood (HIV viral load), one can for instance use descriptive models summarizing the dynamics over time in term of slopes of the trajectories [56]. These slopes can be compared between treatment groups or according to patients' characteristics. Another way for analyzing these data is to define a mathematical model based on the biological knowledge of what drives HIV dynamics. In this case, it is mainly the availability of target cells (the CD4+ T lymphocytes), the production and death rates of infected cells and the clearance of the viral particles that impact the dynamics. Then, a mathematical model most often based on ordinary differential equations (ODE) can be written [48]. Estimating the parameters of this model to fit observed HIV viral load gave a crucial insight in HIV pathogenesis as it revealed the very short half-life of the virions and infected cells and therefore a very high turnover of the virus, making mutations a very frequent event [47].

Having a good mechanistic model in a biomedical context such as HIV infection opens doors to various applications beyond a good understanding of the data. Global and individual predictions can be excellent because of the external validity of a model based on main biological mechanisms. Control theory may serve for defining optimal interventions or optimal designs to evaluate new interventions [40]. Finally, these models can capture explicitly the complex relationship between several processes that change over time and may therefore challenge other proposed approaches such as marginal structural models to deal with causal associations in epidemiology [39].

Therefore, we postulate that this type of model could be very useful in the context of our research that is in complex biological systems. The definition of the model needs to identify the parameter values that fit the data. In clinical research this is challenging because data are sparse, and often unbalanced, coming from populations of subjects. A substantial inter-individual variability is always present and needs to be accounted as this is the main source of information. Although many approaches have been developed to estimate the parameters of non-linear mixed models [51], [59], [43], [49], [44], [58], the difficulty associated with the complexity of ODE models and the sparsity of the data leading to identifiability issues need further research.

Furthermore, the availability of data for each individual (see below) leads to a new challenge in this area. The structural model can easily be much more complex and the observation model may need to integrate much more markers.

## 3.2. High-dimensional statistical learning

With the availability of omics data such as genomics (DNA), transcriptomics (RNA) or proteomics (proteins), but also other types of data, such as those arising from the combination of large observational databases (e.g. in pharmacoepidemiology or environmental epidemiology), high-dimensional data have became increasingly common. Use of molecular biological technics such as Polymerase Chain Reaction (PCR) allows for amplification of DNA or RNA sequences. Nowadays, microarray and Next Generation Sequencing (NGS) techniques give the possibility to explore very large portions of the genome. Furthermore, other assays have also evolved, and traditional measures such as cytometry or imaging have became new sources of big data. Therefore, in the context of HIV research, the dimension of the datasets has much grown in term of number of variables per individual than in term of number of included patients although this latter is also growing thanks to the multi-cohort collaborations such as CASCADE or COHERE organized in the EuroCoord network [0]. As an example, in a phase 1/2 clinical trial evaluating the safety and the immunological response to a dendritic cell-based HIV vaccine, 19 infected patients were included. Bringing together data on cell count, cytokine production,

---

[0] see online at http://www.eurocoord.net

gene expression and viral genome change led to a 20 Go database [55]. This is far from big databases faced in other areas but constitutes a revolution in clinical research where clinical trials of hundred of patients sized few hundred of Ko at most. Therefore, more than the storage and calculation capacities, the challenge is the comprehensive analysis of these data-sets.

The objective is either to select the relevant information or to summarize it for understanding or prediction purposes. When dealing with high-dimensional data, the methodological challenge arises from the fact that data-sets typically contain many variables, much more than observations. Hence, multiple testing is an obvious issue that needs to be taken into account [52]. Furthermore, conventional methods, such as linear models, are inefficient and most of the time even inapplicable. Specific methods have been developed, often derived from the machine learning field, such as regularization methods [57]. The integrative analysis of large data-sets is challenging. For instance, one may want to look at the correlation between two large scale matrices composed by the transcriptome in the one hand and the proteome on the other hand [45]. The comprehensive analysis of these large data-sets concerning several levels from molecular pathways to clinical response of a population of patients needs specific approaches and a very close collaboration with the providers of data that is the immunologists, the virologists, the clinicians...

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

# 3. Research Program

## 3.1. Introduction

The methodological component of <span style="color:red">HIEPACS</span>  concerns the expertise for the design as well as the efficient and scalable implementation of highly parallel numerical algorithms to perform frontier simulations. In order to address these computational challenges a hierarchical organization of the research is considered. In this bottom-up approach, we first consider in Section 3.2 generic topics concerning high performance computational science. The activities described in this section are transversal to the overall project and their outcome will support all the other research activities at various levels in order to ensure the parallel scalability of the algorithms. The aim of this activity is not to study general purpose solution but rather to address these problems in close relation with specialists of the field in order to adapt and tune advanced approaches in our algorithmic designs. The next activity, described in Section 3.3 , is related to the study of parallel linear algebra techniques that currently appear as promising approaches to tackle huge problems on extreme scale platforms. We highlight the linear problems (linear systems or eigenproblems) because they are in many large scale applications the main computational intensive numerical kernels and often the main performance bottleneck. These parallel numerical techniques will be the basis of both academic and industrial collaborations, some are described in Section 4.1 , but will also be closely related to some functionalities developed in the parallel fast multipole activity described in Section 3.4 . Finally, as the accuracy of the physical models increases, there is a real need to go for parallel efficient algorithm implementation for multiphysics and multiscale modeling in particular in the context of code coupling. The challenges associated with this activity will be addressed in the framework of the activity described in Section 3.5 .

Currently, we have one major application (see Section 4.1 ) that is in material physics. We will collaborate to all steps of the design of the parallel simulation tool. More precisely, our applied mathematics skill will contribute to the modelling, our advanced numerical schemes will help in the design and efficient software implementation for very large parallel simulations. We also participate to a few co-design actions in close collaboration with some applicative groups. The objective of this activity is to instantiate our expertise in fields where they are critical for designing scalable simulation tools. We refer to Section 4.2  for a detailed description of these activities.

## 3.2. High-performance computing on next generation architectures

**Participants:** Emmanuel Agullo, Olivier Beaumont, Olivier Coulaud, Pierre Esterie, Lionel Eyraud-Dubois, Mathieu Faverge, Luc Giraud, Abdou Guermouche, Gilles Marait, Pierre Ramet, Jean Roman, Nick Schenkels, Alena Shilova, Mathieu Verite.

The research directions proposed in <span style="color:red">HIEPACS</span>  are strongly influenced by both the applications we are studying and the architectures that we target (i.e., massively parallel heterogeneous many-core architectures, ...). Our main goal is to study the methodology needed to efficiently exploit the new generation of high-performance computers with all the constraints that it induces. To achieve this high-performance with complex applications we have to study both algorithmic problems and the impact of the architectures on the algorithm design.

From the application point of view, the project will be interested in multiresolution, multiscale and hierarchical approaches which lead to multi-level parallelism schemes. This hierarchical parallelism approach is necessary to achieve good performance and high-scalability on modern massively parallel platforms. In this context, more specific algorithmic problems are very important to obtain high performance. Indeed, the kind of applications we are interested in are often based on data redistribution for example (e.g., code coupling applications). This well-known issue becomes very challenging with the increase of both the number of computational nodes and the amount of data. Thus, we have both to study new algorithms and to adapt the

existing ones. In addition, some issues like task scheduling have to be restudied in this new context. It is important to note that the work developed in this area will be applied for example in the context of code coupling (see Section 3.5 ).

Considering the complexity of modern architectures like massively parallel architectures or new generation heterogeneous multicore architectures, task scheduling becomes a challenging problem which is central to obtain a high efficiency. With the recent addition of colleagues from the scheduling community (O. Beaumont and L. Eyraud-Dubois), the team is better equipped than ever to design scheduling algorithms and models specifically tailored to our target problems. It is important to note that this topic is strongly linked to the underlying programming model. Indeed, considering multicore and heterogeneous architectures, it has appeared, in the last five years, that the best programming model is an approach mixing multi-threading within computational nodes and message passing between them. In the last five years, a lot of work has been developed in the high-performance computing community to understand what is critic to efficiently exploit massively multicore platforms that will appear in the near future. It appeared that the key for the performance is firstly the granularity of the computations. Indeed, in such platforms the granularity of the parallelism must be small so that we can feed all the computing units with a sufficient amount of work. It is thus very crucial for us to design new high performance tools for scientific computing in this new context. This will be developed in the context of our solvers, for example, to adapt to this new parallel scheme. Secondly, the larger the number of cores inside a node, the more complex the memory hierarchy. This remark impacts the behavior of the algorithms within the node. Indeed, on this kind of platforms, NUMA effects will be more and more problematic. Thus, it is very important to study and design data-aware algorithms which take into account the affinity between computational threads and the data they access. This is particularly important in the context of our high-performance tools. Note that this work has to be based on an intelligent cooperative underlying run-time (like the tools developed by the Inria STORM Project-Team) which allows a fine management of data distribution within a node.

Another very important issue concerns high-performance computing using "heterogeneous" resources within a computational node. Indeed, with the deployment of the GPU and the use of more specific co-processors, it is important for our algorithms to efficiently exploit these new type of architectures. To adapt our algorithms and tools to these accelerators, we need to identify what can be done on the GPU for example and what cannot. Note that recent results in the field have shown the interest of using both regular cores and GPU to perform computations. Note also that in opposition to the case of the parallelism granularity needed by regular multicore architectures, GPU requires coarser grain parallelism. Thus, making both GPU and regular cores work all together will lead to two types of tasks in terms of granularity. This represents a challenging problem especially in terms of scheduling. From this perspective, we investigate new approaches for composing parallel applications within a runtime system for heterogeneous platforms.

In the context of scaling up, and particularly in the context of minimizing energy consumption, it is generally acknowledged that the solution lies in the use of heterogeneous architectures, where each resource is particularly suited to specific types of tasks, and in a fine control at the algorithmic level of data movements and the trade-offs to be made between computation and communication. In this context, we are particularly interested in the optimization of the training phase of deep convolutional neural networks which consumes a lot of memory and for which it is possible to exchange computations for data movements and memory occupation. We are also interested in the complexity introduced by resource heterogeneity itself, both from a theoretical point of view on the complexity of scheduling problems and from a more practical point of view on the implementation of specific kernels in dense or sparse linear algebra.

In order to achieve an advanced knowledge concerning the design of efficient computational kernels to be used on our high performance algorithms and codes, we will develop research activities first on regular frameworks before extending them to more irregular and complex situations. In particular, we will work first on optimized dense linear algebra kernels and we will use them in our more complicated direct and hybrid solvers for sparse linear algebra and in our fast multipole algorithms for interaction computations. In this context, we will participate to the development of those kernels in collaboration with groups specialized in dense linear algebra. In particular, we intend develop a strong collaboration with the group of Jack Dongarra

at the University of Tennessee and collaborating research groups. The objectives will be to develop dense linear algebra algorithms and libraries for multicore architectures in the context the `PLASMA` project and for `GPU` and hybrid multicore/GPU architectures in the context of the `MAGMA` project. A new solver has emerged from the associate team, Chameleon. While `PLASMA` and `MAGMA` focus on multicore and GPU architectures, respectively, Chameleon makes the most out of heterogeneous architectures thanks to task-based dynamic runtime systems.

A more prospective objective is to study the resiliency in the context of large-scale scientific applications for massively parallel architectures. Indeed, with the increase of the number of computational cores per node, the probability of a hardware crash on a core or of a memory corruption is dramatically increased. This represents a crucial problem that needs to be addressed. However, we will only study it at the algorithmic/application level even if it needed lower-level mechanisms (at OS level or even hardware level). Of course, this work can be performed at lower levels (at operating system) level for example but we do believe that handling faults at the application level provides more knowledge about what has to be done (at application level we know what is critical and what is not). The approach that we will follow will be based on the use of a combination of fault-tolerant implementations of the run-time environments we use (like for example `ULFM`) and an adaptation of our algorithms to try to manage this kind of faults. This topic represents a very long range objective which needs to be addressed to guaranty the robustness of our solvers and applications.

Finally, it is important to note that the main goal of HIEPACS  is to design tools and algorithms that will be used within complex simulation frameworks on next-generation parallel machines. Thus, we intend with our partners to use the proposed approach in complex scientific codes and to validate them within very large scale simulations as well as designing parallel solution in co-design collaborations.

## 3.3. High performance solvers for large linear algebra problems

**Participants:**  Emmanuel Agullo, Olivier Coulaud, Tony Delarue, Mathieu Faverge, Aurélien Falco, Marek Felsoci, Luc Giraud, Abdou Guermouche, Esragul Korkmaz, Gilles Marait, Van Gia Thinh Nguyen, Jean Rene Poirier, Pierre Ramet, Jean Roman, Cristobal Samaniego Alvarado, Guillaume Sylvand, Nicolas Venkovic, Yanfei Xiang.

Starting with the developments of basic linear algebra kernels tuned for various classes of computers, a significant knowledge on the basic concepts for implementations on high-performance scientific computers has been accumulated. Further knowledge has been acquired through the design of more sophisticated linear algebra algorithms fully exploiting those basic intensive computational kernels. In that context, we still look at the development of new computing platforms and their associated programming tools. This enables us to identify the possible bottlenecks of new computer architectures (memory path, various level of caches, inter processor or node network) and to propose ways to overcome them in algorithmic design. With the goal of designing efficient scalable linear algebra solvers for large scale applications, various tracks will be followed in order to investigate different complementary approaches. Sparse direct solvers have been for years the methods of choice for solving linear systems of equations, it is nowadays admitted that classical approaches are not scalable neither from a computational complexity nor from a memory view point for large problems such as those arising from the discretization of large 3D PDE problems. We will continue to work on sparse direct solvers on the one hand to make sure they fully benefit from most advanced computing platforms and on the other hand to attempt to reduce their memory and computational costs for some classes of problems where data sparse ideas can be considered. Furthermore, sparse direct solvers are a key building boxes for the design of some of our parallel algorithms such as the hybrid solvers described in the sequel of this section. Our activities in that context will mainly address preconditioned Krylov subspace methods; both components, preconditioner and Krylov solvers, will be investigated. In this framework, and possibly in relation with the research activity on fast multipole, we intend to study how emerging $\mathcal{H}$-matrix arithmetic can benefit to our solver research efforts.

### 3.3.1. Parallel sparse direct solvers

For the solution of large sparse linear systems, we design numerical schemes and software packages for direct and hybrid parallel solvers. Sparse direct solvers are mandatory when the linear system is very ill-conditioned; such a situation is often encountered in structural mechanics codes, for example. Therefore, to obtain an industrial software tool that must be robust and versatile, high-performance sparse direct solvers are mandatory, and parallelism is then necessary for reasons of memory capability and acceptable solution time. Moreover, in order to solve efficiently 3D problems with more than 50 million unknowns, which is now a reachable challenge with new multicore supercomputers, we must achieve good scalability in time and control memory overhead. Solving a sparse linear system by a direct method is generally a highly irregular problem that induces some challenging algorithmic problems and requires a sophisticated implementation scheme in order to fully exploit the capabilities of modern supercomputers.

New supercomputers incorporate many microprocessors which are composed of one or many computational cores. These new architectures induce strongly hierarchical topologies. These are called NUMA architectures. In the context of distributed NUMA architectures, in collaboration with the Inria STORM team, we study optimization strategies to improve the scheduling of communications, threads and I/O. We have developed dynamic scheduling designed for NUMA architectures in the `PaStiX` solver. The data structures of the solver, as well as the patterns of communication have been modified to meet the needs of these architectures and dynamic scheduling. We are also interested in the dynamic adaptation of the computation grain to use efficiently multi-core architectures and shared memory. Experiments on several numerical test cases have been performed to prove the efficiency of the approach on different architectures. Sparse direct solvers such as `PaStiX` are currently limited by their memory requirements and computational cost. They are competitive for small matrices but are often less efficient than iterative methods for large matrices in terms of memory. We are currently accelerating the dense algebra components of direct solvers using block low-rank compression techniques.

In collaboration with the ICL team from the University of Tennessee, and the STORM team from Inria, we are evaluating the way to replace the embedded scheduling driver of the `PaStiX` solver by one of the generic frameworks, `PaRSEC` or `StarPU`, to execute the task graph corresponding to a sparse factorization. The aim is to design algorithms and parallel programming models for implementing direct methods for the solution of sparse linear systems on emerging computer equipped with GPU accelerators. More generally, this work will be performed in the context of the ANR SOLHARIS project which aims at designing high performance sparse direct solvers for modern heterogeneous systems. This ANR project involves several groups working either on the sparse linear solver aspects (HIEPACS  and ROMA from Inria and APO from IRIT), on runtime systems (STORM from Inria) or scheduling algorithms (HIEPACS  and ROMA from Inria). The results of these efforts will be validated in the applications provided by the industrial project members, namely CEA-CESTA and Airbus Central R & T.

### 3.3.2. *Hybrid direct/iterative solvers based on algebraic domain decomposition techniques*

One route to the parallel scalable solution of large sparse linear systems in parallel scientific computing is the use of hybrid methods that hierarchically combine direct and iterative methods. These techniques inherit the advantages of each approach, namely the limited amount of memory and natural parallelization for the iterative component and the numerical robustness of the direct part. The general underlying ideas are not new since they have been intensively used to design domain decomposition techniques; those approaches cover a fairly large range of computing techniques for the numerical solution of partial differential equations (PDEs) in time and space. Generally speaking, it refers to the splitting of the computational domain into sub-domains with or without overlap. The splitting strategy is generally governed by various constraints/objectives but the main one is to express parallelism. The numerical properties of the PDEs to be solved are usually intensively exploited at the continuous or discrete levels to design the numerical algorithms so that the resulting specialized technique will only work for the class of linear systems associated with the targeted PDE.

In that context, we continue our effort on the design of algebraic non-overlapping domain decomposition techniques that rely on the solution of a Schur complement system defined on the interface introduced by the partitioning of the adjacency graph of the sparse matrix associated with the linear system. Although it is better conditioned than the original system the Schur complement needs to be precondition to be amenable to

a solution using a Krylov subspace method. Different hierarchical preconditioners will be considered, possibly multilevel, to improve the numerical behaviour of the current approaches implemented in our software library `MaPHyS`. This activity will be developed further developped in the H2020 EoCoE2 project. In addition to this numerical studies, advanced parallel implementation will be developed that will involve close collaborations between the hybrid and sparse direct activities.

### 3.3.3. Linear Krylov solvers

Preconditioning is the main focus of the two activities described above. They aim at speeding up the convergence of a Krylov subspace method that is the complementary component involved in the solvers of interest for us. In that framework, we believe that various aspects deserve to be investigated; we will consider the following ones:

- preconditioned block Krylov solvers for multiple right-hand sides. In many large scientific and industrial applications, one has to solve a sequence of linear systems with several right-hand sides given simultaneously or in sequence (radar cross section calculation in electromagnetism, various source locations in seismic, parametric studies in general, ...). For "simultaneous" right-hand sides, the solvers of choice have been for years based on matrix factorizations as the factorization is performed once and simple and cheap block forward/backward substitutions are then performed. In order to effectively propose alternative to such solvers, we need to have efficient preconditioned Krylov subspace solvers. In that framework, block Krylov approaches, where the Krylov spaces associated with each right-hand side are shared to enlarge the search space will be considered. They are not only attractive because of this numerical feature (larger search space), but also from an implementation point of view. Their block-structures exhibit nice features with respect to data locality and re-usability that comply with the memory constraint of multicore architectures. We will continue the numerical study and design of the block GMRES variant that combines inexact break-down detection, deflation at restart and subspace recycling. Beyond new numerical investigations, a software implementation to be included in our linear solver libray `Fabulous` originately developed in the context of the DGA HIBOX project and further developed in the LynCs (Linear Algebra, Krylov-subspace methods, and multi-grid solvers for the discovery of New Physics) sub-project of PRACE-6IP.

- Extension or modification of Krylov subspace algorithms for multicore architectures: finally to match as much as possible to the computer architecture evolution and get as much as possible performance out of the computer, a particular attention will be paid to adapt, extend or develop numerical schemes that comply with the efficiency constraints associated with the available computers. Nowadays, multicore architectures seem to become widely used, where memory latency and bandwidth are the main bottlenecks; investigations on communication avoiding techniques will be undertaken in the framework of preconditioned Krylov subspace solvers as a general guideline for all the items mentioned above.

### 3.3.4. Eigensolvers

Many eigensolvers also rely on Krylov subspace techniques. Naturally some links exist between the Krylov subspace linear solvers and the Krylov subspace eigensolvers. We plan to study the computation of eigenvalue problems with respect to the following two different axes:

- Exploiting the link between Krylov subspace methods for linear system solution and eigensolvers, we intend to develop advanced iterative linear methods based on Krylov subspace methods that use some spectral information to build part of a subspace to be recycled, either though space augmentation or through preconditioner update. This spectral information may correspond to a certain part of the spectrum of the original large matrix or to some approximations of the eigenvalues obtained by solving a reduced eigenproblem. This technique will also be investigated in the framework of block Krylov subspace methods.

- In the context of the calculation of the ground state of an atomistic system, eigenvalue computation is a critical step; more accurate and more efficient parallel and scalable eigensolvers are required.

### *3.3.5. Fast Solvers for FEM/BEM Coupling*

In this research project, we are interested in the design of new advanced techniques to solve large mixed dense/sparse linear systems, the extensive comparison of these new approaches to the existing ones, and the application of these innovative ideas on realistic industrial test cases in the domain of aeroacoustics (in collaboration with Airbus Central R & T).

- The use of $\mathcal{H}$-matrix solvers on these problems has been investigated in the context of the PhD of A. Falco. Airbus CR&T, in collaboration with Inria Bordeaux Sud-Ouest, has developed a task-based $\mathcal{H}$-matrix solver on top of the runtime engine StarPU. Ideas coming from the field of sparse direct solvers (such as nested dissection or symbolic factorization) have been tested within $\mathcal{H}$-matrices.

- The question of parallel scalability of task-based tools is an active subject of research, using new communication engine such as NewMadeleine, that will be investigated during this project, in conjunction with new algorithmic ideas on the task-based writing of $\mathcal{H}$-matrix algorithms.

- Naturally, comparison with existing tools will be performed on large realistic test cases. Coupling schemes between these tools and the hierarchical methods used in $\mathcal{H}$-matrix will be developed and benched as well.

## 3.4. High performance Fast Multipole Method for N-body problems

**Participants:**  Emmanuel Agullo, Olivier Coulaud, Pierre Esterie, Guillaume Sylvand.

In most scientific computing applications considered nowadays as computational challenges (like biological and material systems, astrophysics or electromagnetism), the introduction of hierarchical methods based on an octree structure has dramatically reduced the amount of computation needed to simulate those systems for a given accuracy. For instance, in the N-body problem arising from these application fields, we must compute all pairwise interactions among N objects (particles, lines, ...) at every timestep. Among these methods, the Fast Multipole Method (FMM) developed for gravitational potentials in astrophysics and for electrostatic (coulombic) potentials in molecular simulations solves this N-body problem for any given precision with $O(N)$ runtime complexity against $O(N^2)$ for the direct computation.

The potential field is decomposed in a near field part, directly computed, and a far field part approximated thanks to multipole and local expansions. We introduced a matrix formulation of the FMM that exploits the cache hierarchy on a processor through the Basic Linear Algebra Subprograms (BLAS). Moreover, we developed a parallel adaptive version of the FMM algorithm for heterogeneous particle distributions, which is very efficient on parallel clusters of SMP nodes. Finally on such computers, we developed the first hybrid MPI-thread algorithm, which enables to reach better parallel efficiency and better memory scalability. We plan to work on the following points in HIEPACS .

### *3.4.1. Improvement of calculation efficiency*

Nowadays, the high performance computing community is examining alternative architectures that address the limitations of modern cache-based designs. GPU (Graphics Processing Units) and the Cell processor have thus already been used in astrophysics and in molecular dynamics. The Fast Mutipole Method has also been implemented on GPU. We intend to examine the potential of using these forthcoming processors as a building block for high-end parallel computing in N-body calculations. More precisely, we want to take advantage of our specific underlying BLAS routines to obtain an efficient and easily portable FMM for these new architectures. Algorithmic issues such as dynamic load balancing among heterogeneous cores will also have to be solved in order to gather all the available computation power. This research action will be conduced on close connection with the activity described in Section 3.2 .

### *3.4.2. Non uniform distributions*

In many applications arising from material physics or astrophysics, the distribution of the data is highly non uniform and the data can grow between two time steps. As mentioned previously, we have proposed a hybrid MPI-thread algorithm to exploit the data locality within each node. We plan to further improve the load

balancing for highly non uniform particle distributions with small computation grain thanks to dynamic load balancing at the thread level and thanks to a load balancing correction over several simulation time steps at the process level.

### 3.4.3. Fast multipole method for dislocation operators

The engine that we develop will be extended to new potentials arising from material physics such as those used in dislocation simulations. The interaction between dislocations is long ranged ($O(1/r)$) and anisotropic, leading to severe computational challenges for large-scale simulations. Several approaches based on the FMM or based on spatial decomposition in boxes are proposed to speed-up the computation. In dislocation codes, the calculation of the interaction forces between dislocations is still the most CPU time consuming. This computation has to be improved to obtain faster and more accurate simulations. Moreover, in such simulations, the number of dislocations grows while the phenomenon occurs and these dislocations are not uniformly distributed in the domain. This means that strategies to dynamically balance the computational load are crucial to achieve high performance.

### 3.4.4. Fast multipole method for boundary element methods

The boundary element method (BEM) is a well known solution of boundary value problems appearing in various fields of physics. With this approach, we only have to solve an integral equation on the boundary. This implies an interaction that decreases in space, but results in the solution of a dense linear system with $O(N^3)$ complexity. The FMM calculation that performs the matrix-vector product enables the use of Krylov subspace methods. Based on the parallel data distribution of the underlying octree implemented to perform the FMM, parallel preconditioners can be designed that exploit the local interaction matrices computed at the finest level of the octree. This research action will be conduced on close connection with the activity described in Section 3.3 . Following our earlier experience, we plan to first consider approximate inverse preconditionners that can efficiently exploit these data structures.

## 3.5. Load balancing algorithms for complex simulations

**Participants:** Cyril Bordage, Aurélien Esnard, Pierre Ramet.

Many important physical phenomena in material physics and climatology are inherently complex applications. They often use multi-physics or multi-scale approaches, which couple different models and codes. The key idea is to reuse available legacy codes through a coupling framework instead of merging them into a stand-alone application. There is typically one model per different scale or physics and each model is implemented by a parallel code.

For instance, to model a crack propagation, one uses a molecular dynamic code to represent the atomistic scale and an elasticity code using a finite element method to represent the continuum scale. Indeed, fully microscopic simulations of most domains of interest are not computationally feasible. Combining such different scales or physics is still a challenge to reach high performance and scalability.

Another prominent example is found in the field of aeronautic propulsion: the conjugate heat transfer simulation in complex geometries (as developed by the CFD team of CERFACS) requires to couple a fluid/convection solver (AVBP) with a solid/conduction solver (AVTP). As the AVBP code is much more CPU consuming than the AVTP code, there is an important computational imbalance between the two solvers.

In this context, one crucial issue is undoubtedly the load balancing of the whole coupled simulation that remains an open question. The goal here is to find the best data distribution for the whole coupled simulation and not only for each stand-alone code, as it is most usually done. Indeed, the naive balancing of each code on its own can lead to an important imbalance and to a communication bottleneck during the coupling phase, which can drastically decrease the overall performance. Therefore, we argue that it is required to model the coupling itself in order to ensure a good scalability, especially when running on massively parallel architectures (tens of thousands of processors/cores). In other words, one must develop new algorithms and software implementation to perform a *coupling-aware* partitioning of the whole application. Another related problem is the problem of resource allocation. This is particularly important for the global coupling efficiency and

scalability, because each code involved in the coupling can be more or less computationally intensive, and there is a good trade-off to find between resources assigned to each code to avoid that one of them waits for the other(s). What does furthermore happen if the load of one code dynamically changes relatively to the other one? In such a case, it could be convenient to dynamically adapt the number of resources used during the execution.

There are several open algorithmic problems that we investigate in the HIEPACS  project-team. All these problems uses a similar methodology based upon the graph model and are expressed as variants of the classic graph partitioning problem, using additional constraints or different objectives.

### 3.5.1. Dynamic load-balancing with variable number of processors

As a preliminary step related to the dynamic load balancing of coupled codes, we focus on the problem of dynamic load balancing of a single parallel code, with variable number of processors. Indeed, if the workload varies drastically during the simulation, the load must be redistributed regularly among the processors. Dynamic load balancing is a well studied subject but most studies are limited to an initially fixed number of processors. Adjusting the number of processors at runtime allows one to preserve the parallel code efficiency or keep running the simulation when the current memory resources are exceeded. We call this problem, *MxN graph repartitioning*.

We propose some methods based on graph repartitioning in order to re-balance the load while changing the number of processors. These methods are split in two main steps. Firstly, we study the migration phase and we build a "good" migration matrix minimizing several metrics like the migration volume or the number of exchanged messages. Secondly, we use graph partitioning heuristics to compute a new distribution optimizing the migration according to the previous step results.

### 3.5.2. Load balancing of coupled codes

As stated above, the load balancing of coupled code is a major issue, that determines the performance of the complex simulation, and reaching high performance can be a great challenge. In this context, we develop new graph partitioning techniques, called *co-partitioning*. They address the problem of load balancing for two coupled codes: the key idea is to perform a "coupling-aware" partitioning, instead of partitioning these codes independently, as it is classically done. More precisely, we propose to enrich the classic graph model with *inter-edges*, which represent the coupled code interactions. We describe two new algorithms, and compare them to the naive approach. In the preliminary experiments we perform on synthetically-generated graphs, we notice that our algorithms succeed to balance the computational load in the coupling phase and in some cases they succeed to reduce the coupling communications costs. Surprisingly, we notice that our algorithms do not degrade significantly the global graph edge-cut, despite the additional constraints that they impose.

Besides this, our co-partitioning technique requires to use graph partitioning with *fixed vertices*, that raises serious issues with state-of-the-art software, that are classically based on the well-known recursive bisection paradigm (RB). Indeed, the RB method often fails to produce partitions of good quality. To overcome this issue, we propose a *new* direct $k$-way greedy graph growing algorithm, called KGGGP, that overcomes this issue and succeeds to produce partition with better quality than RB while respecting the constraint of fixed vertices. Experimental results compare KGGGP against state-of-the-art methods, such as `Scotch`, for real-life graphs available from the popular *DIMACS'10* collection.

### 3.5.3. Load balancing strategies for hybrid sparse linear solvers

Graph handling and partitioning play a central role in the activity described here but also in other numerical techniques detailed in sparse linear algebra Section. The Nested Dissection is now a well-known heuristic for sparse matrix ordering to both reduce the fill-in during numerical factorization and to maximize the number of independent computation tasks. By using the block data structure induced by the partition of separators of the original graph, very efficient parallel block solvers have been designed and implemented according to super-nodal or multi-frontal approaches. Considering hybrid methods mixing both direct and iterative solvers such as `MaPHyS`, obtaining a domain decomposition leading to a good balancing of both the size of domain interiors and the size of interfaces is a key point for load balancing and efficiency in a parallel context.

We intend to revisit some well-known graph partitioning techniques in the light of the hybrid solvers and design new algorithms to be tested in the `Scotch` package.

<div align="center">**STORM Project-Team**</div>

# 3. Research Program

## 3.1. Parallel Computing and Architectures

Following the current trends of the evolution of HPC systems architectures, it is expected that future Exascale systems (i.e. Sustaining $10^{18}$ flops) will have millions of cores. Although the exact architectural details and trade-offs of such systems are still unclear, it is anticipated that an overall concurrency level of $O(10^9)$ threads/tasks will probably be required to feed all computing units while hiding memory latencies. It will obviously be a challenge for many applications to scale to that level, making the underlying system sound like "embarrassingly parallel hardware."

From the programming point of view, it becomes a matter of being able to expose extreme parallelism within applications to feed the underlying computing units. However, this increase in the number of cores also comes with architectural constraints that actual hardware evolution prefigures: computing units will feature extra-wide SIMD and SIMT units that will require aggressive code vectorization or "SIMDization", systems will become hybrid by mixing traditional CPUs and accelerators units, possibly on the same chip as the AMD APU solution, the amount of memory per computing unit is constantly decreasing, new levels of memory will appear, with explicit or implicit consistency management, etc. As a result, upcoming extreme-scale system will not only require unprecedented amount of parallelism to be efficiently exploited, but they will also require that applications generate adaptive parallelism capable to map tasks over heterogeneous computing units.

The current situation is already alarming, since European HPC end-users are forced to invest in a difficult and time-consuming process of tuning and optimizing their applications to reach most of current supercomputers' performance. It will go even worse with the emergence of new parallel architectures (tightly integrated accelerators and cores, high vectorization capabilities, etc.) featuring unprecedented degree of parallelism that only too few experts will be able to exploit efficiently. As highlighted by the ETP4HPC initiative, existing programming models and tools won't be able to cope with such a level of heterogeneity, complexity and number of computing units, which may prevent many new application opportunities and new science advances to emerge.

The same conclusion arises from a non-HPC perspective, for single node embedded parallel architectures, combining heterogeneous multicores, such as the ARM big.LITTLE processor and accelerators such as GPUs or DSPs. The need and difficulty to write programs able to run on various parallel heterogeneous architectures has led to initiatives such as HSA, focusing on making it easier to program heterogeneous computing devices. The growing complexity of hardware is a limiting factor to the emergence of new usages relying on new technology.

## 3.2. Scientific and Societal Stakes

In the HPC context, simulation is already considered as a third pillar of science with experiments and theory. Additional computing power means more scientific results, and the possibility to open new fields of simulation requiring more performance, such as multi-scale, multi-physics simulations. Many scientific domains able to take advantage of Exascale computers, these "Grand Challenges" cover large panels of science, from seismic, climate, molecular dynamics, theoretical and astrophysics physics... Besides, embedded applications are also able to take advantage of these performance increase. There is still an on-going trend where dedicated hardware is progressively replaced by off-the-shelf components, adding more adaptability and lowering the cost of devices. For instance, Error Correcting Codes in cell phones are still hardware chips, but with the forthcoming 5G protocol, new software and adaptative solutions relying on low power multicores are also explored. New usages are also appearing, relying on the fact that large computing capacities are becoming more affordable and widespread. This is the case for instance with Deep Neural Networks where the training phase can be done

on supercomputers and then used in embedded mobile systems. The same consideration applies for big data problems, of internet of things, where small sensors provide large amount of data that need to be processed in short amount of time. Even though the computing capacities required for such applications are in general a different scale from HPC infrastructures, there is still a need in the future for high performance computing applications.

However, the outcome of new scientific results and the development of new usages for mobile, embedded systems will be hindered by the complexity and high level of expertise required to tap the performance offered by future parallel heterogeneous architectures.

## 3.3. Towards More Abstraction

As emphasized by initiatives such as the European Exascale Software Initiative (EESI), the European Technology Platform for High Performance Computing (ETP4HPC), or the International Exascale Software Initiative (IESP), the HPC community needs new programming APIs and languages for expressing heterogeneous massive parallelism in a way that provides an abstraction of the system architecture and promotes high performance and efficiency. The same conclusion holds for mobile, embedded applications that require performance on heterogeneous systems.

This crucial challenge given by the evolution of parallel architectures therefore comes from this need to make high performance accessible to the largest number of developers, abstracting away architectural details providing some kind of performance portability, and provided a high level feed-back allowing the user to correct and tune the code. Disruptive uses of the new technology and groundbreaking new scientific results will not come from code optimization or task scheduling, but they require the design of new algorithms that require the technology to be tamed in order to reach unprecedented levels of performance.

Runtime systems and numerical libraries are part of the answer, since they may be seen as building blocks optimized by experts and used as-is by application developers. The first purpose of runtime systems is indeed to provide *abstraction*. Runtime systems offer a uniform programming interface for a specific subset of hardware (e.g., OpenGL or DirectX are well-established examples of runtime systems dedicated to hardware-accelerated graphics) or low-level software entities (e.g., POSIX-thread implementations). They are designed as thin user-level software layers that complement the basic, general purpose functions provided by the operating system calls. Applications then target these uniform programming interfaces in a portable manner. Low-level, hardware dependent details are hidden inside runtime systems. The adaptation of runtime systems is commonly handled through drivers. The abstraction provided by runtime systems thus enables portability. Abstraction alone is however not enough to provide portability of performance, as it does nothing to leverage low-level-specific features to get increased performance and does nothing to help the user tune his code. Consequently, the second role of runtime systems is to *optimize* abstract application requests by dynamically mapping them onto low-level requests and resources as efficiently as possible. This mapping process makes use of scheduling algorithms and heuristics to decide the best actions to take for a given metric and the application state at a given point in its execution time. This allows applications to readily benefit from available underlying low-level capabilities to their full extent without breaking their portability. Thus, optimization together with abstraction allows runtime systems to offer portability of performance. Numerical libraries provide sets of highly optimized kernels for a given field (dense or sparse linear algebra, FFT, etc.) either in an autonomous fashion or using an underlying runtime system.

Application domains cannot resort to libraries for all codes however, computation patterns such as stencils are a representative example of such difficulty. The compiler technology plays here a central role, in managing high level semantics, either through templates, domain specific languages or annotations. Compiler optimizations, and the same applies for runtime optimizations, are limited by the level of semantics they manage. Providing part of the algorithmic knowledge of an application, for instance knowing that it computes a 5-point stencil and then performs a dot product, would lead to more opportunities to adapt parallelism, memory structures, and is a way to leverage the evolving hardware. Besides, with the need for automatic optimization comes the need for *feed-back* to the user, corresponding to the need to debug the code and also to understand what the runtime

has performed. Here the compiler plays also a central role in the analysis of the code, and the instrumentation of the program given to the runtime.

Compilers and runtime play a crucial role in the future of high performance applications, by defining the input language for users, and optimizing/transforming it into high performance code. The objective of STORM is to propose better interactions between compiler and runtime and more semantics for both approaches.

The results of the team on-going research in 2019 reflect this focus. Results presented in Sections 7.11 , 7.15 , 7.10  and 7.9  correspond to efforts for higher abstractions through DSL or libraries, and decouple algorithmics from parallel optimizations. Results in Section 7.8  correspond to efforts on parallelism expression and again abstraction, starting from standard parallel programming languages. Results described in Sections 7.1  and 7.16 provide feed-back information, through visualization and deadlock detection for parallel executions. The work described in Sections 7.3 , 7.4 , 7.5 , 7.6 ,7.12 , 7.7  and 7.13  focus in particular on StarPU and its development in order to better abstract architecture, resilience and optimizations. The work presented Section 7.2  aims to help developers with optimization.

Finally, Sections 7.14   and 7.17   present an on-going effort on improving the Chameleon library and strengthening its relation with StarPU and the NewMadeleine communication library. They represent real-life applications for the runtime methods we develop.

# TADAAM Project-Team

# 3. Research Program

## 3.1. Need for System-Scale Optimization

Firstly, in order for applications to make the best possible use of the available resources, it is impossible to expose all the low-level details of the hardware to the program, as it would make impossible to achieve portability. Hence, the standard approach is to add intermediate layers (programming models, libraries, compilers, runtime systems, etc.) to the software stack so as to bridge the gap between the application and the hardware. With this approach, optimizing the application requires to express its parallelism (within the imposed programming model), organize the code, schedule and load-balance the computations, etc. In other words, in this approach, the way the code is written and the way it is executed and interpreted by the lower layers drives the optimization. In any case, this approach is centered on how computations are performed. Such an approach is therefore no longer sufficient, as the way an application is executing does depend less and less on the organization of computation and more and more on the way its data is managed.

Secondly, modern large-scale parallel platforms comprise tens to hundreds of thousand nodes [0]. However, very few applications use the whole machine. In general, an application runs only on a subset of the nodes [0]. Therefore, most of the time, an application shares the network, the storage and other resources with other applications running concurrently during its execution. Depending on the allocated resources, it is not uncommon that the execution of one application interferes with the execution of a neighboring one.

Lastly, even if an application is running alone, each element of the software stack often performs its own optimization independently. For instance, when considering an hybrid MPI/OpenMP application, one may realize that threads are concurrently used within the OpenMP runtime system, within the MPI library for communication progression, and possibly within the computation library (BLAS) and even within the application itself (pthreads). However, none of these different classes of threads are aware of the existence of the others. Consequently, the way they are executed, scheduled, prioritized does not depend on their relative roles, their locations in the software stack nor on the state of the application.

The above remarks show that in order to go beyond the state-of-the-art, it is necessary to design a new set of mechanisms allowing cross-layer and system-wide optimizations so as to optimize the way data is allocated, accessed and transferred by the application.

## 3.2. Scientific Challenges and Research Issues

In TADAAM, we will tackle the problem of efficiently executing an application, at system-scale, on an HPC machine. We assume that the application is already optimized (efficient data layout, use of effective libraries, usage of state-of-the-art compilation techniques, etc.). Nevertheless, even a statically optimized application will not be able to be executed at scale without considering the following dynamic constraints: machine topology, allocated resources, data movement and contention, other running applications, access to storage, etc. Thanks to the proposed layer, we will provide a simple and efficient way for already existing applications, as well as new ones, to express their needs in terms of resource usage, locality and topology, using a high-level semantic.

It is important to note that we target the optimization of each application independently but also several applications at the same time and at system-scale, taking into account their resource requirement, their network usage or their storage access. Furthermore, dealing with code-coupling application is an intermediate use-case that will also be considered.

---

[0]More than 22,500 XE6 compute node for the BlueWaters system; 5040 B510 Bullx Nodes for the Curie machine; more than 49,000 BGQ nodes for the MIRA machine.

[0]In 2014, the median case was 2048 nodes for the BlueWaters system and, for the first year of the Curie machine, the median case was 256 nodes

Several issues have to be considered. The first one consists in providing relevant **abstractions and models to describe the topology** of the available resources **and the application behavior**.

Therefore, the first question we want to answer is: **"How to build scalable models and efficient abstractions enabling to understand the impact of data movement, topology and locality on performance?"** These models must be sufficiently precise to grasp the reality, tractable enough to enable efficient solutions and algorithms, and simple enough to remain usable by non-hardware experts. We will work on (1) better describing the memory hierarchy, considering new memory technologies; (2) providing an integrated view of the nodes, the network and the storage; (3) exhibiting qualitative knowledge; (4) providing ways to express the multi-scale properties of the machine. Concerning abstractions, we will work on providing general concepts to be integrated at the application or programming model layers. The goal is to offer means, for the application, to express its high-level requirements in terms of data access, locality and communication, by providing abstractions on the notion of hierarchy, mesh, affinity, traffic metrics, etc.

In addition to the abstractions and the aforementioned models we need to **define a clean and expressive API in a scalable way**, in order for applications to express their needs (memory usage, affinity, network, storage access, model refinement, etc.).

Therefore, the second question we need to answer is: "**how to build a system-scale, stateful, shared layer that can gather applications needs expressed with a high-level semantic?**". This work will require not only to define a clean API where applications will express their needs, but also to define how such a layer will be shared across applications and will scale on future systems. The API will provide a simple yet effective way to express different needs such as: memory usage of a given portion of the code; start of a compute intensive part; phase where the network is accessed intensively; topology-aware affinity management; usage of storage (in read and/or write mode); change of the data layout after mesh refinement, etc. From an engineering point of view, the layer will have a hierarchical design matching the hardware hierarchy, so as to achieve scalability.

Once this has been done, the service layer, will have all the information about the environment characteristics and application requirements. We therefore need to design a set of **mechanisms to optimize applications execution**: communication, mapping, thread scheduling, data partitioning/mapping/movement, etc.

Hence, the last scientific question we will address is: "**How to design fast and efficient algorithms, mechanisms and tools to enable execution of applications at system-scale, in full a HPC ecosystem, taking into account topology and locality?**" A first set of research is related to thread and process placement according to the topology and the affinity. Another large field of study is related to data placement, allocation and partitioning: optimizing the way data is accessed and processed especially for mesh-based applications. The issues of transferring data across the network will also be tackled, thanks to the global knowledge we have on the application behavior and the data layout. Concerning the interaction with other applications, several directions will be tackled. Among these directions we will deal with matching process placement with resource allocation given by the batch scheduler or with the storage management: switching from a best-effort application centric strategy to global optimization scheme.

# Auctus Team

# 3. Research Program

## 3.1. Analysis and modelling of human behavior

### 3.1.1. Scientific Context

The purpose of this axis is to provide metrics to assess human behavior. We place ourselves here from the point of view of the human being and more precisely of the industrial operator. We assume the following working hypotheses: the operator's task and environmental conditions are known and circumscribed; the operator is trained in the task, production tools and safety instructions; the task is repeated with more or less frequent intervals. We focus our proposals on assessing:

- the physical and cognitive fragility of operators in order to meet assistance needs;
- cognitive biases and physical constraints leading to a loss of operator safety;
- ergonomic, performance and acceptance of the production tool.

In the industrial context, the fields that best answer these questions are work ergonomics and cognitive sciences. Scientists typically work on 4 axes: physiological/biomechanical, cognitive, psychological and sociological. More specifically, we focus on biomechanical, cognitive and psychological aspects, as described by the ANACT [24], [26]. The aim here is to translate these factors into metrics, optimality criteria or constraints in order to implement them in our methodologies for analysis, design and control of the collaborative robot.

To understand our desired contributions in robotics, we must review the current state of ergonomic workstation evaluation, particularly at the biomechanical level. The ergonomist evaluates the gesture through the observation of workstations and, generally, through questionnaires. This requires long periods of field observation, followed by analyses based on ergonomic grids (e.g. RULA [42], REBA [32], LUBA [37], OWAS [36], ROSA [60],...). Until then, the use of more complex measurement systems was reserved for laboratories, particularly biomechanical laboratories. The appearance of inexpensive sensors such as IMUs (Inertial Measurement Units) or RGB-D cameras makes it possible to consider a digitalized, and therefore objective, observation of the gesture, postures and more generally of human movement. Thanks to these sensors, which are more or less intrusive, it is now possible to permanently install observation systems on production lines. This completely changes paradigms and opens the door to longitudinal observations. It should be noted that this is comparable to the evolution of maintenance, which becomes predictive.

On the strength of this new paradigm, *ergonomic robotics* has recently taken an interest in this type of evaluation to adapt the robot's movements in order to reduce ergonomic risk scores. This approach complements the more traditional approaches that only consider the performance of the action produced by the human in interaction with the robot. However, we must go further. Indeed, the ergonomic criteria are based on the principle that the comfort positions are distant from the human articular stops. In addition, the notation must be compatible with an observation of the human being through the eye of the ergonomist. In practice, evaluations are inaccurate and subjective [63]. Moreover, they are made for quasi-static human positions without taking into account the evolution of the person's physical, physiological and psychological state. The repetition of gestures, the solicitation of muscles and joints is one of the questions that must complete these analyses. One of the methods used by ergonomists to limit biomechanical exposures is to increase variations in motor stress by rotating tasks [61]. However, this type of extrinsic method is not always possible in the industrial context [40].

One of Auctus' objectives is to show how, through a cobot, the operator's environment can be varied to encourage more appropriate motor strategies. To do so, we must focus on a field of biomechanics that studies the intrinsic variability of the motor system allowed by the joint redundancy of the human body. This motor variability refers to the natural alternation of postures, movements and muscle activity observed in the individual to respond to a requested task [61]. This natural variation leads to differences between the motor coordinates used by individuals, which evokes the notion of motor strategy [33].

As shown by the cognitive dimension of ergonomics (see above), we believe that some of these motor strategies are a physically quantifiable reflection of the operator's cognitive state. For example, fatigue [57] and its anticipation or the manual expertise (dexterous and cognitive) of the operator which allows him to anticipate his movements over long periods of time in order to preserve his body, his performance and his pain.

### 3.1.2. *Methodology*

How can we observe, understand and quantify these human motor strategies to better design and control the behavior of the cobotic assistant? When we study the systems of equations considered (kinematic, static, dynamic, musculoskeletal), several problems appear and explain our methodological choices:

- the large dimensions of the problems to be considered, due to joint, muscle and placement redundancy,
- the variabilities of the parameters, for example: physiological (consider not an operator, but a set of operators), geometric (consider a set of possible placements of the operator) and static (consider a set of forces that the operator must produce);
- the uncertainties of measurement, model approximation.

The idea is to start from a description of redundant workspaces (geometric, static, dynamic...). To do this, we use set theory approaches, based on interval analysis [3], [48], which allow us to respond to the uncertainties and variability issues previously mentioned. In addition, one of the advantages of these techniques is that they allow the results to be certified, which is essential to address safety issues. Some members of the team has already achieved success in mechanical design for performance certification and robot design [44]. The adaptation of these approaches allows us to obtain a mapping of ergonomic and efficient movements in which we can project the operators' motor strategies and thus define a metric quantifying the sensorimotor commands chosen with regard to the cognitive criteria studied.

It is therefore necessary to:

- propose new indices linking different types of performance (ergonomic biomechanical robotics, but also influence of fatigue, stress, level of expertise on the evolution of performance);
- divide the gesture into homogeneous phases: this process is complex and depends on the type of index used and the techniques used. We are exploring several ways: inverse optimal control, learning methods, or the use of techniques from signal processing.
- develop interval extensions of the identified indices. These indices are not necessarily the result of a direct model, and algorithms need to be developed or adapted (calculation of manipulability, UCM, etc.).
- Aggregate proposals into a dedicated interval analysis library (use of and contribution to the existing ALIAS-Inria and the open source IBEX library).

The originality and contribution of the methodology is to allow an analysis taking into account in the same model the measurement uncertainties (important for on-site use of analytical equipment), the variability of tasks and trajectories, and the physiological characteristics of the operators.

Other avenues of research are being explored, particularly around the inverse optimal control [49] which allows us to project human movement on the basis of performance indices and thus to offer a possible interpretation in the analysis of behaviors.

We also use automatic classification techniques: 1) to propose cognitive models that will be learned experimentally 2) for segmentation or motion recognition, for example by testing Reservoir Computing [34] approaches.

## 3.2. Operator / robot coupling

### 3.2.1. *Scientific Context*

Thanks to the progress made in recent years in the field of p-HRI (Physical Human-Robot Interaction), robotic systems are beginning to operate in the same workspace as humans, which is profoundly changing industrial

issues and allowing a wide variety of human-robot coupling solutions to be considered to perform the same task [25]. Different types of interactions exist. They can be classified in different ways: according to the degree of autonomy of the robot and its proximity to the user [31] with particularities for "wearable robots" [30], [29], or for collaborative robotics [62], or according to the role of the human being [58]. From a cognitive point of view, classifications are more concerned with autonomy, the complexity of information processing and the type of communication and representation of the human being by the robot [47], [64].

We proposed a classification of cobotic systems according to the configuration of the schema of interactions between humans, robots and the environment [45], [55].

The parameters of the coupling being numerous and complex, the determination of the most appropriate type of coupling for a type of problem is an open problem [50], [2], [46], [41]. The traditional approach consists in trying to identify and classify the various possible options and to select the one that seems most relevant with regard to the feasibility, efficiency, budget envelope and acceptability of the operator. One of the main objectives of our research project is to define a typology of cobots or cobotic systems in order to specify the methodology for developing the best solution: what are the criteria for defining the best robotic architecture, what type of coupling, what autonomy of the robot, what role for the operator, what risks for the human, what overall performance? These are the key issues that need to be addressed. To meet this methodological need, we propose an approach guided by experience on use cases obtained thanks to our industrial partners.

### *3.2.2. Methodology*

Task analysis and human behavior modelling, discussed in the previous sections, should help to characterize the different types of coupling and interaction modalities, their advantages and disadvantages, in order to assist in the decision-making process. One of the ideas we would like to develop is to try to break down the task into a sequence of elementary gestures corresponding to simple motor actions performed in a clearly identified context and to evaluate for each of them the degree of feasibility in automatic mode or in robot assistance mode. The assessment must take into account a large number of parameters that relate to physical interactions, human-robot communication, reliability and human factors, including acceptability and impact on the valuation or devaluation of the operator's work. Concerning the evaluation of human factors, we have already begun to work on the subject within the more general framework of human systems interactions by operating Bayesian networks, drawing inspiration from the work of [28], [56].

The adoption of assessment criteria for a single domain (e. g. robotics or ergonomics) cannot guarantee that the performance of this coupling will be maximized. From design to evaluation, cross-effects must be constantly considered:

- impact of the cobot design on the user's performance: intuitiveness, adaptation to intra- and inter-individual variations, affordance, stress factors (noise, vibrations,...), fatigue factors (control laws, necessary attention,...) and motivation factors (effectiveness, efficiency, aesthetics,...);
- impact of user performance on cobot exploitation: risks of human error (attention error, perseveration, circumvention of procedures, syndrome outside the loop) [28].

In addition to purely physical assistance, some cobotic systems are designed to assist the operator in his decision-making. The issues of trust, acceptance, sharing of representations and co-construction of a shared awareness of the situation are then to be addressed [59].

## 3.3. Design of cobotic systems

### *3.3.1. Architectural design*

Is it necessary to cobotize, robotize or assist the human being? Which mechanical architecture meets the task challenges (a serial cobot, a specific mechanism, an exoskeleton)? What type of interaction (H/R cohabitation, comanipulation, teleoperation)? These questions are the first requests from our industrial partners. For the moment, we have few comprehensive methodological answers to provide them. Choosing a collaborative robot architecture is a difficult problem [38]. It is all the more when the questions are approached from both a cognitive ergonomics and robotics perspective. There are indeed major methodological and conceptual

differences in these areas. It is therefore necessary to bridge these representational gaps and to propose an approach that takes into consideration the expectations of the robotician to model and formalize the general properties of a cobotic system as well as those of the ergonomist to define the expectations in terms of an assistance tool.

To do this, we propose a user-centered design approach, with a particular focus on human-system interactions. From a methodological point of view, this requires first of all the development of a structured experimental approach aimed at characterizing the task to be carried out through a "system" analysis but also at capturing the physical markers of its realization: movements and efforts required, ergonomic stress. This characterization must be done through the prism of the systematic study of the exchange of information (and their nature) by humans in their performance of the considered task. On the basis of these analyses, the main challenge is to define a decision support tool for the choice of the robotic architecture and for the specifications of the role assigned to the robot and the operator as well as their interactions.

The evolution of the chosen methodology is for the moment empirical, based on the user cases regularly treated in the team (see sections on contracts and partnerships).

It can be summarized for the moment as:

- identify difficult jobs on industrial sites. This is done through visits and exchanges with our partners (manager, production manager, ergonomist...);

- select some of them, then observe the human in its ecological environment. Our tools allow us to produce a motion analysis, currently based on ergonomic criteria. In parallel we carry out a physical evaluation of the task in terms of expected performance and an evaluation of the operator by means of questionnaires.

- Synthesize these first results to deduce the robotic architectures to be initiated, the key points of human-robot interaction to be developed, the difficulties in terms of human factors to be taken into account.

In addition, the different human and task analyses take advantage of the different expertise available within the team. We would like to gradually introduce the evaluation criteria presented above. Indeed, the team has already worked on the current dominant approach: the use of a virtual human to design the cobotic cell through virtual tools [1]. However, the very large dimensions of the problems treated (modelling of the body's ddl and the constraints applied to it) makes it difficult to carry out a certified analysis. We then choose to go through the calculation of the body's workspace, representing its different performances, which is not yet done in this field. The idea here is to apply set theory approaches, using interval analysis and already discussed in section 3.1.2 . The goal is then to extend to intervals the constraints played in virtual reality during the simulation. This would allow the operator to check his trajectories and scenarios not only for a single case study but also for sets of cases. For example, it can be verified that, regardless of the bounded sets of simulated operator physiologies, the physical constraints of a simulated trajectory are not violated. Thus, the assisted design tools certify cases of use as a whole. Moreover, the intersection between the human and robot workspaces provides the necessary constraints to certify the feasibility of a task. This allows us to better design a cobotic system to integrate physical constraints. In the same way, we will look for ways in which human cognitive markers can be included in this approach.

Thus, we summarize here the contributions of the other research axes, from the analysis of human behavior in its environment for an identified task, to the choice of a mechanical architecture, via an evaluation of torque and interactions. All the previous analyses provide design constraints. This methodological approach is perfectly integrated into an Appropriate Design approach used for the dimensional design of robots, again based on interval analysis. Indeed, to the desired performance of the human-robot couple in relation to a task, it is sufficient to add the constraints limiting the difficulty of the operator's gesture as described above. The challenges are then the change of scale in models that symbiotically consider the human-robot pair, the uncertain, flexible and uncontrollable nature of human behavior and the many evaluation indices needed to describe them.

### *3.3.2. Control design*

The control laws of collaborative robots from the major robot manufacturers differ little or not at all from the existing control laws in the field of conventional industrial robotics. Security is managed a posteriori, as an exception, by a security PLC / PC. It is therefore not an intrinsic property of the controller. This quite strongly restricts the possibilities of physical interaction [0] and collaboration and leads to sub-optimal operation of the robotic system. It is difficult in this context to envision real human-robot collaboration. Collaborative operation requires, in this case, a control calculation that integrates safety and ergonomics as a priori constraints.

The control of truly collaborative robots in an industrial context is, from our point of view, underpinned by two main issues. The first is related to the macroscopic adaptation of the robot's behaviour according to the phases of the production process. The second is related to the fine adaptation of the degree and/or nature of the robot's assistance according to the ergonomic state of the operator. If this second problem is part of a historical dynamic in robotics that consists in placing safety constraints, particularly those related to the presence of a human being, at the heart of the control problem [31] [43], [35], it is not approached from the more subtle point of view of ergonomics where the objective cannot be translated only in terms of human life or death, but rather in terms of long-term respect for their physical and mental integrity. Thus, the simple and progressive adoption by a human operator of the collaborative robot intended to assist him in his gesture requires a self-adaptation in the time of the command. This self-adaptation is a fairly new subject in the literature [51], [52].

At the macroscopic level, the task plan to be performed for a given industrial operation can be represented by a finite state machine. In order not to increase the human's cognitive load by explicitly asking him to manage transitions for the robot, we propose to develop a decision algorithm to ensure discrete transitions from one task (and the associated assistance mode) to another based on an online estimate of the current state of the human-robot couple. The associated scientific challenge requires establishing a link between the robot's involvement and a given working situation. To do so, we propose an incremental approach to learning this complex relationship. The first stage of this work will consist in identifying the general and relevant control variables to conduct this learning in an efficient and reusable way, regardless of the particular method of calculating the control action. Physically realistic simulations and real word experiments will be used to feed this learning process.

In order to handle mode transitions, we propose to explore the richness of the multi-tasking control formalism under constraints [39] in order to ensure a continuous transition from one control mode to another while guaranteeing compliance with a certain number of control constraints. Some of these constraints are based on ergonomic specifications and are dependent on the state of the robot and of the human operator, which, by nature, is difficult to predict accurately. We propose to exploit the interval analysis paradigm to efficiently formulate ergonomic constraints robust to the various existing uncertainties.

Purely discrete or reactive adaptation of the control law would make no sense given the slow dynamics of certain physiological phenomena such as fatigue. Thus, we propose to formulate the control problem as a predictive problem where the impact of the control decision at a time $t$ is anticipated at different time horizons. This requires a prediction of human movement and knowledge of the motor variability strategies it employs. This prediction is possible on the basis of the supervision at all times of the operational objectives (task in progress) in the short term. However, it requires the use of a virtual human model and possibly a dynamic simulation to quantify the impact of these potential movements in terms of performance, including ergonomics. It is impractical to use a predictive command with simulation in the loop with an advanced virtual manikin model. We therefore propose to adapt the prediction horizon and the complexity of the corresponding model in order to guarantee a reasonable computational complexity.

The planned developments require both an approach to modelling human sensorimotor behaviour, particularly in terms of accommodating fatigue via motor variability, and validating related models in an experimental framework based on observation of movement and quantification of ergonomic performance. Experimental developments must also focus on the validation of proposed control approaches in concrete contexts. To begin with, the Woobot project related to gesture assistance for carpenters (Nassim Benhahib's thesis) and

---

[0]In the ISO TS 15066 technical specification on collaborative robotics, human-robot physical interaction is allowed but perceived as a situation to be avoided.

a collaboration currently being set up with Safran on assistance to operators in shrink-wrapping tasks (manual knotting) in aeronautics are rich enough background elements to support the research conducted. Collaborative research projects with PSA will also soon provide a larger set of contexts in which the proposed research can be validated.

# 3. Research Program

## 3.1. Research Program

Research in artificial intelligence, machine learning and pattern recognition has produced a tremendous amount of results and concepts in the last decades. A blooming number of learning paradigms - supervised, unsupervised, reinforcement, active, associative, symbolic, connectionist, situated, hybrid, distributed learning... - nourished the elaboration of highly sophisticated algorithms for tasks such as visual object recognition, speech recognition, robot walking, grasping or navigation, the prediction of stock prices, the evaluation of risk for insurances, adaptive data routing on the internet, etc... Yet, we are still very far from being able to build machines capable of adapting to the physical and social environment with the flexibility, robustness, and versatility of a one-year-old human child.

Indeed, one striking characteristic of human children is the nearly open-ended diversity of the skills they learn. They not only can improve existing skills, but also continuously learn new ones. If evolution certainly provided them with specific pre-wiring for certain activities such as feeding or visual object tracking, evidence shows that there are also numerous skills that they learn smoothly but could not be "anticipated" by biological evolution, for example learning to drive a tricycle, using an electronic piano toy or using a video game joystick. On the contrary, existing learning machines, and robots in particular, are typically only able to learn a single pre-specified task or a single kind of skill. Once this task is learnt, for example walking with two legs, learning is over. If one wants the robot to learn a second task, for example grasping objects in its visual field, then an engineer needs to re-program manually its learning structures: traditional approaches to task-specific machine/robot learning typically include engineer choices of the relevant sensorimotor channels, specific design of the reward function, choices about when learning begins and ends, and what learning algorithms and associated parameters shall be optimized.

As can be seen, this requires a lot of important choices from the engineer, and one could hardly use the term "autonomous" learning. On the contrary, human children do not learn following anything looking like that process, at least during their very first years. Babies develop and explore the world by themselves, focusing their interest on various activities driven both by internal motives and social guidance from adults who only have a folk understanding of their brains. Adults provide learning opportunities and scaffolding, but eventually young babies always decide for themselves what activity to practice or not. Specific tasks are rarely imposed to them. Yet, they steadily discover and learn how to use their body as well as its relationships with the physical and social environment. Also, the spectrum of skills that they learn continuously expands in an organized manner: they undergo a developmental trajectory in which simple skills are learnt first, and skills of progressively increasing complexity are subsequently learnt.

A link can be made to educational systems where research in several domains have tried to study how to provide a good learning experience to learners. This includes the experiences that allow better learning, and in which sequence they must be experienced. This problem is complementary to that of the learner that tries to learn efficiently, and the teacher here has to use as efficiently the limited time and motivational resources of the learner. Several results from psychology [59] and neuroscience [85] have argued that the human brain feels intrinsic pleasure in practicing activities of optimal difficulty or challenge. A teacher must exploit such activities to create positive psychological states of flow [73].

A grand challenge is thus to be able to build machines that possess this capability to discover, adapt and develop continuously new know-how and new knowledge in unknown and changing environments, like human children. In 1950, Turing wrote that the child's brain would show us the way to intelligence: "Instead of trying to produce a program to simulate the adult mind, why not rather try to produce one which simulates the child's" [154]. Maybe, in opposition to work in the field of Artificial Intelligence who has focused on mechanisms trying to match the capabilities of "intelligent" human adults such as chess playing or natural language

dialogue [91], it is time to take the advice of Turing seriously. This is what a new field, called developmental (or epigenetic) robotics, is trying to achieve [109] [158]. The approach of developmental robotics consists in importing and implementing concepts and mechanisms from developmental psychology [117], cognitive linguistics [72], and developmental cognitive neuroscience [96] where there has been a considerable amount of research and theories to understand and explain how children learn and develop. A number of general principles are underlying this research agenda: embodiment [63] [131], grounding [89], situatedness [49], self-organization [150] [132], enaction [156], and incremental learning [67].

Among the many issues and challenges of developmental robotics, two of them are of paramount importance: exploration mechanisms and mechanisms for abstracting and making sense of initially unknown sensorimotor channels. Indeed, the typical space of sensorimotor skills that can be encountered and learnt by a developmental robot, as those encountered by human infants, is immensely vast and inhomogeneous. With a sufficiently rich environment and multimodal set of sensors and effectors, the space of possible sensorimotor activities is simply too large to be explored exhaustively in any robot's life time: it is impossible to learn all possible skills and represent all conceivable sensory percepts. Moreover, some skills are very basic to learn, some other very complicated, and many of them require the mastery of others in order to be learnt. For example, learning to manipulate a piano toy requires first to know how to move one's hand to reach the piano and how to touch specific parts of the toy with the fingers. And knowing how to move the hand might require to know how to track it visually.

Exploring such a space of skills randomly is bound to fail or result at best on very inefficient learning [128]. Thus, exploration needs to be organized and guided. The approach of epigenetic robotics is to take inspiration from the mechanisms that allow human infants to be progressively guided, i.e. to develop. There are two broad classes of guiding mechanisms which control exploration:

1. **internal guiding mechanisms,** and in particular intrinsic motivation, responsible of spontaneous exploration and curiosity in humans, which is one of the central mechanisms investigated in FLOWERS, and technically amounts to achieve online active self-regulation of the growth of complexity in learning situations;

2. **social learning and guidance,** a learning mechanisms that exploits the knowledge of other agents in the environment and/or that is guided by those same agents. These mechanisms exist in many different forms like emotional reinforcement, stimulus enhancement, social motivation, guidance, feedback or imitation, some of which being also investigated in FLOWERS;

### 3.1.1. Internal guiding mechanisms

In infant development, one observes a progressive increase of the complexity of activities with an associated progressive increase of capabilities [117], children do not learn everything at one time: for example, they first learn to roll over, then to crawl and sit, and only when these skills are operational, they begin to learn how to stand. The perceptual system also gradually develops, increasing children perceptual capabilities other time while they engage in activities like throwing or manipulating objects. This make it possible to learn to identify objects in more and more complex situations and to learn more and more of their physical characteristics.

Development is therefore progressive and incremental, and this might be a crucial feature explaining the efficiency with which children explore and learn so fast. Taking inspiration from these observations, some roboticists and researchers in machine learning have argued that learning a given task could be made much easier for a robot if it followed a developmental sequence and "started simple" [53] [79]. However, in these experiments, the developmental sequence was crafted by hand: roboticists manually build simpler versions of a complex task and put the robot successively in versions of the task of increasing complexity. And when they wanted the robot to learn a new task, they had to design a novel reward function.

Thus, there is a need for mechanisms that allow the autonomous control and generation of the developmental trajectory. Psychologists have proposed that intrinsic motivations play a crucial role. Intrinsic motivations are mechanisms that push humans to explore activities or situations that have intermediate/optimal levels of novelty, cognitive dissonance, or challenge [59] [73] [75]. The role and structure of intrinsic motivation in humans have been made more precise thanks to recent discoveries in neuroscience showing the implication

of dopaminergic circuits and in exploration behaviours and curiosity [74] [93] [147]. Based on this, a number of researchers have began in the past few years to build computational implementation of intrinsic motivation [128] [129] [145] [57] [94] [112] [146]. While initial models were developed for simple simulated worlds, a current challenge is to manage to build intrinsic motivation systems that can efficiently drive exploratory behaviour in high-dimensional unprepared real world robotic sensorimotor spaces [129], [128], [130], [143]. Specific and complex problems are posed by real sensorimotor spaces, in particular due to the fact that they are both high-dimensional as well as (usually) deeply inhomogeneous. As an example for the latter issue, some regions of real sensorimotor spaces are often unlearnable due to inherent stochasticity or difficulty, in which case heuristics based on the incentive to explore zones of maximal unpredictability or uncertainty, which are often used in the field of active learning [70] [90] typically lead to catastrophic results. The issue of high dimensionality does not only concern motor spaces, but also sensory spaces, leading to the problem of correctly identifying, among typically thousands of quantities, those latent variables that have links to behavioral choices. In FLOWERS, we aim at developing intrinsically motivated exploration mechanisms that scale in those spaces, by studying suitable abstraction processes in conjunction with exploration strategies.

### 3.1.2. *Socially Guided and Interactive Learning*

Social guidance is as important as intrinsic motivation in the cognitive development of human babies [117]. There is a vast literature on learning by demonstration in robots where the actions of humans in the environment are recognized and transferred to robots [52]. Most such approaches are completely passive: the human executes actions and the robot learns from the acquired data. Recently, the notion of interactive learning has been introduced in [151], [60], motivated by the various mechanisms that allow humans to socially guide a robot [139]. In an interactive context the steps of self-exploration and social guidance are not separated and a robot learns by self exploration and by receiving extra feedback from the social context [151], [99], [113].

Social guidance is also particularly important for learning to segment and categorize the perceptual space. Indeed, parents interact a lot with infants, for example teaching them to recognize and name objects or characteristics of these objects. Their role is particularly important in directing the infant attention towards objects of interest that will make it possible to simplify at first the perceptual space by pointing out a segment of the environment that can be isolated, named and acted upon. These interactions will then be complemented by the children own experiments on the objects chosen according to intrinsic motivation in order to improve the knowledge of the object, its physical properties and the actions that could be performed with it.

In FLOWERS, we are aiming at including intrinsic motivation system in the self-exploration part thus combining efficient self-learning with social guidance [122], [123]. We also work on developing perceptual capabilities by gradually segmenting the perceptual space and identifying objects and their characteristics through interaction with the user [110] and robots experiments [95]. Another challenge is to allow for more flexible interaction protocols with the user in terms of what type of feedback is provided and how it is provided [107].

Exploration mechanisms are combined with research in the following directions:

### 3.1.3. *Cumulative learning, reinforcement learning and optimization of autonomous skill learning*

FLOWERS develops machine learning algorithms that can allow embodied machines to acquire cumulatively sensorimotor skills. In particular, we develop optimization and reinforcement learning systems which allow robots to discover and learn dictionaries of motor primitives, and then combine them to form higher-level sensorimotor skills.

### 3.1.4. *Autonomous perceptual and representation learning*

In order to harness the complexity of perceptual and motor spaces, as well as to pave the way to higher-level cognitive skills, developmental learning requires abstraction mechanisms that can infer structural information out of sets of sensorimotor channels whose semantics is unknown, discovering for example the topology of the body or the sensorimotor contingencies (proprioceptive, visual and acoustic). This process is meant to

be open- ended, progressing in continuous operation from initially simple representations towards abstract concepts and categories similar to those used by humans. Our work focuses on the study of various techniques for:

- autonomous multimodal dimensionality reduction and concept discovery;
- incremental discovery and learning of objects using vision and active exploration, as well as of auditory speech invariants;
- learning of dictionaries of motion primitives with combinatorial structures, in combination with linguistic description;
- active learning of visual descriptors useful for action (e.g. grasping);

### 3.1.5. Embodiment and maturational constraints

FLOWERS studies how adequate morphologies and materials (i.e. morphological computation), associated to relevant dynamical motor primitives, can importantly simplify the acquisition of apparently very complex skills such as full-body dynamic walking in biped. FLOWERS also studies maturational constraints, which are mechanisms that allow for the progressive and controlled release of new degrees of freedoms in the sensorimotor space of robots.

### 3.1.6. Discovering and abstracting the structure of sets of uninterpreted sensors and motors

FLOWERS studies mechanisms that allow a robot to infer structural information out of sets of sensorimotor channels whose semantics is unknown, for example the topology of the body and the sensorimotor contingencies (proprioceptive, visual and acoustic). This process is meant to be open-ended, progressing in continuous operation from initially simple representations to abstract concepts and categories similar to those used by humans.

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

# 3. Research Program
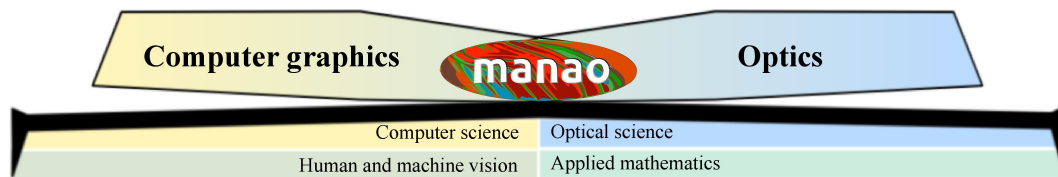
## 3.1. Related Scientific Domains



*Figure 3. Related scientific domains of the MANAO project.*

The *MANAO* project aims at studying, acquiring, modeling, and rendering the interactions between the three components that are light, shape, and matter from the viewpoint of an observer. As detailed more lengthily in the next section, such a work will be done using the following approach: first, we will tend to consider that these three components do not have strict frontiers when considering their impacts on the final observers; then, we will not only work in **computer graphics**, but also at the intersection of computer graphics and **optics**, exploring the mutual benefits that the two domains may provide. It is thus intrinsically a **transdisciplinary** project (as illustrated in Figure 3 ) and we expect results in both domains.

Thus, the proposed team-project aims at establishing a close collaboration between computer graphics (e.g., 3D modeling, geometry processing, shading techniques, vector graphics, and GPU programming) and optics (e.g., design of optical instruments, and theories of light propagation). The following examples illustrate the strengths of such a partnership. First, in addition to simpler radiative transfer equations  [36] commonly used in computer graphics, research in the later will be based on state-of-the-art understanding of light propagation and scattering in real environments. Furthermore, research will rely on appropriate instrumentation expertise for the measurement  [48], [49] and display  [47] of the different phenomena. Reciprocally, optics researches may benefit from the expertise of computer graphics scientists on efficient processing to investigate interactive simulation, visualization, and design. Furthermore, new systems may be developed by unifying optical and digital processing capabilities. Currently, the scientific background of most of the team members is related to computer graphics and computer vision. A large part of their work have been focused on simulating and analyzing optical phenomena as well as in acquiring and visualizing them. Combined with the close collaboration with the optics laboratory LP2N (<span style="color:red">http://www.lp2n.fr</span>) and with the students issued from the "Institut d'Optique" (<span style="color:red">http://www.institutoptique.fr</span>), this background ensures that we can expect the following results from the project: the construction of a common vocabulary for tightening the collaboration between the two scientific domains and creating new research topics. By creating this context, we expect to attract (and even train) more trans-disciplinary researchers.

At the boundaries of the *MANAO* project lie issues in **human and machine vision**. We have to deal with the former whenever a human observer is taken into account. On one side, computational models of human vision are likely to guide the design of our algorithms. On the other side, the study of interactions between light, shape, and matter may shed some light on the understanding of visual perception. The same kind of connections are expected with machine vision. On the one hand, traditional computational methods for acquisition (such as photogrammetry) are going to be part of our toolbox. On the other hand, new display technologies (such as the ones used for augmented reality) are likely to benefit from our integrated approach

and systems. In the *MANAO* project we are mostly users of results from human vision. When required, some experimentation might be done in collaboration with experts from this domain, like with the European PRISM project. For machine vision, provided the tight collaboration between optical and digital systems, research will be carried out inside the *MANAO* project.

Analysis and modeling rely on **tools from applied mathematics** such as differential and projective geometry, multi-scale models, frequency analysis [38] or differential analysis [70], linear and non-linear approximation techniques, stochastic and deterministic integrations, and linear algebra. We not only rely on classical tools, but also investigate and adapt recent techniques (e.g., improvements in approximation techniques), focusing on their ability to run on modern hardware: the development of our own tools (such as Eigen) is essential to control their performances and their abilities to be integrated into real-time solutions or into new instruments.

## 3.2. Research axes

The *MANAO* project is organized around four research axes that cover the large range of expertise of its members and associated members. We briefly introduce these four axes in this section. More details and their inter-influences that are illustrated in the Figure 2 will be given in the following sections.

Axis 1 is the theoretical foundation of the project. Its main goal is to increase the understanding of light, shape, and matter interactions by combining expertise from different domains: optics and human/machine vision for the analysis and computer graphics for the simulation aspect. The goal of our analyses is to identify the different layers/phenomena that compose the observed signal. In a second step, the development of physical simulations and numerical models of these identified phenomena is a way to validate the pertinence of the proposed decompositions.

In Axis 2, the final observers are mainly physical captors. Our goal is thus the development of new acquisition and display technologies that combine optical and digital processes in order to reach fast transfers between real and digital worlds, in order to increase the convergence of these two worlds.

Axes 3 and 4 focus on two aspects of computer graphics: rendering, visualization and illustration in Axis 3, and editing and modeling (content creation) in Axis 4. In these two axes, the final observers are mainly human users, either generic users or expert ones (e.g., archaeologist [74], computer graphics artists).

## 3.3. Axis 1: Analysis and Simulation

**Challenge:** Definition and understanding of phenomena resulting from interactions between light, shape, and matter as seen from an observer point of view.

**Results:** Theoretical tools and numerical models for analyzing and simulating the observed optical phenomena.

To reach the goals of the *MANAO* project, we need to **increase our understanding** of how light, shape, and matter act together in synergy and how the resulting signal is finally observed. For this purpose, we need to identify the different phenomena that may be captured by the targeted observers. This is the main objective of this research axis, and it is achieved by using three approaches: the simulation of interactions between light, shape, and matter, their analysis and the development of new numerical models. This resulting improved knowledge is a foundation for the researches done in the three other axes, and the simulation tools together with the numerical models serve the development of the joint optical/digital systems in Axis 2 and their validation.

One of the main and earliest goals in computer graphics is to faithfully reproduce the real world, focusing mainly on light transport. Compared to researchers in physics, researchers in computer graphics rely on a subset of physical laws (mostly radiative transfer and geometric optics), and their main concern is to efficiently use the limited available computational resources while developing as fast as possible algorithms. For this purpose, a large set of theoretical as well as computational tools has been introduced to take a **maximum benefit of hardware** specificities. These tools are often dedicated to specific phenomena (e.g., direct or indirect lighting, color bleeding, shadows, caustics). An efficiency-driven approach needs such a classification of light paths [44] in order to develop tailored strategies [86]. For instance, starting from simple direct lighting,

more complex phenomena have been progressively introduced: first diffuse indirect illumination [42], [78], then more generic inter-reflections [51], [36] and volumetric scattering [75], [33]. Thanks to this search for efficiency and this classification, researchers in computer graphics have developed a now recognized expertise in fast-simulation of light propagation. Based on finite elements (radiosity techniques) or on unbiased Monte Carlo integration schemes (ray-tracing, particle-tracing, ...), the resulting algorithms and their combination are now sufficiently accurate to be used-back in physical simulations. The *MANAO* project will continue the search for **efficient and accurate simulation** techniques, but extending it from computer graphics to optics. Thanks to the close collaboration with scientific researchers from optics, new phenomena beyond radiative transfer and geometric optics will be explored.

Search for algorithmic efficiency and accuracy has to be done in parallel with **numerical models**. The goal of visual fidelity (generalized to accuracy from an observer point of view in the project) combined with the goal of efficiency leads to the development of alternative representations. For instance, common classical finite-element techniques compute only basis coefficients for each discretization element: the required discretization density would be too large and to computationally expensive to obtain detailed spatial variations and thus visual fidelity. Examples includes texture for decorrelating surface details from surface geometry and high-order wavelets for a multi-scale representation of lighting [32]. The numerical complexity explodes when considering directional properties of light transport such as radiance intensity (Watt per square meter and per steradian - $W.m^{-2}.sr^{-1}$), reducing the possibility to simulate or accurately represent some optical phenomena. For instance, Haar wavelets have been extended to the spherical domain [77] but are difficult to extend to non-piecewise-constant data [80]. More recently, researches prefer the use of Spherical Radial Basis Functions [83] or Spherical Harmonics [69]. For more complex data, such as reflective properties (e.g., BRDF [63], [52] - 4D), ray-space (e.g., Light-Field [60] - 4D), spatially varying reflective properties (6D - [73]), new models, and representations are still investigated such as rational functions [66] or dedicated models [20] and parameterizations [76], [81]. For each (newly) defined phenomena, we thus explore the space of possible numerical representations to determine the **most suited one for a given application**, like we have done for BRDF [66].
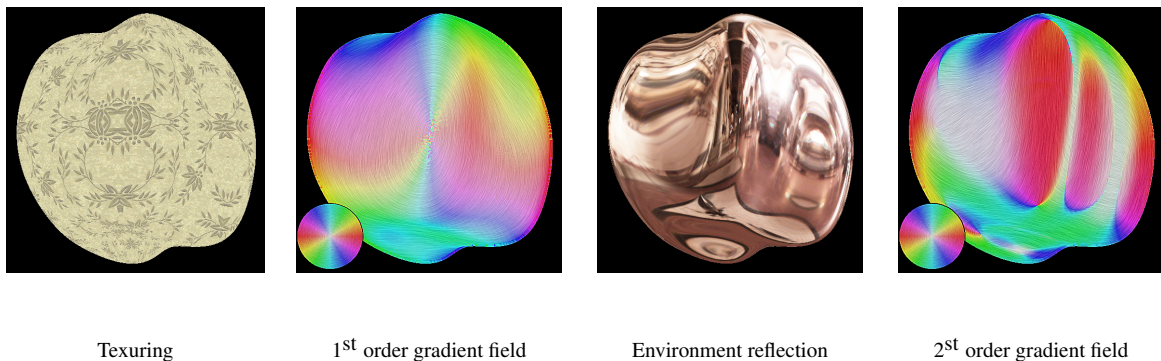


Texuring                1<sup>st</sup> order gradient field        Environment reflection        2<sup>st</sup> order gradient field

*Figure 4. First-oder analysis [87] have shown that shading variations are caused by depth variations (first-order gradient field) and by normal variations (second-order fields). These fields are visualized using hue and saturation to indicate direction and magnitude of the flow respectively.*

Before being able to simulate or to represent the different **observed phenomena**, we need to define and describe them. To understand the difference between an observed phenomenon and the classical light, shape, and matter decomposition, we can take the example of a highlight. Its observed shape (by a human user or a sensor) is the resulting process of the interaction of these three components, and can be simulated this way. However, this does not provide any intuitive understanding of their relative influence on the final shape: an artist will directly describe the resulting shape, and not each of the three properties. We thus want to decompose the observed signal into models for each scale that can be easily understandable, representable,

and manipulable. For this purpose, we will rely on the **analysis** of the resulting interaction of light, shape, and matter as observed by a human or a physical sensor. We first consider this analysis from an **optical point of view**, trying to identify the different phenomena and their scale according to their mathematical properties (e.g., differential [70] and frequency analysis [38]). Such an approach has leaded us to exhibit the influence of surfaces flows (depth and normal gradients) into lighting pattern deformation (see Figure 4 ). For a **human observer**, this correspond to one recent trend in computer graphics that takes into account the human visual systems [39] both to evaluate the results and to guide the simulations.

## 3.4. Axis 2: From Acquisition to Display

**Challenge:** Convergence of optical and digital systems to blend real and virtual worlds.

**Results:** Instruments to acquire real world, to display virtual world, and to make both of them interact.
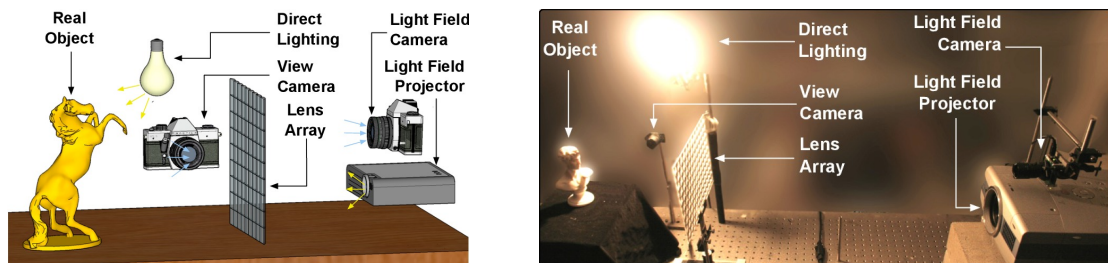


*Figure 5. Light-Field transfer: global illumination between real and synthetic objects [31]*

In this axis, we investigate *unified acquisition and display systems*, that is systems which combine optical instruments with digital processing. From digital to real, we investigate new display approaches [60], [47]. We consider projecting systems and surfaces [27], for personal use, virtual reality and augmented reality [22]. From the real world to the digital world, we favor direct measurements of parameters for models and representations, using (new) optical systems unless digitization is required [41], [40]. These resulting systems have to acquire the different phenomena described in Axis 1 and to display them, in an efficient manner [45], [21], [46], [49]. By efficient, we mean that we want to shorten the path between the real world and the virtual world by increasing the data bandwidth between the real (analog) and the virtual (digital) worlds, and by reducing the latency for real-time interactions (we have to prevent unnecessary conversions, and to reduce processing time). To reach this goal, the systems have to be designed as a whole, not by a simple concatenation of optical systems and digital processes, nor by considering each component independently [50].

To increase data bandwidth, one solution is to **parallelize more and more the physical systems**. One possible solution is to multiply the number of simultaneous acquisitions (e.g., simultaneous images from multiple viewpoints [49], [68]). Similarly, increasing the number of viewpoints is a way toward the creation of full 3D displays [60]. However, full acquisition or display of 3D real environments theoretically requires a continuous field of viewpoints, leading to huge data size. Despite the current belief that the increase of computational power will fill the missing gap, when it comes to visual or physical realism, if you double the processing power, people may want four times more accuracy, thus increasing data size as well. To reach the best performances, a trade-off has to be found between the amount of data required to represent accurately the reality and the amount of required processing. This trade-off may be achieved using **compressive sensing**. Compressive sensing is a new trend issued from the applied mathematics community that provides tools to accurately reconstruct a signal from a small set of measurements assuming that it is sparse in a transform domain (e.g., [67], [92]).

We prefer to achieve this goal by avoiding as much as possible the classical approach where acquisition is followed by a fitting step: this requires in general a large amount of measurements and the fitting itself may consume consequently too much memory and preprocessing time. By **preventing unnecessary conversion** through fitting techniques, such an approach increase the speed and reduce the data transfer for acquisition but also for display. One of the best recent examples is the work of Cossairt et al. [31]. The whole system is designed around a unique representation of the energy-field issued from (or leaving) a 3D object, either virtual or real: the Light-Field. A Light-Field encodes the light emitted in any direction from any position on an object. It is acquired thanks to a lens-array that leads to the capture of, and projection from, multiple simultaneous viewpoints. A unique representation is used for all the steps of this system. Lens-arrays, parallax barriers, and coded-aperture [57] are one of the key technologies to develop such acquisition (e.g., Light-Field camera [0] [50] and acquisition of light-sources [41]), projection systems (e.g., auto-stereoscopic displays). Such an approach is versatile and may be applied to improve classical optical instruments [55]. More generally, by designing unified optical and digital systems [64], it is possible to leverage the requirement of processing power, the memory footprint, and the cost of optical instruments.

Those are only some examples of what we investigate. We also consider the following approaches to develop new unified systems. First, similar to (and based on) the analysis goal of Axis 1, we have to take into account as much as possible the characteristics of the measurement setup. For instance, when fitting cannot be avoided, integrating them may improve both the processing efficiency and accuracy [66]. Second, we have to integrate signals from multiple sensors (such as GPS, accelerometer, ...) to prevent some computation (e.g., [58]). Finally, the experience of the group in surface modeling help the design of optical surfaces [53] for light sources or head-mounted displays.

## 3.5. Axis 3: Rendering, Visualization and Illustration

**Challenge:** How to offer the most legible signal to the final observer in real-time?

**Results:** High-level shading primitives, expressive rendering techniques for object depiction, real-time realistic rendering algorithms

| Realistic | Rendering | Visualization | and Illustration |
|---|---|---|---|



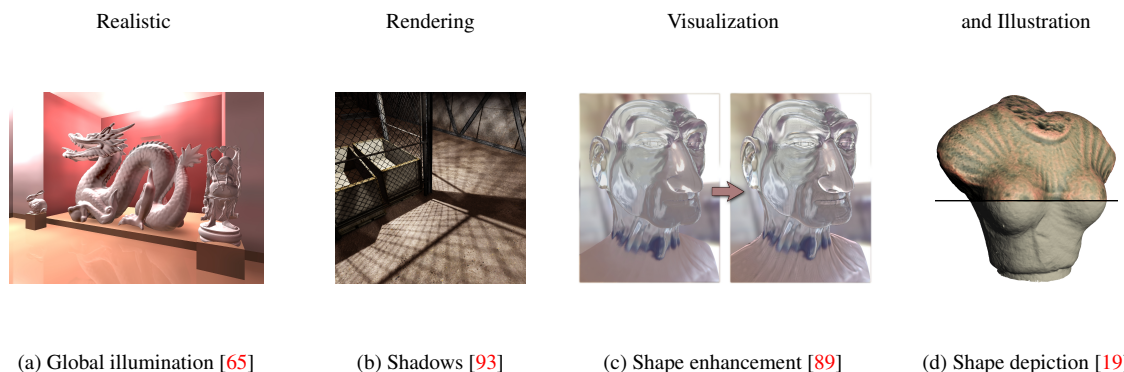(a) Global illumination [65]    (b) Shadows [93]    (c) Shape enhancement [89]    (d) Shape depiction [19]

*Figure 6. In the MANAO project, we are investigating rendering techniques from realistic solutions (e.g., inter-reflections (a) and shadows (b)) to more expressive ones (shape enhancement (c) with realistic style and shape depiction (d) with stylized style) for visualization.*

The main goal of this axis is to offer to the final observer, in this case mostly a human user, the most legible signal in real-time. Thanks to the analysis and to the decomposition in different phenomena resulting from interactions between light, shape, and matter (Axis 1), and their perception, we can use them to convey essential information in the most pertinent way. Here, the word *pertinent* can take various forms depending on the application.

---

[0]Lytro, http://www.lytro.com/

In the context of scientific illustration and visualization, we are primarily interested in tools to convey shape or material characteristics of objects in animated 3D scenes. **Expressive rendering** techniques (see Figure 6 c,d) provide means for users to depict such features with their own style. To introduce our approach, we detail it from a shape-depiction point of view, domain where we have acquired a recognized expertise. Prior work in this area mostly focused on stylization primitives to achieve line-based rendering [90], [54] or stylized shading [25], [89] with various levels of abstraction. A clear representation of important 3D **object features** remains a major challenge for better shape depiction, stylization and abstraction purposes. Most existing representations provide only local properties (e.g., curvature), and thus lack characterization of broader shape features. To overcome this limitation, we are developing higher level descriptions of shape [18] with increased robustness to sparsity, noise, and outliers. This is achieved in close collaboration with Axis 1 by the use of higher-order local fitting methods, multi-scale analysis, and global regularization techniques. In order not to neglect the observer and the material characteristics of the objects, we couple this approach with an analysis of the appearance model. To our knowledge, this is an approach which has not been considered yet. This research direction is at the heart of the *MANAO* project, and has a strong connection with the analysis we plan to conduct in Axis 1. Material characteristics are always considered at the light ray level, but an understanding of **higher-level primitives** (like the shape of highlights and their motion) would help us to produce more legible renderings and permit novel stylizations; for instance, there is no method that is today able to create stylized renderings that follow the motion of highlights or shadows. We also believe such tools also play a fundamental role for geometry processing purposes (such as shape matching, reassembly, simplification), as well as for editing purposes as discussed in Axis 4.

In the context of **real-time photo-realistic rendering** ((see Figure 6 a,b), the challenge is to compute the most plausible images with minimal effort. During the last decade, a lot of work has been devoted to design approximate but real-time rendering algorithms of complex lighting phenomena such as soft-shadows [91], motion blur [38], depth of field [79], reflexions, refractions, and inter-reflexions. For most of these effects it becomes harder to discover fundamentally new and faster methods. On the other hand, we believe that significant speedup can still be achieved through more clever use of **massively parallel architectures** of the current and upcoming hardware, and/or through more clever tuning of the current algorithms. In particular, regarding the second aspect, we remark that most of the proposed algorithms depend on several parameters which can be used to **trade the speed over the quality**. Significant speed-up could thus be achieved by identifying effects that would be masked or facilitated and thus devote appropriate computational resources to the rendering [56], [37]. Indeed, the algorithm parameters controlling the quality vs speed are numerous without a direct mapping between their values and their effect. Moreover, their ideal values vary over space and time, and to be effective such an auto-tuning mechanism has to be extremely fast such that its cost is largely compensated by its gain. We believe that our various work on the analysis of the appearance such as in Axis 1 could be beneficial for such purpose too.

Realistic and real-time rendering is closely related to Axis 2: real-time rendering is a requirement to close the loop between real world and digital world. We have to thus develop algorithms and rendering primitives that allow the integration of the acquired data into real-time techniques. We have also to take care of that these real-time techniques have to work with new display systems. For instance, stereo, and more generally multi-view displays are based on the multiplication of simultaneous images. Brute force solutions consist in independent rendering pipeline for each viewpoint. A more energy-efficient solution would take advantages of the computation parts that may be factorized. Another example is the rendering techniques based on image processing, such as our work on augmented reality [29]. Independent image processing for each viewpoint may disturb the feeling of depth by introducing inconsistent information in each images. Finally, more dedicated displays [47] would require new rendering pipelines.

## 3.6. Axis 4: Editing and Modeling

**Challenge:** Editing and modeling appearance using drawing- or sculpting-like tools through high level representations.

**Results:** High-level primitives and hybrid representations for appearance and shape.

During the last decade, the domain of computer graphics has exhibited tremendous improvements in image quality, both for 2D applications and 3D engines. This is mainly due to the availability of an ever increasing amount of shape details, and sophisticated appearance effects including complex lighting environments. Unfortunately, with such a growth in visual richness, even so-called *vectorial* representations (e.g., subdivision surfaces, Bézier curves, gradient meshes, etc.) become very dense and unmanageable for the end user who has to deal with a huge mass of control points, color labels, and other parameters. This is becoming a major challenge, with a necessity for novel representations. This Axis is thus complementary of Axis 3: the focus is the development of primitives that are easy to use for modeling and editing.

More specifically, we plan to investigate *vectorial representations* that would be amenable to the production of rich shapes with a minimal set of primitives and/or parameters. To this end we plan to build upon our insights on dynamic local reconstruction techniques and implicit surfaces [30] [24]. When working in 3D, an interesting approach to produce detailed shapes is by means of procedural geometry generation. For instance, many natural phenomena like waves or clouds may be modeled using a combination of procedural functions. Turning such functions into triangle meshes (main rendering primitives of GPUs) is a tedious process that appears not to be necessary with an adapted vectorial shape representation where one could directly turn procedural functions into implicit geometric primitives. Since we want to prevent unnecessary conversions in the whole pipeline (here, between modeling and rendering steps), we will also consider *hybrid representations* mixing meshes and implicit representations. Such research has thus to be conducted while considering the associated editing tools as well as performance issues. It is indeed important to keep *real-time performance* (cf. Axis 2) throughout the interaction loop, from user inputs to display, via editing and rendering operations. Finally, it would be interesting to add *semantic information* into 2D or 3D geometric representations. Semantic geometry appears to be particularly useful for many applications such as the design of more efficient manipulation and animation tools, for automatic simplification and abstraction, or even for automatic indexing and searching. This constitutes a complementary but longer term research direction.

In the *MANAO* project, we want to investigate representations beyond the classical light, shape, and matter decomposition. We thus want to directly control the appearance of objects both in 2D and 3D applications (e.g., [84]): this is a core topic of computer graphics. When working with 2D vector graphics, digital artists must carefully set up color gradients and textures: examples range from the creation of 2D logos to the photo-realistic imitation of object materials. Classic vector primitives quickly become impractical for creating illusions of complex materials and illuminations, and as a result an increasing amount of time and skill is required. This is only for still images. For animations, vector graphics are only used to create legible appearances composed of simple lines and color gradients. There is thus a need for more complex primitives that are able to accommodate complex reflection or texture patterns, while keeping the ease of use of vector graphics. For instance, instead of drawing color gradients directly, it is more advantageous to draw flow lines that represent local surface concavities and convexities. Going through such an intermediate structure then allows to deform simple material gradients and textures in a coherent way (see Figure 7 ), and animate them all at once. The manipulation of 3D object materials also raises important issues. Most existing material models are tailored to faithfully reproduce physical behaviors, not to be *easily controllable* by artists. Therefore artists learn to tweak model parameters to satisfy the needs of a particular shading appearance, which can quickly become cumbersome as the complexity of a 3D scene increases. We believe that an alternative approach is required, whereby material appearance of an object in a typical lighting environment is directly input (e.g., painted or drawn), and adapted to match a plausible material behavior. This way, artists will be able to create their own appearance (e.g., by using our shading primitives [84]), and replicate it to novel illumination environments and 3D models. For this purpose, we will rely on the decompositions and tools issued from Axis 1.
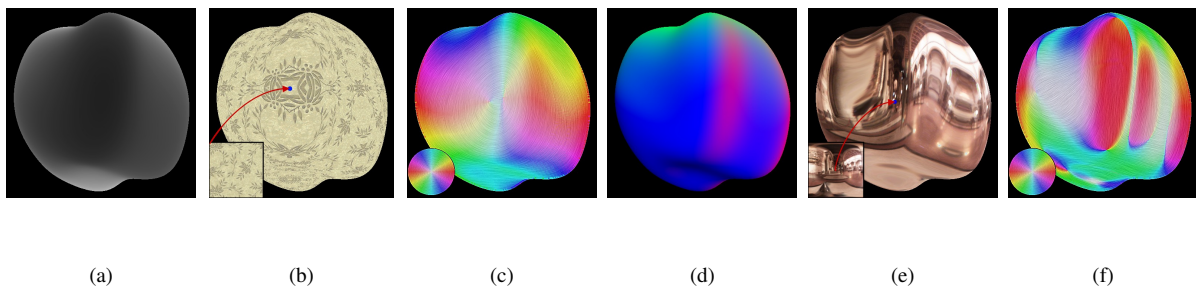
(a)  (b)  (c)  (d)  (e)  (f)

*Figure 7. Based on our analysis [87] (Axis 1), we have designed a system that mimics texture (left) and shading (right) effects using image processing alone. It takes depth (a) and normal (d) images as input, and uses them to deform images (b-e) in ways that closely approximate surface flows (c-f). It provides a convincing, yet artistically controllable illusion of 3D shape conveyed through texture or shading cues.*

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

# 3. Research Program

## 3.1. Research Program

To achieve our overall objective, we follow two main research axes, plus one transverse axis, as illustrated in Figure 2 .
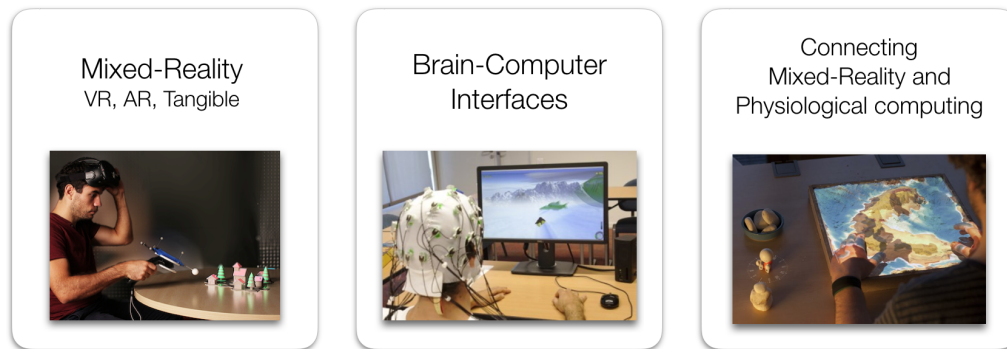


*Figure 2. Main research axes of Potioc.*

In the first axis dedicated to **Interaction in Mixed-Reality spaces**, we explore interaction paradigms that encompass virtual and/or physical objects. We are notably interested in hybrid environments that co-locate virtual and physical spaces, and we also explore approaches that allow one to move from one space to the other.

The second axis is dedicated to **Brain-Computer Interfaces (BCI)**, i.e., systems enabling user to interact by means of brain activity only. We target BCI systems that are reliable and accessible to a large number of people. To do so, we work on brain signal processing algorithms as well as on understanding and improving the way we train our users to control these BCIs.

Finally, in the **transverse** axis, we explore new approaches that involve both mixed-reality and neuro-physiological signals. In particular, tangible and augmented objects allow us to explore interactive physical visualizations of human inner states. Physiological signals also enable us to better assess user interaction, and consequently, to refine the proposed interaction techniques and metaphors.

From a methodological point of view, for these three axes, we work at three different interconnected levels. The first level is centered on the human sensori-motor and cognitive abilities, as well as user strategies and preferences, for completing interaction tasks. We target, in a fundamental way, a better understanding of humans interacting with interactive systems. The second level is about the creation of interactive systems. This notably includes development of hardware and software components that will allow us to explore new input and output modalities, and to propose adapted interaction techniques. Finally, in a last higher level, we are interested in specific application domains. We want to contribute to the emergence of new applications and usages, with a societal impact.