

*Inria*

RESEARCH CENTER

FIELD

**Networks, Systems and Services,  
Distributed Computing**

Activity Report 2019

**Section Scientific Foundations**

Edition: 2020-03-21



DISTRIBUTED SYSTEMS AND MIDDLEWARE

1. COAST Project-Team ..... 5  
2. CTRL-A Project-Team ..... 7  
3. DELYS Project-Team ..... 8  
4. MIMOVE Project-Team ..... 10  
5. Myriads Project-Team ..... 12  
6. SPIRALS Project-Team ..... 17  
7. STACK Project-Team ..... 22  
8. WHISPER Project-Team ..... 30  
9. WIDE Project-Team ..... 34

DISTRIBUTED AND HIGH PERFORMANCE COMPUTING

10. ALPINES Project-Team ..... 41  
11. AVALON Project-Team ..... 43  
12. DATAMOVE Project-Team ..... 46  
13. HIEPACS Project-Team ..... 49  
14. KERDATA Project-Team ..... 58  
15. POLARIS Project-Team ..... 62  
16. ROMA Project-Team ..... 66  
17. STORM Project-Team ..... 71  
18. TADAAM Project-Team ..... 74

DISTRIBUTED PROGRAMMING AND SOFTWARE ENGINEERING

19. DIVERSE Project-Team ..... 76  
20. EASE Project-Team ..... 86  
21. FOCUS Project-Team ..... 88  
22. INDES Project-Team ..... 89  
23. RMOD Project-Team ..... 90

NETWORKS AND TELECOMMUNICATIONS

24. AGORA Project-Team ..... 94  
25. COATI Project-Team ..... 98  
26. DANTE Project-Team ..... 99  
27. DIANA Project-Team ..... 102  
28. DIONYSOS Project-Team ..... 104  
29. DYOGENE Project-Team ..... 106  
30. EVA Project-Team ..... 107  
31. FUN Project-Team ..... 110  
32. GANG Project-Team ..... 114  
33. GANG Project-Team ..... 117  
34. MARACAS Team ..... 120  
35. NEO Project-Team ..... 124  
36. RESIST Team ..... 125  
37. SOCRATE Project-Team ..... 128

38. TRIBE Project-Team ..... 131

## COAST Project-Team

### 3. Research Program

#### 3.1. Introduction

Our scientific foundations are grounded on distributed collaborative systems supported by sophisticated data sharing mechanisms and on service oriented computing with an emphasis on orchestration and on non-functional properties. Distributed collaborative systems enable distributed group work supported by computer technologies. Designing such systems requires an expertise in Distributed Systems and in Computer-supported collaborative Work research area. Besides theoretical and technical aspects of distributed systems, the design of distributed collaborative systems must take into account the human factor to offer solutions suitable for users and groups. The Coast team vision is to move away from a centralized authority based collaboration toward a decentralized collaboration. Users will have full control over their data. They can store them locally and decide with whom to share them. The Coast team investigates the issues related to the management of distributed shared data and coordination between users and groups. Service oriented Computing [16] is an established domain on which the ECOO, Score and now the Coast teams have been contributing for a long time. It refers to the general discipline that studies the development of computer applications on the web. A service is an independent software program with a specific functional context and capabilities published as a service contract (or more traditionally an API). A service composition aggregates a set of services and coordinates their interactions. The scale, the autonomy of services, the heterogeneity and some design principles underlying Service Oriented Computing open new research questions that are at the basis of our research. They span the disciplines of **distributed computing**, **software engineering** and **computer supported collaborative work** (CSCW). Our approach to contribute to the general vision of Service Oriented Computing is to focus on the issue of the efficient and flexible construction of reliable and secure high-level services. We aim to achieve it through the coordination/orchestration/composition of other services provided by distributed organizations or people.

#### 3.2. Consistency Models for Distributed Collaborative Systems

Collaborative systems are distributed systems that allow users to share data. One important issue is to manage consistency of shared data according to concurrent access. Traditional consistency criteria such as serializability, linearizability are not adequate for collaborative systems. Causality, Convergence and Intention preservation (CCI) [21] are more suitable for developing middleware for collaborative applications. We develop algorithms for ensuring CCI properties on collaborative distributed systems. Constraints on the algorithms are different according to the kind of distributed system and to the data structure. The distributed system can be centralized, decentralized or peer-to-peer. The type of data can include strings, growable arrays, ordered trees, semantic graphs and multimedia data.

#### 3.3. Optimistic Replication

Replication of data among different nodes of a network promotes reliability, fault tolerance, and availability. When data are mutable, consistency among the different replicas must be ensured. Pessimistic replication is based on the principle of single-copy consistency while optimistic replication allows the replicas to diverge during a short time period. The consistency model for optimistic replication [19] is called eventual consistency, meaning that replicas are guaranteed to converge to the same value when the system is idle. Our research focuses on the two most promising families of optimistic replication algorithms for ensuring CCI:

- operational transformation (OT) algorithms [14]
- algorithms based on commutative replicated data types (CRDT) [18].

Operational transformation algorithms are based on the application of a transformation function when a remote modification is integrated into the local document. Integration algorithms are generic, being parametrised by operational transformation functions which depend on replicated document types. The advantage of these algorithms is their genericity. These algorithms can be applied to any data type and they can merge heterogeneous data in a uniform manner. Commutative replicated data types is a new class of algorithms initiated by WooT [15], the first algorithm designed WithOut Operational Transformations. They ensure consistency of highly dynamic content on peer-to-peer networks. Unlike traditional optimistic replication algorithms, they can ensure consistency without concurrency control. CRDT algorithms rely on natively commutative operations defined on abstract data types such as lists or ordered trees. Thus, they do not require a merge algorithm or an integration procedure.

### **3.4. Process Orchestration and Management**

Process Orchestration and Management is considered as a core discipline behind Service Management and Computing. It includes the analysis, the modelling, the execution, the monitoring and the continuous improvement of enterprise processes and is for us a central domain of studies. Many efforts have been devoted establishing standard business process models founded on well-grounded theories (e.g. Petri Nets) that meet the needs of business analysts, software engineers and software integrator. This led to heated debate in the Business Process Management (BPM) community as the two points of view are very difficult to reconcile. On one side, business people in general require models that are easy to use and understand and that can be quickly adapted to exceptional situations. On the other side, IT people need models with an operational semantic in order to be able transform them into executable artifacts. Part of our work has been an attempt to reconcile these points of view. This resulted in the development of the Bonita BPM system. It resulted also more recently on our work in crisis management where the same people are designing, executing and monitoring the process as it executes. More generally, and at a larger scale, we have been considering the problem of processes spanning the barriers of organizations. This leads to the more general problem of service composition as a way to coordinate inter organizational construction of applications. These applications provide value, based on the composition of lower level services [12].

### **3.5. Service Composition**

Recently, we started a study on service composition for software architects where services are coming from different providers with different plans (capacity, degree of resilience...). The objective is to support the architects to select the most accurate services (wrt. to their requirements, both functional and non-functional) and plans for building their software. We also compute the properties that we enforce for the composition of these services.

## CTRL-A Project-Team

### 3. Research Program

#### 3.1. Modeling and control techniques for autonomic computing

The main objective of CTRL-A translates into a number of scientific challenges, the most important of these are:

- (i) programming language support, on the two facets of model-oriented languages, based on automata [5], and of domain specific languages, following e.g., a component-based approach [4], [1] or related to rule-based or HMI languages ;
- (ii) design methods for reconfiguration controller design in computing systems, proposing generic systems architectures and models based on transition systems [3], classical continuous control or controlled stochastic systems.

We adopt a strategy of constant experimental identification of needs and validation of proposals, in application domains like middleware platforms for Cloud systems [3], multi-core HPC architectures [10], Dynamic Partial Reconfiguration in FPGA-based hardware [2] and the IoT and smart environments [8].

Achieving the goals of CTRL-A requires multidisciplinary and expertise from several domains. The expertise in Autonomic Computing and programming languages is covered internally by members of the Ctrl-A team. On the side of theoretical aspects of control, we have active external collaborations with researchers specialized in Control Theory, in the domain of Discrete Event Systems as well as in classical, continuous control. Additionally, an important requirement for our research to have impact is to have access to concrete, real-world computing systems requiring reconfiguration control. We target autonomic computing at different scales, in embedded systems or in cloud infrastructures, which are traditionally different domains. This is addressed by external collaborations, with experts in either hardware or software platforms, who are generally missing our competences on model-based control of reconfigurations.

## DELYS Project-Team

### 3. Research Program

#### 3.1. Research rationale

DELYS addresses both theoretical and practical issues of *Computer Systems*, leveraging our dual expertise in theoretical and experimental research. Our approach is a “virtuous cycle,” triggered by issues with real systems, of algorithm design which we prove correct and evaluate theoretically, and then implement and test experimentally feeding back to theory. The major challenges addressed by DELYS are the sharing of information and guaranteeing correct execution of highly-dynamic computer systems. Our research covers a large spectrum of distributed computer systems: multicore computers, mobile networks, cloud computing systems, and dynamic communicating entities. This holistic approach enables handling related problems at different levels. Among such problems we can highlight consensus, fault detection, scalability, search of information, resource allocation, replication and consistency of shared data, dynamic content distribution, and concurrent and parallel algorithms.

Two main evolutions in the Computer Systems area strongly influence our research project:

(1) Modern computer systems are **increasingly distributed, dynamic** and composed of multiple devices **geographically spread over heterogeneous platforms**, spanning multiple management domains. Years of research in the field are now coming to fruition, and are being used by millions of users of web systems, peer-to-peer systems, gaming and social applications, cloud computing, and now fog computing. These new uses bring new challenges, such as *adaptation to dynamically-changing conditions*, where knowledge of the system state can only be partial and incomplete.

(2) **Heterogeneous architectures and virtualisation are everywhere**. The parallelism offered by distributed clusters and *multicore* architectures is opening highly parallel computing to new application areas. To be successful, however, many issues need to be addressed. Challenges include obtaining a consistent view of shared resources, such as memory, and optimally distributing computations among heterogeneous architectures. These issues arise at a more fine-grained level than before, leading to the need for different solutions down to OS level itself.

The scientific challenges of the distributed computing systems are subject to many important features which include scalability, fault tolerance, dynamics, emergent behaviour, heterogeneity, and virtualisation at many levels. Algorithms designed for traditional distributed systems, such as resource allocation, data storage and placement, and concurrent access to shared data, need to be redefined or revisited in order to work properly under the constraints of these new environments. Sometimes, classical “*static*” problems, (*e.g.*, Election Leader, Spanning Tree Construction, ...) even need to be redefined to consider the unstable nature of the distributed system. In particular, DELYS will focus on a number of key challenges:

Consistency in geo-scale systems. Distributed systems need to scale to large geographies and large numbers of attached devices, while executing in an untamed, unstable environment. This poses difficult scientific challenges, which are all the more pressing as the cloud moves more and more towards the edge, IoT and mobile computing. A key issue is how to share data effectively and consistently across the whole spectrum. DELYS has made several key contributions, including CRDTs, the Transactional Causal Consistency Plus model, the AntidoteDB geo-distributed database, and its edge extension EdgeAnt.

Rethinking distributed algorithms. From a theoretical point of view the key question is how to adapt the fundamental building blocks to new architectures. More specifically, how to rethink the classical algorithms to take into account the dynamics of advanced modern systems. Since a recent past, there have been several papers that propose models for dynamic systems: there is practically a different model for each setting and currently there is no unification of models. Furthermore, models



often suffer of lack of realism. One of the key challenge is to identify which assumptions make sense in new distributed systems. DELYS's objectives are then (1) to identify under which realistic assumptions a given fundamental problem such as mutual exclusion, consensus or leader election can be solved and (2) to design efficient algorithms under these assumptions.

**Resource management in heterogeneous systems.** The key question is how to manage resources on large and heterogeneous configurations. Managing resources in such systems requires fully decentralized solutions, and to rethink the way various platforms can collaborate and interoperate with each other. In this context, data management is a key component. The fundamental issue we address is how to efficiently and reliably share information in highly distributed environments.

**Adaptation of runtimes.** One of the main challenge of the OS community is how to adapt runtime supports to new architectures. With the increasingly widespread use of multicore architectures and virtualised environments, internal runtime protocols need to be revisited. Especially, memory management is crucial in OS and virtualisation technologies have highly impact on it. On one hand, the isolation property of virtualisation has severe side effects on the efficiency of memory allocation since it needs to be constantly balanced between hosted OSs. On the other hand, by hiding the physical machine to OSs, virtualisation prevents them to efficiently place their data in memory on different cores. Our research will thus focus on providing solutions to efficiently share memory between OSs without jeopardizing isolation properties.

## MIMOVE Project-Team

### 3. Research Program

#### 3.1. Introduction

The research questions identified above call for radically new ways in conceiving, developing and running mobile distributed systems. In response to this challenge, MiMove's research aims at enabling next-generation mobile distributed systems that are the focus of the following research topics.

#### 3.2. Emergent mobile distributed systems

Uncertainty in the execution environment calls for designing mobile distributed systems that are able to run in a beforehand unknown, ever-changing context. Nevertheless, the complexity of such change cannot be tackled at system design-time. Emergent mobile distributed systems are systems which, due to their automated, dynamic, environment-dependent composition and execution, *emerge* in a possibly non-anticipated way and manifest *emergent properties*, i.e., both systems and their properties take their complete form only at runtime and may evolve afterwards. This contrasts with the typical software engineering process, where a system is finalized during its design phase. MiMove's research focuses on enabling the emergence of mobile distributed systems while assuring that their required properties are met. This objective builds upon pioneering research effort in the area of *emergent middleware* initiated by members of the team and collaborators [1], [3].

#### 3.3. Large-scale mobile sensing and actuation

The extremely large scale and dynamicity expected in future mobile sensing and actuation systems lead to the clear need for algorithms and protocols for addressing the resulting challenges. More specifically, since connected devices will have the capability to sense physical phenomena, perform computations to arrive at decisions based on the sensed data, and drive actuation to change the environment, enabling proper coordination among them will be key to unlocking their true potential. Although similar challenges have been addressed in the domain of networked sensing, including by members of the team [7], the specific challenges arising from the *extremely large scale* of mobile devices – a great number of which will be attached to people, with uncontrolled mobility behavior – are expected to require a significant rethink in this domain. MiMove's research investigates techniques for efficient coordination of future mobile sensing and actuation systems with a special focus on their dependability.

#### 3.4. Mobile social crowd-sensing

While mobile social sensing opens up the ability of sensing phenomena that may be costly or impossible to sense using embedded sensors (e.g., subjective crowdedness causing discomfort or joyfulness, as in a bus or in a concert) and leading to a feeling of being more socially involved for the citizens, there are unique consequent challenges. Specifically, MiMove's research focuses on the problems involved in the combination of the physically sensed data, which are quantitative and objective, with the mostly qualitative and subjective data arising from social sensing. Enabling the latter calls for introducing mechanisms for incentivising user participation and ensuring the privacy of user data, as well as running empirical studies for understanding the complex social behaviors involved. These objectives build upon previous research work by members of the team on mobile social ecosystems and privacy, as well as a number of efforts and collaborations in the domain of smart cities and transport that have resulted in novel mobile applications enabling empirical studies of social sensing systems.

### 3.5. Active and passive probing methods

We are developing methods that actively introduce probes in the network to discover properties of the connected devices and network segments. We are focusing in particular on methods to discover properties of home networks (connected devices and their types) and to distinguish if performance bottlenecks lie within the home network versus in the different network segments outside (e.g., Internet access provider, interconnects, or content provider). Our goal is to develop adaptive methods that can leverage the collaboration of the set of available devices (including end-user devices and the home router, depending on which devices are running the measurement software).

We are also developing passive methods that simply observe network traffic to infer the performance of networked applications and the location of performance bottlenecks, as well as to extract patterns of web content consumption. We are working on techniques to collect network traffic both at user's end-devices and at home routers. We also have access to network traffic traces collected on a campus network and on a large European broadband access provider.

### 3.6. Inferring user online experience

We are developing hybrid measurement methods that combine passive network measurement techniques to infer application performance with techniques from HCI to measure user perception as well as methods to directly measure application quality. We later use the resulting datasets to build models of user perception of network performance based only on data that we can obtain automatically from the user device or from user's traffic observed in the network.

### 3.7. Real time data analytics

The challenge of deriving insights from the Internet of Things (IoT) has been recognized as one of the most exciting and key opportunities for both academia and industry. The time value of data is crucial for many IoT-based systems requiring *real-time* (or near real-time) *control* and *automation*. Such systems typically collect data continuously produced by "things" (i.e., devices), and analyze them in (sub-) seconds in order to act promptly, e.g., for detecting security breaches of digital systems, for spotting malfunctions of physical assets, for recommending goods and services based on the proximity of potential clients, etc. Hence, they require to both *ingest* and *analyze in real-time* data arriving with different velocity from various IoT data streams.

Existing incremental (online or streaming) techniques for descriptive statistics (e.g., frequency distributions, frequent patterns, etc.) or predictive statistics (e.g., classification, regression) usually assume a good enough quality dataset for mining patterns or training models. However, IoT raw data produced in the wild by sensors embedded in the environment or wearable by users are prone to errors and noise. Effective and efficient algorithms are needed for *detecting* and *repairing data impurities* (for controlling data quality) as well as *understanding data dynamics* (for defining alerts) in real-time, for collections of IoT data streams that might be geographically distributed. Moreover, supervised deep learning and data analytics techniques are challenged by the presence of sparse ground truth data in real IoT applications. Lightweight and adaptive semi-supervised or unsupervised techniques are needed to power real-time anomaly and novelty detection in IoT data streams. The effectiveness of these techniques should be able to reach a useful level through training on a relatively small amount of (preferably unlabeled) data while they can cope distributional characteristics of data evolving over time.

## Myriads Project-Team

### 3. Research Program

#### 3.1. Introduction

In this section, we present our research challenges along four work directions: resource and application management in distributed cloud and fog computing architectures for scaling clouds in Section 3.2, energy management strategies for greening clouds in Section 3.3, security and data protection aspects for securing cloud-based information systems and applications in Section 3.4, and methods for experimenting with clouds in Section 3.5.

#### 3.2. Scaling fogs and clouds

##### 3.2.1. Resource management in hierarchical clouds

The next generation of utility computing appears to be an evolution from highly centralized clouds towards more decentralized platforms. Today, cloud computing platforms mostly rely on large data centers servicing a multitude of clients from the edge of the Internet. Servicing cloud clients in this manner suggests that locality patterns are ignored: wherever the client issues his/her request from, the request will have to go through the backbone of the Internet provider to the other side of the network where the data center relies. Besides this extra network traffic and this latency overhead that could be avoided, other common centralization drawbacks in this context stand in limitations in terms of security/legal issues and resilience.

At the same time, it appears that network backbones are over-provisioned for most of their usage. This advocates for placing computing resources directly within the backbone network. The general challenge of resource management for such clouds stands in trying to be locality-aware: for the needs of an application, several virtual machines may exchange data. Placing them *close* to each others can significantly improve the performance of the application they compose. More generally, building an overlay network which takes the hierarchical aspects of the platform without being a hierarchical overlay – which comes with load balancing and resilience issues is a challenge by itself.

We expect to integrate the results of these works in the Discovery initiative [33] which aims at revisiting OpenStack to offer a cloud stack able to manage utility computing platforms where computing resources are located in small computing centers in the backbone's PoPs (Point of Presence) and interconnected through the backbone's internal links.

##### 3.2.2. Resource management in fog computing architectures

Fog computing infrastructures are composed of compute, storage and networking resources located at the edge of wide-area networks, in immediate proximity to the end users. Instead of treating the mobile operator's network as a high-latency dumb pipe between the end users and the external service providers, fog platforms aim at deploying cloud functionalities *within* the mobile phone network, inside or close to the mobile access points. Doing so is expected to deliver added value to the content providers and the end users by enabling new types of applications ranging from Internet-of-Things applications to extremely interactive systems (e.g., augmented reality). Simultaneously, it will generate extra revenue streams for the mobile network operators, by allowing them to position themselves as cloud computing operators and to rent their already-deployed infrastructure to content and application providers.

Fog computing platforms have very different geographical distribution compared to traditional clouds. While traditional clouds are composed of many reliable and powerful machines located in a very small number of data centers and interconnected by very high-speed networks, mobile edge cloud are composed of a very large number of points-of-presence with a couple of weak and potentially unreliable servers, interconnected with each other by commodity long-distance networks. This creates new demands for the organization of a scalable mobile edge computing infrastructure, and opens new directions for research.

The main challenges that we plan to address are:

- How should an edge cloud infrastructure be designed such that it remains scalable, fault-tolerant, controllable, energy-efficient, etc.?
- How should applications making use of edge clouds be organized? One promising direction is to explore the extent to which stream-data processing platforms such as Apache Spark and Apache Flink can be adapted to become one of the main application programming paradigms in such environments.

### 3.2.3. *Self-optimizing applications in multi-cloud environments*

As the use of cloud computing becomes pervasive, the ability to deploy an application on a multi-cloud infrastructure becomes increasingly important. Potential benefits include avoiding dependence on a single vendor, taking advantage of lower resource prices or resource proximity, and enhancing application availability. Supporting multi-cloud application management involves two tasks. First, it involves selecting an initial multi-cloud application deployment that best satisfies application objectives and optimizes performance and cost. Second, it involves dynamically adapting the application deployment in order to react to changes in execution conditions, application objectives, cloud provider offerings, or resource prices. Handling price changes in particular is becoming increasingly complex. The reason is the growing trend of providers offering sophisticated, dynamic pricing models that allow buying and selling resources of finer granularities for shorter time durations with varying prices.

Although multi-cloud platforms are starting to emerge, these platforms impose a considerable amount of effort on developers and operations engineers, provide no support for dynamic pricing, and lack the responsiveness and scalability necessary for handling highly-distributed, dynamic applications with strict quality requirements. The goal of this work is to develop techniques and mechanisms for automating application management, enabling applications to cope with and take advantage of the dynamic, diverse, multi-cloud environment in which they operate.

The main challenges arising in this context are:

- selecting effective decision-making approaches for application adaptation,
- supporting scalable monitoring and adaptation across multiple clouds,
- performing adaptation actions in a cost-efficient and safe manner.

## 3.3. Greening clouds

The ICT (Information and Communications Technologies) ecosystem now approaches 5% of world electricity consumption and this ICT energy use will continue to grow fast because of the information appetite of Big Data, large networks and large infrastructures as Clouds that unavoidably leads to large power.

### 3.3.1. *Smart grids and clouds*

We propose exploiting Smart Grid technologies to come to the rescue of energy-hungry Clouds. Unlike in traditional electrical distribution networks, where power can only be moved and scheduled in very limited ways, Smart Grids dynamically and effectively adapt supply to demand and limit electricity losses (currently 10% of produced energy is lost during transmission and distribution).

For instance, when a user submits a Cloud request (such as a Google search for instance), it is routed to a data center that processes it, computes the answer and sends it back to the user. Google owns several data centers spread across the world and for performance reasons, the center answering the user's request is more likely to be the one closest to the user. However, this data center may be less energy efficient. This request may have consumed less energy, or a different kind of energy (renewable or not), if it had been sent to this further data center. In this case, the response time would have been increased but maybe not noticeably: a different trade-off between quality of service (QoS) and energy-efficiency could have been adopted.

While Clouds come naturally to the rescue of Smart Grids for dealing with this big data issue, little attention has been paid to the benefits that Smart Grids could bring to distributed Clouds. To our knowledge, no previous work has exploited the Smart Grids potential to obtain and control the energy consumption of entire Cloud infrastructures from underlying facilities such as air conditioning equipment (which accounts for 30% to 50% of a data center's electricity bill) to network resources (which are often operated by several actors) and to computing resources (with their heterogeneity and distribution across multiple data centers). We aim at taking advantage of the opportunity brought by the Smart Grids to exploit renewable energy availability and to optimize energy management in distributed Clouds.

### **3.3.2. Energy cost models**

Cloud computing allows users to outsource the computer resources required for their applications instead of using a local installation. It offers on-demand access to the resources through the Internet with a pay-as-you-go pricing model. However, this model hides the electricity cost of running these infrastructures.

The costs of current data centers are mostly driven by their energy consumption (specifically by the air conditioning, computing and networking infrastructures). Yet, current pricing models are usually static and rarely consider the facilities' energy consumption per user. The challenge is to provide a fair and predictable model to attribute the overall energy costs per virtual machine and to increase energy-awareness of users.

Another goal consists in better understanding the energy consumption of computing and networking resources of Clouds in order to provide energy cost models for the entire infrastructure including incentivizing cost models for both Cloud providers and energy suppliers. These models will be based on experimental measurement campaigns on heterogeneous devices. Inferring a cost model from energy measurements is an arduous task since simple models are not convincing, as shown in our previous work. We aim at proposing and validating energy cost models for the heterogeneous Cloud infrastructures in one hand, and the energy distribution grid on the other hand. These models will be integrated into simulation frameworks in order to validate our energy-efficient algorithms at larger scale.

### **3.3.3. Energy-aware users**

In a moderately loaded Cloud, some servers may be turned off when not used for energy saving purpose. Cloud providers can apply resource management strategies to favor idle servers. Some of the existing solutions propose mechanisms to optimize VM scheduling in the Cloud. A common solution is to consolidate the mapping of the VMs in the Cloud by grouping them in a fewer number of servers. The unused servers can then be turned off in order to lower the global electricity consumption.

Indeed, current work focuses on possible levers at the virtual machine suppliers and/or services. However, users are not involved in the choice of using these levers while significant energy savings could be achieved with their help. For example, they might agree to delay slightly the calculation of the response to their applications on the Cloud or accept that it is supported by a remote data center, to save energy or wait for the availability of renewable energy. The VMs are black boxes from the Cloud provider point of view. So, the user is the only one to know the applications running on her VMs.

We plan to explore possible collaborations between virtual machine suppliers, service providers and users of Clouds in order to provide users with ways of participating in the reduction of the Clouds energy consumption. This work will follow two directions: 1) to investigate compromises between power and performance/service quality that cloud providers can offer to their users and to propose them a variety of options adapted to their workload; and 2) to develop mechanisms for each layer of the Cloud software stack to provide users with a quantification of the energy consumed by each of their options as an incentive to become greener.

## **3.4. Securing clouds**

### **3.4.1. Security monitoring SLO**

While the trend for companies to outsource their information system in clouds is confirmed, the problem of securing an information system becomes more difficult. Indeed, in the case of infrastructure clouds, physical

resources are shared between companies (also called tenants) but each tenant controls only parts of the shared resources, and, thanks to virtualization, the information system can be dynamically and automatically reconfigured with added or removed resources (for example starting or stopping virtual machines), or even moved between physical resources (for example using virtual machine migration). Partial control of shared resources brings new classes of attacks between tenants, and security monitoring mechanisms to detect such attacks are better placed out of the tenant-controlled virtual information systems, that is under control of the cloud provider. Dynamic and automatic reconfigurations of the information system make it unfeasible for a tenant's security administrator to setup the security monitoring components to detect attacks, and thus an automated self-adaptable security monitoring service is required.

Combining the two previous statements, there is a need for a dependable, automatic security monitoring service provided to tenants by the cloud provider. Our goal is to address the following challenges to design such a security monitoring service:

1. to define relevant Service-Level Objectives (SLOs) of a security monitoring service, that can figure in the Service-Level Agreement (SLA) signed between a cloud provider and a tenant;
2. to design heuristics to automatically configure provider-controlled security monitoring software components and devices so that SLOs are reached, even during automatic reconfigurations of tenants' information systems;
3. to design evaluation methods for tenants to check that SLOs are reached.

Moreover in challenges 2 and 3 the following sub-challenges must be addressed:

- although SLAs are bi-lateral contracts between the provider and each tenant, the implementation of the contracts is based on shared resources, and thus we must study methods to combine the SLOs;
- the designed methods should have a minimal impact on performance.

### **3.4.2. Data protection in Cloud-based IoT services**

The Internet of Things is becoming a reality. Individuals have their own swarm of connected devices (e.g. smartphone, wearables, and home connected objects) continually collecting personal data. A novel generation of services is emerging exploiting data streams produced by the devices' sensors. People are deprived of control of their personal data as they don't know precisely what data are collected by service providers operating on Internet (oISP), for which purpose they could be used, for how long they are stored, and to whom they are disclosed. In response to privacy concerns the European Union has introduced, with the Global Data Protection Regulation (GDPR), new rules aimed at enforcing the people's rights to personal data protection. The GDPR also gives strong incentives to oISPs to comply. However, today, oISPs can't make their systems GDPR-compliant since they don't have the required technologies. We argue that a new generation of system is mandatory for enabling oISPs to conform to the GDPR. We plan to design an open source distributed operating system for native implementation of new GDPR rules and ease the programming of compliant cloud-based IoT services. Among the new rules, transparency, right of erasure, and accountability are the most challenging ones to be implemented in IoT environments but could fundamentally increase people's confidence in oISPs. Deployed on individuals' swarms of devices and oISPs' cloud-hosted servers, it will enforce detailed data protection agreements and accountability of oISPs' data processing activities. Ultimately we will show to what extent the new GDPR rules can be implemented for cloud-based IoT services.

## **3.5. Experimenting with Clouds**

Cloud platforms are challenging to evaluate and study with a sound scientific methodology. As with any distributed platform, it is very difficult to gather a global and precise view of the system state. Experiments are not reproducible by default since these systems are shared between several stakeholders. This is even worsened by the fact that microscopic differences in the experimental conditions can lead to drastic changes since typical Cloud applications continuously adapt their behavior to the system conditions.

### 3.5.1. Experimentation methodologies for clouds

We propose to combine two complementary experimental approaches: direct execution on testbeds such as Grid'5000, that are eminently convincing but rather labor intensive, and simulations (using *e.g.*, SimGrid) that are much more light-weighted, but requires careful assessment. One specificity of the Myriads team is that we are working on these experimental methodologies *per se*, raising the standards of *good experiments* in our community.

We plan to make SimGrid widely usable beyond research laboratories, in order to evaluate industrial systems and to teach the future generations of cloud practitioners. This requires to frame the specific concepts of Cloud systems and platforms in actionable interfaces. The challenge is to make the framework both easy to use for simple studies in educational settings while modular and extensible to suit the specific needs of every advanced industrial-class users.

We aim at leveraging the convergence opportunities between methodologies by further bridging simulation and real testbeds. The predictions obtained from the simulator should be validated against some real-world experiments obtained on the target production platform, or on a similar platform. This (in)validation of the predicted results often improves the understanding of the modeled system. On the other side, it may even happen that the measured discrepancies are due to some mis-configuration of the real platform that would have been undetected without this (in)validation study. In that sense, the simulator constitutes a precious tool for the quality assurance of real testbeds such as Grid'5000.

Scientists need more help to make their Cloud experiments fully reproducible, in the spirit of Open Science exemplified by the HAL Open Archive, actively backed by Inria. Users still need practical solutions to archive, share and compare the whole experimental settings, including the raw data production (particularly in the case of real testbeds) and their statistical analysis. This is a long lasting task to which we plan to collaborate through the research communities gathered around the Grid'5000 and SimGrid scientific instruments.

Finally, since correction and performance can constitute contradictory goals, it is particularly important to study them jointly. To that extend, we want to bridge the performance studies, that constitute our main scientific heritage, to correction studies leveraging formal techniques. SimGrid already includes support to exhaustively explore the possible executions. We plan to continue this work to ease the use of the relevant formal methods to the experimenter studying Cloud systems.

### 3.5.2. Use cases

In system research it is important to work on real-world use cases from which we extract requirements inspiring new research directions and with which we can validate the system services and mechanisms we propose. In the framework of our close collaboration with the Data Science Technology department of the LBNL, we will investigate cloud usage for scientific data management. Next-generation scientific discoveries are at the boundaries of datasets, *e.g.*, across multiple science disciplines, institutions and spatial and temporal scales. Today, data integration processes and methods are largely adhoc or manual. A generalized resource infrastructure that integrates knowledge of the data and the processing tasks being performed by the user in the context of the data and resource lifecycle is needed. Clouds provide an important infrastructure platform that can be leveraged by including knowledge for distributed data integration.



## SPIRALS Project-Team

### 3. Research Program

#### 3.1. Introduction

Our research program on self-adaptive software targets two key properties that are detailed in the remainder of this section: *self-healing* and *self-optimization*.

#### 3.2. Objective #1: Self-healing - Mining software artifacts to automatically evolve systems

Software systems are under the pressure of changes all along their lifecycle. Agile development blurs the frontier between design and execution and requires constant adaptation. The size of systems (millions of lines of code) multiplies the number of bugs by the same order of magnitude. More and more systems, such as sensor network devices, live in "surviving" mode, in the sense that they are neither rebootable nor upgradable.

Software bugs are hidden in source code and show up at development-time, testing-time or worse, once deployed in production. Except for very specific application domains where formal proofs are achievable, bugs can not be eradicated. As an order of magnitude, on 16 Dec 2011, the Eclipse bug repository contains 366 922 bug reports. Software engineers and developers work on bug fixing on a daily basis. Not all developers spend the same time on bug fixing. In large companies, this is sometimes a full-time role to manage bugs, often referred to as *Quality Assurance* (QA) software engineers. Also, not all bugs are equal, some bugs are analyzed and fixed within minutes, others may take months to be solved [79].

In terms of research, this means that: (i) one needs means to automatically adapt the design of the software system through automated refactoring and API extraction, (ii) one needs approaches to automate the process of adapting source code in order to fix certain bugs, (iii) one needs to revisit the notion of error-handling so that instead of crashing in presence of errors, software adapts itself to continue with its execution, *e.g.*, in degraded mode.

There is no one-size-fits-all solution for each of these points. However, we think that novel solutions can be found by using **data mining and machine learning techniques tailored for software engineering** [80]. This body of research consists of mining some knowledge about a software system by analyzing the source code, the version control systems, the execution traces, documentation and all kinds of software development and execution artifacts in general. This knowledge is then used within recommendation systems for software development, auditing tools, runtime monitors, frameworks for resilient computing, etc.

The novelty of our approach consists of using and tailoring data mining techniques for analyzing software artifacts (source code, execution traces) in order to achieve the **next level of automated adaptation** (*e.g.*, automated bug fixing). Technically, we plan to mix unsupervised statistical learning techniques (*e.g.* frequent item set mining) and supervised ones (*e.g.* training classifiers such as decision trees). This research is currently not being performed by data mining research teams since it requires a high level of domain expertise in software engineering, while software engineering researchers can use off-the-shelf data mining libraries, such as Weka [58].

We now detail the two directions that we propose to follow to achieve this objective.

##### 3.2.1. Learning from software history how to design software and fix bugs

The first direction is about mining techniques in software repositories (*e.g.*, CVS, SVN, Git). Best practices can be extracted by data mining source code and the version control history of existing software systems. The design and code of expert developers significantly vary from the artifacts of novice developers. We will learn to differentiate those design characteristics by comparing different code bases, and by observing the semantic refactoring actions from version control history. Those design rules can then feed the test-develop-refactor constant adaptation cycle of agile development.

**Fault localization of bugs reported in bug repositories.** We will build a solid foundation on empirical knowledge about bugs reported in bug repository. We will perform an empirical study on a set of representative bug repositories to identify classes of bugs and patterns of bug data. For this, we will build a tool to browse and annotate bug reports. Browsing will be helped with two kinds of indexing: first, the tool will index all textual artifacts for each bug report; second it will index the semantic information that is not present by default in bug management software—*i.e.*, “contains a stacktrace”). Both indexes will be used to find particular subsets of bug reports, for instance “all bugs mentioning invariants and containing a stacktrace”. Note that queries with this kind of complexity and higher are mostly not possible with the state-of-the-art of bug management software. Then, analysts will use annotation features to annotate bug reports. The main outcome of the empirical study will be the identification of classes of bugs that are appropriate for automated localization. Then, we will run machine learning algorithms to identify the latent links between the bug report content and source code features. Those algorithms would use as training data the existing traceability links between bug reports and source code modifications from version control systems. We will start by using decision trees since they produce a model that is explicit and understandable by expert developers. Depending on the results, other machine learning algorithms will be used. The resulting system will be able to locate elements in source code related to a certain bug report with a certain confidence.

**Automated bug fix generation with search-based techniques.** Once a location in code is identified as being the cause of the bug, we can try to automatically find a potential fix. We envision different techniques: (1) infer fixes from existing contracts and specifications that are violated; (2) infer fixes from the software behavior specified as a test suite; (3) try different fix types one-by-one from a list of identified bug fix patterns; (4) search fixes in a fix space that consists of combinations of atomic bug fixes. Techniques 1 and 2 are explored in [54] and [78]. We will focus on the latter techniques. To identify bug fix patterns and atomic bug fixes, we will perform a large-scale empirical study on software changes (also known as changesets when referring to changes across multiple files). We will develop tools to navigate, query and annotate changesets in a version control system. Then, a grounded theory will be built to master the nature of fixes. Eventually, we will decompose change sets in atomic actions using clustering on changeset actions. We will then use this body of empirical knowledge to feed search-based algorithms (*e.g.* genetic algorithms) that will look for meaningful fixes in a large fix space. To sum up, our research on automated bug fixing will try not only to point to source code locations responsible of a bug, but to search for code patterns and snippets that may constitute the skeleton of a valid patch. Ultimately, a blend of expert heuristics and learned rules will be able to produce valid source code that can be validated by developers and committed to the code base.

### 3.2.2. Run-time self-healing

The second proposed research direction is about inventing a self-healing capability at run-time. This is complementary to the previous objective that mainly deals with development time issues. We will achieve this in two steps. First, we want to define frameworks for resilient software systems. Those frameworks will help to maintain the execution even in the presence of bugs—*i.e.* to let the system survive. As exposed below, this may mean for example to switch to some degraded modes. Next, we want to go a step further and to define solutions for automated runtime repair, that is, not simply compensating the erroneous behavior, but also determining the correct repair actions and applying them at run-time.

**Mining best effort values.** A well-known principle of software engineering is the “fail-fast” principle. In a nutshell, it states that as soon as something goes wrong, software should stop the execution before entering incorrect states. This is fine when a human user is in the loop, capable of understanding the error or at least rebooting the system. However, the notion of “failure-oblivious computing” [71] shows that in certain domains, software should run in a resilient mode (*i.e.* capable of recovering from errors) and/or best-effort mode—*i.e.* a slightly imprecise computation is better than stopping. Hence, we plan to investigate data mining techniques in order to learn best-effort values from past executions (*i.e.* somehow learning what is a correct state, or the opposite what is not a completely incorrect state). This knowledge will then be used to adapt the software state and flow in order to mitigate the error consequences, the exact opposite of fail-fast for systems with long-running cycles.

**Embedding search based algorithms at runtime.** Harman recently described the field of search-based software engineering [59]. We believe that certain search based approaches can be embedded at runtime with the goal of automatically finding solutions that avoid crashing. We will create software infrastructures that allow automatically detecting and repairing faults at run-time. The methodology for achieving this task is based on three points: (1) empirical study of runtime faults; (2) learning approaches to characterize runtime faults; (3) learning algorithms to produce valid changes to the software runtime state. An empirical study will be performed to analyze those bug reports that are associated with runtime information (*e.g.* core dumps or stacktraces). After this empirical study, we will create a system that learns on previous repairs how to produce small changes that solve standard runtime bugs (*e.g.* adding an array bound check to throw a handled domain exception rather than a spurious language exception). To achieve this task, component models will be used to (1) encapsulate the monitoring and reparation meta-programs in appropriate components and (2) support runtime code modification using scripting, reflective or bytecode generation techniques.

### 3.3. Objective #2: Self-optimization - Sharing runtime behaviors to continuously adapt software

Complex distributed systems have to seamlessly adapt to a wide variety of deployment targets. This is due to the fact that developers cannot anticipate all the runtime conditions under which these systems are immersed. A major challenge for these software systems is to develop their capability to continuously reason about themselves and to take appropriate decisions and actions on the optimizations they can apply to improve themselves. This challenge encompasses research contributions in different areas, from environmental monitoring to real-time symptoms diagnosis, to automated decision making. The variety of distributed systems, the number of optimization parameters, and the complexity of decisions often resign the practitioners to design monolithic and static middleware solutions. However, it is now globally acknowledged that the development of dedicated building blocks does not contribute to the adoption of sustainable solutions. This is confirmed by the scale of actual distributed systems, which can—for example—connect several thousands of devices to a set of services hosted in the Cloud. In such a context, the lack of support for smart behaviors at different levels of the systems can inevitably lead to its instability or its unavailability. In June 2012, an outage of Amazon’s Elastic Compute Cloud in North Virginia has taken down Netflix, Pinterest, and Instagram services. During hours, all these services failed to satisfy their millions of customers due to the lack of integration of a self-optimization mechanism going beyond the boundaries of Amazon.

The research contributions we envision within this area will therefore be organized as a reference model for engineering **self-optimized distributed systems** autonomously driven by *adaptive feedback control loops*, which will automatically enlarge their scope to cope with the complexity of the decisions to be taken. This solution introduces a multi-scale approach, which first privileges local and fast decisions to ensure the homeostasis<sup>0</sup> property of a single node, and then progressively propagates symptoms in the network in order to reason on a longer term and a larger number of nodes. Ultimately, domain experts and software developers can be automatically involved in the decision process if the system fails to find a satisfying solution. The research program for this objective will therefore focus on the study of mechanisms for **monitoring, taking decisions, and automatically reconfiguring software at runtime and at various scales**. As stated in the self-healing objective, we believe that there is no one-size-fits-all mechanism that can span all the scales of the system. We will therefore study and identify an optimal composition of various adaptation mechanisms in order to produce long-living software systems.

The novelty of this objective is to exploit the wisdom of crowds to define new middleware solutions that are able to continuously adapt software deployed in the wild. We intend to demonstrate the applicability of this approach to distributed systems that are deployed from mobile phones to cloud infrastructures. The key scientific challenges to address can be summarized as follows: *How does software behave once deployed in the wild? Is it possible to automatically infer the quality of experience, as it is perceived by users? Can the*

<sup>0</sup>Homeostasis is the property of a system that regulates its internal environment and tends to maintain a stable, relatively constant condition of properties [Wikipedia].

*runtime optimizations be shared across a wide variety of software? How optimizations can be safely operated on large populations of software instances?*

The remainder of this section further elaborates on the opportunities that can be considered within the frame of this objective.

### 3.3.1. *Monitoring software in the wild*

Once deployed, developers are generally no longer aware of how their software behave. Even if they heavily use testbeds and benchmarks during the development phase, they mostly rely on the bugs explicitly reported by users to monitor the efficiency of their applications. However, it has been shown that contextual artifacts collected at runtime can help to understand performance leaks and optimize the resilience of software systems [81]. Monitoring and understanding the context of software at runtime therefore represent the first building block of this research challenge. Practically, we intend to investigate crowd-sensing approaches, to smartly collect and process runtime metrics (*e.g.*, request throughput, energy consumption, user context). Crowd-sensing can be seen as a specific kind of **crowdsourcing** activity, which refers to the capability of lifting a (large) diffuse group of participants to delegate the task of retrieving trustable data from the field. In particular, crowd-sensing covers not only *participatory sensing* to involve the user in the sensing task (*e.g.*, surveys), but also *opportunistic sensing* to exploit mobile sensors carried by the user (*e.g.*, smartphones).

While reported metrics generally enclose raw data, the monitoring layer intends to produce meaningful indicators like the *Quality of Experience* (QoE) perceived by users. This QoE reflects representative symptoms of software requiring to trigger appropriate decisions in order to improve its efficiency. To diagnose these symptoms, the system has to process a huge variety of data including runtime metrics, but also history of logs to explore the sources of the reported problems and identify opportunities for optimizations. The techniques we envision at this level encompass **machine learning**, **principal component analysis**, and fuzzy logic [70] to provide enriched information to the decision level.

### 3.3.2. *Collaborative decision-making approaches*

Beyond the symptoms analysis, decisions should be taken in order to improve the *Quality of Service* (QoS). In our opinion, collaborative approaches represent a promising solution to effectively converge towards the most appropriate optimization to apply for a given symptom. In particular, we believe that exploiting the **wisdom of the crowd** can help the software to optimize itself by sharing its experience with other software instances exhibiting similar symptoms. The intuition here is that the body of knowledge that supports the optimization process cannot be specific to a single software instance as this would restrain the opportunities for improving the quality and the performance of applications. Rather, we think that any software instance can learn from the experience of others.

With regard to the state-of-the-art, we believe that a multi-levels decision infrastructure, inspired from distributed systems like Spotify [57], can be used to build a decentralized decision-making algorithm involving the surrounding peers before requesting a decision to be taken by more central control entity. In the context of collaborative decision-making, peer-based approaches therefore consist in quickly reaching a consensus on the decision to be adopted by a majority of software instances. Software instances can share their knowledge through a micro-economic model [51], that would weight the recommendations of experienced instances, assuming their age reflects an optimal configuration.

Beyond the peer level, the adoption of algorithms inspired from evolutionary computations, such as **genetic programming**, at an upper level of decision can offer an opportunity to test and compare several alternative decisions for a given symptom and to observe how does the crowd of applications evolves. By introducing some diversity within this population of applications, some instances will not only provide a satisfying QoS, but will also become naturally resilient to unforeseen situations.

### 3.3.3. *Smart reconfigurations in the large*

Any decision taken by the crowd requires to propagate back to and then operated by the software instances. While simplest decisions tend to impact software instances located on a single host (*e.g.*, laptop, smartphone),

this process can also exhibit more complex reconfiguration scenarios that require the orchestration of various actions that have to be safely coordinated across a large number of hosts. While it is generally acknowledged that centralized approaches raise scalability issues, we think that self-optimization should investigate different reconfiguration strategies to propagate and apply the appropriate actions. The investigation of such strategies can be addressed in two steps: the consideration of *scalable data propagation protocols* and the identification of *smart reconfiguration mechanisms*.

With regard to the challenge of scalable data propagation protocols, we think that research opportunities encompass not only the exploitation of gossip-based protocols [56], but also the adoption of publish/subscribe abstractions [64] in order to decouple the decision process from the reconfiguration. The fundamental issue here is the definition of a communication substrate that can accommodate the propagation of decisions with relaxed properties, inspired by *Delay Tolerant Networks* (DTN), in order to reach weakly connected software instances. We believe that the adoption of asynchronous communication protocols can provide the sustainable foundations for addressing various execution environments including harsh environments, such as developing countries, which suffer from a partial connectivity to the network. Additionally, we are interested in developing the principle of *social networks of applications* in order to seamlessly group and organize software instances according to their similarities and acquaintances. The underlying idea is that grouping application instances can contribute to the identification of optimization profiles not only contributing to the monitoring layer, but also interested in similar reconfigurations. Social networks of applications can contribute to the anticipation of reconfigurations by exploiting the symptoms of similar applications to improve the performance of others before that problems actually happen.

With regard to the challenge of smart reconfiguration mechanisms, we are interested in building on our established experience of adaptive middleware [75] in order to investigate novel approaches to efficient application reconfigurations. In particular, we are interested in adopting seamless micro-updates and micro-reboot techniques to provide in-situ reconfiguration of pieces of software. Additionally, the provision of safe and secured reconfiguration mechanisms is clearly a key issue that requires to be carefully addressed in order to avoid malicious exploitation of dynamic reconfiguration mechanisms against the software itself. In this area, although some reconfiguration mechanisms integrate transaction models [65], most of them are restricted to local reconfigurations, without providing any support for executing distributed reconfiguration transactions. Additionally, none of the approached published in the literature include security mechanisms to preserve from unauthorized or malicious reconfigurations.

## STACK Project-Team

### 3. Research Program

#### 3.1. Overview

STACK research activities have been organized around four research topics. The two first ones are related to the resource management mechanisms and the programming support that are mandatory to operate and use ICT geo-distributed resources (compute, storage, network). They are transverse to the System/Middleware/Application layers, which generally composed a software stack, and nurture each other (*i.e.*, the resource management mechanisms will leverage abstractions/concepts proposed by the programming support axis and reciprocally). The third and fourth research topics are related to the Energy and Security dimensions (both also crosscutting the three software layers). Although they could have been merged with the two first axes, we identified them as independent research directions due to their critical aspects with respect to the societal challenges they represent. In the following, we detail the actions we plan to do in each research direction.

#### 3.2. Resource Management

The challenge in this axis is to identify, design or revise mechanisms that are mandatory to operate and use a set of massively geo-distributed resources in an efficient manner [50]. This covers considering challenges at the scale of nodes, within one site (*i.e.*, one geographical location) and throughout the whole geo-distributed ICT infrastructure. It is noteworthy that the network community has been investigating similar challenges for the last few years [69]. To benefit from their expertise, in particular on how to deal with intermittent networks, STACK members have recently initiated exchanges and collaborative actions with some network research groups and telcos (see Sections 8.1 and 9.1). We emphasize, however, that we do not deliver contributions related to network equipments/protocols. The scientific and technical achievements we aim to deliver are related to the (distributed) system aspects.

##### 3.2.1. Performance Characterization of Low-Level Building Blocks

Although Cloud Computing has enabled the consolidation of services and applications into a subset of servers, current operating system mechanisms do not provide appropriate abstractions to prevent (or at least control) the performance degradation that occurs when several workloads compete for the same resources [101]. Keeping in mind that server density is going to increase with physical machines composed of more and more cores and that applications will be more and more data intensive, it is mandatory to identify interferences that appear at a low level on each dimension (compute, memory, network, and storage) and propose countermeasures. In particular, previous studies [101], [61] on pros and cons of current technologies – virtual machines (VMs) [75], [83], containers and microservices – which are used to consolidate applications on the same server, should be extended: In addition to evaluating the performance we can expect from each of these technologies on a single node, it is important to investigate interferences that may result from cross-layer and remote communications [102]. We will consider in particular all interactions related to geo-distributed systems mechanisms/services that are mandatory to operate and use geo-distributed ICT infrastructures.

##### 3.2.2. Geo-Distributed System Mechanisms

Although several studies have been highlighting the advantages of geo-distributed ICT infrastructures in various domains (see Section 3.), progress on how to operate and use such infrastructures is marginal. Current solutions [32] [33] are rather close to the initial Cisco Fog Computing proposal that only allows running domain-specific applications on edge resources and centralized Cloud platforms [40] (in other words, these solutions do not allow running stateful workloads in isolated environments such as containers or VMs). More recently, solutions leveraging the idea of federating VIMs (as the aforementioned ETSI MEC proposal [88]) have been proposed. ONAP [95], an industry-driven solution, enables the orchestration and



automation of virtual network functions across distinct VIMs. From the academic side, FogBow [42] aims to support federations of Infrastructure-as-a-Service (IaaS) providers. Finally, NIST initiated a collaborative effort with IEEE to advance Federated Cloud platforms through the development of a conceptual architecture and a vocabulary<sup>0</sup>. Although all these projects provide valuable contributions, they face the aforementioned orchestration limitations (*i.e.*, they do not manage decisions taken in each VIM). Moreover, they all have been designed by only considering the developer/user's perspective. They provide abstractions to manage the life cycle of geo-distributed applications, but do not address administrative requirements.

To cope with specifics of Wide-Area networks while delivering most features that made Cloud Computing solutions successful also at the edge, our community should first identify limitations/drawbacks of current resource management system mechanisms with respect to the Fog/Edge requirements and propose revisions when needed [68], [81].

To achieve this aim, STACK members propose to conduct first a series of studies aiming at understanding the software architecture and footprint of major services that are mandatory for operating and using Fog/Edge infrastructures (storage backends, monitoring services, deployment/reconfiguration mechanisms, etc.). Leveraging these studies, we will investigate how these services should be deployed in order to deal with resources constraints, performance variability, and network split brains. We will rely on contributions that have been accomplished in distributed algorithms and self-\* approach for the last decade. In the short and medium term, we plan to evaluate the relevance of NewSQL systems [53] to store internal states of distributed system mechanisms in an edge context, and extend our proposals on new storage backends such as key/value stores [52], [94], and burst buffers [103]. We also plan to conduct new investigations on data-stream frameworks for Fog and Edge infrastructures [47]. These initial contributions should enable us to identify general rules to deliver other advanced system mechanisms that will be mandatory at the higher levels in particular for the deployment and reconfiguration manager in charge of orchestrating all resources.

### 3.2.3. Capacity Planning and Placement Strategies

An objective shared by users and providers of ICT infrastructures is to limit as much as possible the operational costs while providing the expected and requested quality of service (QoS). To optimize this cost while meeting QoS requirements, data and applications have to be placed in the best possible way onto physical resources according to data sources, data types (stream, graphs), application constraints (real-time requirements) and objective functions. Furthermore, the placement of applications must evolve through time to cope with the fluctuations in terms of application resource needs as well as the physical events that occur at the infrastructure level (resource creation/removals, hardware failures, etc.). This placement problem, *a.k.a.* the deployment and reconfiguration challenge as it will be described in Section 3.3, can be modeled in many different ways, most of the time by multi-dimensional and multi-objective bin-packing problems or by scheduling problems which are known to be difficult to solve. Many studies have been performed, for example, to optimize the placement of virtual machines onto ICT infrastructures [77]. STACK will inherit the knowledge acquired through previous activities in this domain, particularly its use of constraint programming strategies in autonomic managers [73], [72], relying on MAPE (monitor, analyze, plan, and execute) control loops. While constraint programming approaches are known to hardly scale, they enable the composition of various constraints without requiring to change heuristic algorithms each time a new constraint has to be considered [71]. We believe it is a strong advantage to deal with the diversity brought by geo-distributed ICT infrastructures. Moreover, we have shown in previous work that decentralized approaches can tackle the scalability issue while delivering placement decisions good enough and sometimes close to the optimal [87].

Leveraging this expertise, we propose, first, to identify new constraints raised by massively geo-distributed infrastructures (*e.g.*, data locality, energy, security, reliability and the heterogeneity and mobility of the underlying infrastructure). Based on this preliminary study, we will explore new placement strategies not only for computation sandboxes but for data (location, replication, streams, etc.) in order to benefit from the geo-distribution of resources and meet the required QoS. These investigations should lead to collaborations with operational research and optimization groups such as TASC, another research group from IMT Atlantique.

<sup>0</sup><https://collaborate.nist.gov/twiki-cloud-computing/bin/view/CloudComputing/FederatedCloudPWGFC> (Dec 2018).

Second, we will leverage contributions made on the previous axis “Performance Characterization of Low-Level Building Blocks” to determine how the deployment of the different units (software components and data sets) should be executed in order to reduce as much as possible the time to reconfigure the system (*i.e.*, the *Execution* phase in the control loop). In some recent work [83], we have shown that the provisioning of a new virtual machine should be done carefully to mitigate boot penalties. More generally, proposing an efficient action plan for the *Execution* phase will be a major point as Wide-Area-Network specifics may lead to significant delays, in particular when the amount of data to be manipulated is important.

Finally, we will investigate new approaches to decentralize the placement process while considering the geo-distributed context. Among the different challenges to address, we will study how a combination of autonomic managers, at both the infrastructure and application levels [60], could be proposed in a decentralized manner. Our first idea is to geo-distribute a fleet of small control loops over the whole infrastructure. By improving the locality of data collection and weakening combinatorics, these loops would allow the system to address responsiveness and quality expectations.

### 3.3. Programming Support

We pursue two main research directions relative to new programming support: first, developing new programming models with appropriate support in existing languages (libraries, embedded DSLs, etc.) and, second, providing new means for deployment and reconfiguration in geo-distributed ICT environments, principally supporting the mapping of software onto the infrastructure. For both directions two levels of challenges are considered. On the one hand, the *generic* level refers to efforts on programming support that can be applied to any kind of distributed software, application or system. On this level, contributions could thus be applied to any of the three layers addressed by STACK (*i.e.*, system, middleware or application). On the other hand, the corresponding generic programming means may not be appropriate in practice (*e.g.*, requirements for more dedicated support, performance constraints, etc.), even if they may lead to interesting general properties. For this reason, a *specific* level is also considered. This level could be based on the generic one but addresses specific cases or domains.

#### 3.3.1. Programming Models and Languages Extensions

The current landscape of programming support for cloud applications is fragmented. This fragmentation is based on apparently different needs for various kinds of applications, in particular, web-based, computation-based, focusing on the organization of the computation, and data-based applications, within the last case a quite strong dichotomy between applications considering data as sets or relations, close to traditional database applications and applications considering data as real-time streams. This has led to various programming models, in a loose sense, including for instance microservices, graph processing, dataflows, streams, etc. These programming models have mostly been offered to the application programmer in the guise of frameworks, each offering subtle variants of the programming models with various implementation decisions favoring particular application and infrastructure settings. Whereas most frameworks are dedicated to a given programming model, *e.g.*, basic Pregel [82], Hive [97], Hadoop [98], some of them are more general-purpose through the provision of several programming models, *e.g.*, Flink [46] and Spark [79]. Finally, some dedicated language support has been considered for some models (*e.g.*, the language SPL underlying IBM Streams [74]) as well as core languages and calculi (*e.g.*, [43], [92]).

This situation raises a number of challenges on its own, related to a better structuring of the landscape. It is necessary to better understand the various programming models and their possible relations, with the aim of facilitating, if not their complete integration, at least their composition, at the conceptual level but also with respect to their implementations, as specific languages and frameworks.

Switching to massively geo-distributed infrastructures adds to these challenges by leading to a new range of applications (*e.g.*, smart-\* applications) that, by nature, require mixing these various programming models, together with a much more dynamic management of their runtime.



In this context, STACK would like to explore two directions:

- First, we propose to contribute to generic programming models and languages to address composability of different programming models [55]. For example, providing a generic stream data processing model that can operate under both data stream [46] and operation stream [104] modes, thus streams can be processed in micro batches to favour high throughput or record by record to sustain low latency. Software engineering properties such as separation of concerns and composition should help address such challenges [35], [93]. They should also facilitate the software deployment and reconfiguration challenges discussed below.
- Second, we plan to revise relevant programming models, the associated specific languages, and their implementation according to the massive geo-distribution of the underlying infrastructure, the data sources, and application end-users. For example, although SPL is extensible and distributed, it has been designed to run on multi-cores and clusters [74]. It does not provide the level of dynamicity required by geo-distributed applications (*e.g.*, to handle topology changes, loss of connectivity at the edge, etc.). Moreover, as more network data transfers will happen within a massively geo-distributed infrastructure, correctness of data transfers should be guaranteed. This has potential impact from the programming models to their implementations.

### 3.3.2. Deployment and Reconfiguration Challenges

The second research direction deals with the complexity of deploying distributed software (whatever the layer, application, middleware or system) onto an underlying infrastructure. As both the deployed pieces of software and the infrastructures addressed by STACK are large, massively distributed, heterogeneous and highly dynamic, the deployment process cannot be handled manually by developers or administrators. Furthermore, and as already mentioned in Section 3.2, the initial deployment of some distributed software will evolve through time because of the dynamicity of both the deployed software and the underlying infrastructures. When considering reconfiguration, which encompasses deployment as a specific case, the problem becomes more difficult for two main reasons: (1) the current state of both the deployed software and the infrastructure has to be taken into account when deciding on a reconfiguration plan, (2) as the software is already running the reconfiguration should minimize disruption time, while avoiding inconsistencies [80], [85]. Many deployment tools have been proposed both in academia and industry [57]. For example, Ansible<sup>0</sup>, Chef<sup>0</sup> and Puppet<sup>0</sup> are very well-known generic tools to automate the deployment process through a set of batch instructions organized in groups (*e.g.*, *playbooks* in Ansible). Some tools are specific to a given environment, like Kolla to deploy OpenStack, or the embedded deployment manager within Spark. Few reconfiguration capabilities are available in production tools such as *scaling* and *restart* after a fault<sup>0 0</sup>. Academia has contributed to generic deployment and reconfiguration models. Most of these contributions are component-based. Component models divide a distributed software as a set of component instances (or modules) and their assembly, where components are connected through well defined interfaces [93]. Thus, modeling the reconfiguration process consists in describing the life cycle of different components and their interactions. Most component-based approaches offer a fixed life cycle, *i.e.*, identical for any component [62]. Two main contributions are able to customize life cycles, Fractal [45], [38] and its evolutions [35], [36], [59], and Aeolus [54]. In Fractal, the *control* part of a component (*e.g.*, its life cycle) is modeled itself as a component assembly that is highly flexible. Aeolus, on the other hand, offers a finer control on both the evolution and the synchronization of the deployment process by modeling each component life cycle with a finite state machine.

A reconfiguration raises at least five questions, all of them are correlated: (1) *why software has to be reconfigured?* (monitoring, modeling and analysis) (2) *what should be reconfigured?* (software modeling and analysis), (3) *how should it be reconfigured?* (software modeling and planning decisions), (4) *where should it*

<sup>0</sup><https://www.ansible.com/>

<sup>0</sup><https://www.chef.io/chef/>

<sup>0</sup><https://puppet.com/>

<sup>0</sup><https://kubernetes.io/>

<sup>0</sup><https://jjujucharms.com/>

*be reconfigured?* (infrastructure modeling and planning decisions), and (5) *when to reconfigure it?* (scheduling algorithms). STACK will contribute to all aspects of a reconfiguration process as described above. However, according to the expertise of STACK members, we will focus mainly on the three first questions: *why*, *what* and *how*, leaving questions *where* and *when* to collaborations with operational research and optimization teams.

First of all, we would like to investigate *why software has to be reconfigured?* Many reasons could be mentioned, such as hardware or software fault tolerance, mobile users, dynamicity of software services, etc. All those reasons are related somehow to the Quality of Service (QoS) or the Service Level Agreement (SLA) between the user and the Cloud provider. We first would like to explore the specificities of QoS and SLAs in the case of massively geo-distributed ICT environments [89]. By being able to formalize this question, analyzing the requirement of a reconfiguration will be facilitated.

Second, we think that four important properties should be enhanced when deploying and reconfiguring models in massively geo-distributed ICT environments. First, as low-latency applications and systems will be subject to deployment and reconfiguration, the performance and the ability to scale are important. Second, as many different kinds of deployments and reconfigurations will concurrently hold within the infrastructure, processes have to be reliable, which is facilitated by a fine-grained control of the process. Finally, as many different software elements will be subject to deployment and reconfiguration, common generic models and engines for deployment and reconfiguration should be designed [44]. For these reasons, we intend to go beyond Aeolus by: first, leveraging the expression of parallelism within the deployment process, which should lead to better performance; second, improving the separation of concerns between the component developer and the reconfiguration developer; third, enhancing the possibility to perform concurrent and decentralized reconfigurations.

Research challenges relative to programming support have been presented above. Many of these challenges are related, in different manners, to the resource management level of STACK or to crosscutting challenges, *i.e.*, energy and security. First, one can notice that any programming model or deployment and reconfiguration implementation should be based on mechanisms related to resource management challenges. For this reason, all challenges addressed within this section are linked with lower level building blocks presented in Section 3.2. Second, as detailed above, deployment and reconfiguration address at least five questions. The question *what?* is naturally related to programming support. However, questions *why*, *how?*, *where?* and *when?* are also related to Section 3.2, for example, to monitoring and capacity planning. Moreover, regarding the deployment and reconfiguration challenges, one can note that the same goals recursively happen when deploying the control building blocks themselves (bootstrap issue). This comforts the need to design generic deployment and reconfiguration models and frameworks. These low-level models should then be used as back-ends to higher-level solutions. Finally, as *energy* and *security* are crosscutting themes within the STACK project, many additional energy and security considerations could be added to the above challenges. For example, our deployment and reconfiguration frameworks and solutions could be used to guarantee the deployment of end-to-end security policies or to answer specific energy constraints [70] as detailed in the next section.

### 3.4. Energy

The overall electrical consumption of DCs grows according to the demand of Utility Computing. Considering that the latter has been continuously increasing since 2008, the energy footprint of Cloud services overall is nowadays critical with about 91 billion kilowatt-hours of electricity [91]. Besides the ecological impact, the energy consumption is a predominant criterion for providers since it determines a large part of the operational cost of their infrastructure. Among the different approaches that have been investigated to reduce the energy footprint, some studies have been investigating the use of renewable energy sources to power microDCs [64]. Workload distribution for geo-distributed DCs is also another promising approach [66], [78], [99]. Our research will extend these results with the ultimate goal of considering the different opportunities to control the energy footprint across the whole stack (hardware and software opportunities, renewable energy, thermal management, etc.). In particular, we identified several challenges that we will address in this context within the STACK framework.

First, we propose to evaluate the energy efficiency of low-level building blocks, from the viewpoints of computation (VMs, containers, microkernel, microservices) [58] and data (hard drives, SSD, in-memory storage, distributed file systems). For computations, in the continuity of our previous work [56], [73], we will investigate workload placement policies according to energy (minimizing energy consumption, power capping, thermal load balancing, etc.). Regarding the data dimension, we will investigate, in particular, the trade-offs between energy consumption and data availability, durability and consistency [51], [94]. Our ambition is to propose an adaptive energy-aware data layout and replication scheme to ensure data availability with minimal energy consumption. It is noteworthy that these new activities will also consider our previous work on DCs partially powered by renewable energy (see the SeDuCe project, in Section 6.7), with the ultimate goal of reducing the CO<sub>2</sub> footprint.

Second, we will complete current studies to understand pros and cons of massively geo-distributed infrastructures from the energy perspective. Addressing the energy challenge is a complex task that involves considering several dimensions such as the energy consumption due to the physical resources (CPU, memory, disk, network), the performance of the applications (from the computation and data viewpoints), and the thermal dissipation caused by air conditioning in each DC. Each of these aspects can be influenced by each level of the software stack (*i.e.*, low-level building blocks, coordination and autonomous loops, and finally application life cycle). In previous projects, we have studied and modeled the consumption of the main components, notably the network, as part of a single microDC. We plan to extend these models to deal with geo-distribution. The objective is to propose models that will enable us to refine our placement algorithms as discussed in the next paragraph. These models should be able to consider the energy consumption induced by all WAN data exchanges, including site-to-site data movements as well as the end users' communications for accessing virtualized resources.

Third, we expect to implement green-energy-aware balancing strategies, leveraging the aforementioned contributions. Although the infrastructures we envision increase complexity (because WAN aspects should also be taken into account), the geo-distribution of resources brings several opportunities from the energy viewpoint. For instance, it is possible to define several workload/data placement policies according to renewable energy availability. Moreover, a tightly-coupled software stack allows users to benefit from such a widely distributed infrastructure in a transparent way while enabling administrators to balance resources in order to benefit from green energy sources when available. An important difficulty, compared to centralized infrastructures, is related to data sharing between software instances. In particular, we will study issues raised by the distribution and replication of services across several microDCs. In this new context, many challenges must be addressed: where to place the data (Cloud, Edge) in order to mitigate data movements? What is the impact in terms of energy consumption, network and response time of these two approaches? How to manage the consistency of replicated data/services? All these aspects must be studied and integrated into our placement algorithms.

Fourth, we will investigate the energy footprint of the current techniques that address failure and performance variability in large-scale systems. For instance, *stragglers* (*i.e.*, tasks that take a significantly longer time to finish than the normal execution time) are natural results of performance variability, they cause extra resource and energy consumption. Our goal is to understand the energy overhead of these techniques and introduce new handling techniques that take into consideration the energy efficiency of the platform [86].

Finally, in order to answer specific energy constraints, we want to reify energy aspects at the application level and propose a metric related to the use of energy (Green SLA [34]), for example to describe the maximum allowed CO<sub>2</sub> emissions of a Fog/Edge service. Unlike other approaches [67], [39], [65] that attempt to identify the best trade-off, we want to offer to developers/end-users the opportunity to select the best choice between application performance, correctness and energy footprint. Such a capability will require reifying the energy dimension at the level of big-data and interactive applications. Besides, with the emergence of renewable energy (*e.g.*, solar panels for microDC), investigating the energy consumption *vs* performance trade-off [70] and the smart usage of green energy for ICT geo-distributed services seems promising. For example, we want to offer the opportunity to developers/end-users to control the scaling of the applications based on this trade-

off instead of current approaches that only considered application load. Providing such a capability will also require appropriate software abstractions.

### 3.5. Security

Because of its large size and complex software structure, geo-distributed applications and infrastructures are particularly exposed to security and privacy issues [90]. They are subject to numerous security vulnerabilities that are frequently exploited by malicious attackers in order to exfiltrate personal, institutional or corporate data. Securing these systems require security and privacy models and corresponding techniques that are applicable at all software layers in order to guard interactions at each level but also between levels. However, very few security models exist for the lower layers of the software stack and no model enables the handling of interactions involving the complete software stack. Any modification to its implementation, deployment status, configuration, etc., may introduce new or trigger existing security and privacy issues. Finally, applications that execute on top of the software stack may introduce security issues or be affected by vulnerabilities of the stack. Overall, security and privacy issues are therefore interdependent with all other activities of the STACK team and constitute an important research topic for the team.

As part of the STACK activities, we consider principally security and privacy issues related to the vertical and horizontal compositions of software components forming the software stack and the distributed applications running on top of it. Modifications to the *vertical composition* of the software stack affect different software levels at once. As an example, side-channel attacks often target virtualized services (*i.e.*, services running within VMs); attackers may exploit insecure hardware caches at the system level to exfiltrate data from computations at the higher level of VM services [84], [100]. Security and privacy issues also affect *horizontal compositions*, that is, compositions of software abstractions on one level: most frequently horizontal compositions are considered on the level of applications/services but they are also relevant on the system level or the middleware level, such as compositions involving encryption and database fragmentation services.

The STACK members aim at addressing two main research issues: enabling full-stack (vertical) security and per-layer (horizontal) security. Both of these challenges are particularly hard in the context of large geo-distributed systems because they are often executed on heterogeneous infrastructures and are part of different administrative domains and governed by heterogeneous security and privacy policies. For these reasons they typically lack centralized control, are frequently subject to high latency and are prone to failures.

Concretely, we will consider two classes of security and privacy issues in this context. First, on a general level, we strive for a method for the programming and reasoning about compositions of security and privacy mechanisms including, but not limited to, encryption, database fragmentation and watermarking techniques. Currently, no such general method exists, compositions have only been devised for specific and limited cases, for example, compositions that support the commutation of specific encryption and watermarking techniques [76], [48]. We provided preliminary results on such compositions [49] and have extended them to biomedical, notably genetic, analyses in the e-health domain [41]. Second, on the level of security and privacy properties, we will focus on isolation properties that can be guaranteed through vertical and horizontal composition techniques. We have proposed first results in this context in form of a compositional notion of distributed side channel attacks that operate on the system and middleware levels [37].

It is noteworthy that the STACK members do not have to be experts on the individual security and privacy mechanisms, such as watermarking and database fragmentation. We are, however, well-versed in their main properties so that we can integrate them into our composition model. We also interact closely with experts in these techniques and the corresponding application domains, notably e-health for instance, in the context of the PRIVGEN project<sup>0</sup>, see Section 9.1 .

More generally, we highlight that security issues in distributed systems are very closely related to the other STACK challenges, dimensions and research directions. Guaranteeing security properties across the software stack and throughout software layers in highly volatile and heterogeneous geo-distributed systems is expected to harness and contribute results to the self-management capabilities investigated as part of the team's resource

<sup>0</sup>Privacy-preserving sharing and processing of genetic data, <https://privgen.cominlabs.u-bretagne.fr/fr>

management challenges. Furthermore, security and privacy properties are crosscutting concerns that are intimately related to the challenges of application life cycle management. Similarly, the security issues are also closely related to the team's work on programming support. This includes new means for programming, notably in terms of event and stream programming, but also the deployment and reconfiguration challenges, notably concerning automated deployment. As a crosscutting functionality, the security challenges introduced above must be met in an integrated fashion when designing, constructing, executing and adapting distributed applications as well as managing distributed resources.

## WHISPER Project-Team

### 3. Research Program

#### 3.1. Scientific Foundations

##### 3.1.1. Program analysis

A fundamental goal of the research in the Whisper team is to elicit and exploit the knowledge found in existing code. To do this in a way that scales to a large code base, systematic methods are needed to infer code properties. We may build on either static [36], [38], [40] or dynamic analysis [58], [60], [65]. Static analysis consists of approximating the behavior of the source code from the source code alone, while dynamic analysis draws conclusions from observations of sample executions, typically of test cases. While dynamic analysis can be more accurate, because it has access to information about actual program behavior, obtaining adequate test cases is difficult. This difficulty is compounded for infrastructure software, where many, often obscure, cases must be handled, and external effects such as timing can have a significant impact. Thus, we expect to primarily use static analyses. Static analyses come in a range of flavors, varying in the extent to which the analysis is *sound*, *i.e.*, the extent to which the results are guaranteed to reflect possible run-time behaviors.

One form of sound static analysis is *abstract interpretation* [38]. In abstract interpretation, atomic terms are interpreted as sound abstractions of their values, and operators are interpreted as functions that soundly manipulate these abstract values. The analysis is then performed by interpreting the program in a compositional manner using these abstracted values and operators. Alternatively, *dataflow analysis* [49] iteratively infers connections between variable definitions and uses, in terms of local transition rules that describe how various kinds of program constructs may impact variable values. Schmidt has explored the relationship between abstract interpretation and dataflow analysis [73]. More recently, more general forms of symbolic execution [36] have emerged as a means of understanding complex code. In symbolic execution, concrete values are used when available, and these are complemented by constraints that are inferred from terms for which only partial information is available. Reasoning about these constraints is then used to prune infeasible paths, and obtain more precise results. A number of works apply symbolic execution to operating systems code [33], [34].

While sound approaches are guaranteed to give correct results, they typically do not scale to the very diverse code bases that are prevalent in infrastructure software. An important insight of Engler et al. [42] was that valuable information could be obtained even when sacrificing soundness, and that sacrificing soundness could make it possible to treat software at the scales of the kernels of the Linux or BSD operating systems. Indeed, for certain types of problems, on certain code bases, that may mostly follow certain coding conventions, it may mostly be safe to *e.g.*, ignore the effects of aliases, assume that variable values are unchanged by calls to unanalyzed functions, etc. Real code has to be understood by developers and thus cannot be too complicated, so such simplifying assumptions are likely to hold in practice. Nevertheless, approaches that sacrifice soundness also require the user to manually validate the results. Still, it is likely to be much more efficient for the user to perform a potentially complex manual analysis in a specific case, rather than to implement all possible required analyses and apply them everywhere in the code base. A refinement of unsound analysis is the CEGAR approach [37], in which a highly approximate analysis is complemented by a sound analysis that checks the individual reports of the approximate analysis, and then any errors in reasoning detected by the sound analysis are used to refine the approximate analysis. The CEGAR approach has been applied effectively on device driver code in tools developed at Microsoft [25]. The environment in which the driver executes, however, is still represented by possibly unsound approximations.



Going further in the direction of sacrificing soundness for scalability, the software engineering community has recently explored a number of approaches to code understanding based on techniques developed in the areas of natural language understanding, data mining, and information retrieval. These approaches view code, as well as other software-related artifacts, such as documentation and postings on mailing lists, as bags of words structured in various ways. Statistical methods are then used to collect words or phrases that seem to be highly correlated, independently of the semantics of the program constructs that connect them. The obliviousness to program semantics can lead to many false positives (invalid conclusions) [55], but can also highlight trends that are not apparent at the low level of individual program statements. We have previously explored combining such statistical methods with more traditional static analysis in identifying faults in the usage of constants in Linux kernel code [53].

### 3.1.2. Domain Specific Languages

Writing low-level infrastructure code is tedious and difficult, and verifying it is even more so. To produce non-trivial programs, we could benefit from moving up the abstraction stack to enable both programming and proving as quickly as possible. Domain-specific languages (DSLs), also known as *little languages*, are a means to that end [6] [61].

#### 3.1.2.1. Traditional approach.

Using little languages to aid in software development is a tried-and-trusted technique [75] by which programmers can express high-level ideas about the system at hand and avoid writing large quantities of formulaic C boilerplate.

This approach is typified by the Devil language for hardware access [8]. An OS programmer describes the register set of a hardware device in the high-level Devil language, which is then compiled into a library providing C functions to read and write values from the device registers. In doing so, Devil frees the programmer from having to write extensive bit-manipulation macros or inline functions to map between the values the OS code deals with, and the bit-representation used by the hardware: Devil generates code to do this automatically.

However, DSLs are not restricted to being “stub” compilers from declarative specifications. The Bossa language [7] is a prime example of a DSL involving imperative code (syntactically close to C) while offering a high-level of abstraction. This design of Bossa enables the developer to implement new process scheduling policies at a level of abstraction tailored to the application domain.

Conceptually, a DSL both abstracts away low-level details and justifies the abstraction by its semantics. In principle, it reduces development time by allowing the programmer to focus on high-level abstractions. The programmer needs to write less code, in a language with syntax and type checks adapted to the problem at hand, thus reducing the likelihood of errors.

#### 3.1.2.2. Embedding DSLs.

The idea of a DSL has yet to realize its full potential in the OS community. Indeed, with the notable exception of interface definition languages for remote procedure call (RPC) stubs, most OS code is still written in a low-level language, such as C. Where DSL code generators are used in an OS, they tend to be extremely simple in both syntax and semantics. We conjecture that the effort to implement a given DSL usually outweighs its benefit. We identify several serious obstacles to using DSLs to build a modern OS: specifying what the generated code will look like, evolving the DSL over time, debugging generated code, implementing a bug-free code generator, and testing the DSL compiler.

Filet-o-Fish (FoF) [39] addresses these issues by providing a framework in which to build correct code generators from semantic specifications. This framework is presented as a Haskell library, enabling DSL writers to *embed* their languages within Haskell. DSL compilers built using FoF are quick to write, simple, and compact, but encode rigorous semantics for the generated code. They allow formal proofs of the runtime behavior of generated code, and automated testing of the code generator based on randomized inputs, providing greater test coverage than is usually feasible in a DSL. The use of FoF results in DSL compilers that OS developers can quickly implement and evolve, and that generate provably correct code. FoF has been used

to build a number of domain-specific languages used in Barrelfish, [26] an OS for heterogeneous multicore systems developed at ETH Zurich.

The development of an embedded DSL requires a few supporting abstractions in the host programming language. FoF was developed in the purely functional language Haskell, thus benefiting from the type class mechanism for overloading, a flexible parser offering convenient syntactic sugar, and purity enabling a more algebraic approach based on small, composable combinators. Object-oriented languages – such as Smalltalk [43] and its descendant Pharo [30] – or multi-paradigm languages – such as the Scala programming language [63] – also offer a wide range of mechanisms enabling the development of embedded DSLs. Perhaps surprisingly, a low-level imperative language – such as C – can also be extended so as to enable the development of embedded compilers [27].

### 3.1.2.3. Certifying DSLs.

Whilst automated and interactive software verification tools are progressively being applied to larger and larger programs, we have not yet reached the point where large-scale, legacy software – such as the Linux kernel – could formally be proved “correct”. DSLs enable a pragmatic approach, by which one could realistically strengthen a large legacy software by first narrowing down its critical component(s) and then focus our verification efforts onto these components.

Dependently-typed languages, such as Coq or Idris, offer an ideal environment for embedding DSLs [35], [31] in a unified framework enabling verification. Dependent types support the type-safe embedding of object languages and Coq’s mixfix notation system enables reasonably idiomatic domain-specific concrete syntax. Coq’s powerful abstraction facilities provide a flexible framework in which to not only implement and verify a range of domain-specific compilers [39], but also to combine them, and reason about their combination.

Working with many DSLs optimizes the “horizontal” compositionality of systems, and favors reuse of building blocks, by contrast with the “vertical” composition of the traditional compiler pipeline, involving a stack of comparatively large intermediate languages that are harder to reuse the higher one goes. The idea of building compilers from reusable building blocks is a common one, of course. But the interface contracts of such blocks tend to be complex, so combinations are hard to get right. We believe that being able to write and verify formal specifications for the pieces will make it possible to know when components can be combined, and should help in designing good interfaces.

Furthermore, the fact that Coq is also a system for formalizing mathematics enables one to establish a close, formal connection between embedded DSLs and non-trivial domain-specific models. The possibility of developing software in a truly “model-driven” way is an exciting one. Following this methodology, we have implemented a certified compiler from regular expressions to x86 machine code [48]. Interestingly, our development crucially relied on an existing Coq formalization, due to Braibant and Pous, [32] of the theory of Kleene algebras.

While these individual experiments seem to converge toward embedding domain-specific languages in rich type theories, further experimental validation is required. Indeed, Barrelfish is an extremely small software compared to the Linux kernel. The challenge lies in scaling this methodology up to large software systems. Doing so calls for a unified platform enabling the development of a myriad of DSLs, supporting code reuse across DSLs as well as providing support for mechanically-verified proofs.

## 3.2. Research direction: Tools for improving legacy infrastructure software

A cornerstone of our work on legacy infrastructure software is the Coccinelle program matching and transformation tool for C code. Coccinelle has been in continuous development since 2005. Today, Coccinelle is extensively used in the context of Linux kernel development, as well as in the development of other software, such as wine, python, kvm, and systemd. Currently, Coccinelle is a mature software project, and no research is being conducted on Coccinelle itself. Instead, we leverage Coccinelle in other research projects [28], [29], [64], [66], [70], [72], [74], [59], [54], both for code exploration, to better understand at a large scale problems in Linux development, and as an essential component in tools that require program matching and transformation. The continuing development and use of Coccinelle is also a source of visibility in the Linux kernel developer



community. We submitted the first patches to the Linux kernel based on Coccinelle in 2007. Since then, over 5500 patches have been accepted into the Linux kernel based on the use of Coccinelle, including around 3000 by over 500 developers from outside our research group.

Our recent work has focused on driver porting. Specifically, we have considered the problem of porting a Linux device driver across versions, particularly backporting, in which a modern driver needs to be used by a client who, typically for reasons of stability, is not able to update their Linux kernel to the most recent version. When multiple drivers need to be backported, they typically need many common changes, suggesting that Coccinelle could be applicable. Using Coccinelle, however, requires writing backporting transformation rules. In order to more fully automate the backporting (or symmetrically forward porting) process, these rules should be generated automatically. We have carried out a preliminary study in this direction with David Lo of Singapore Management University; this work, published at ICSME 2016 [77], is limited to a port from one version to the next one, in the case where the amount of change required is limited to a single line of code. Whisper has been awarded an ANR PRCI grant to collaborate with the group of David Lo on scaling up the rule inference process and proposing a fully automatic porting solution.

### 3.3. Research direction: developing infrastructure software using Domain Specific Languages

We wish to pursue a *declarative* approach to developing infrastructure software. Indeed, there exists a significant gap between the high-level objectives of these systems and their implementation in low-level, imperative programming languages. To bridge that gap, we propose an approach based on domain-specific languages (DSLs). By abstracting away boilerplate code, DSLs increase the productivity of systems programmers. By providing a more declarative language, DSLs reduce the complexity of code, thus the likelihood of bugs.

Traditionally, systems are built by accretion of several, independent DSLs. For example, one might use Devil [8] to interact with devices, Bossa [7] to implement the scheduling policies. However, much effort is duplicated in implementing the back-ends of the individual DSLs. Our long term goal is to design a unified framework for developing and composing DSLs, following our work on Filet-o-Fish [39]. By providing a single conceptual framework, we hope to amortize the development cost of a myriad of DSLs through a principled approach to reusing and composing them.

Beyond the software engineering aspects, a unified platform brings us closer to the implementation of mechanically-verified DSLs. Using the Coq proof assistant as an x86 macro-assembler [48] is a step in that direction, which belongs to a larger trend of hosting DSLs in dependent type theories [31], [35], [62]. A key benefit of those approaches is to provide – by construction – a formal, mechanized semantics to the DSLs thus developed. This semantics offers a foundation on which to base further verification efforts, whilst allowing interaction with non-verified code. We advocate a methodology based on incremental, piece-wise verification. Whilst building fully-certified systems from the top-down is a worthwhile endeavor [50], we wish to explore a bottom-up approach by which one focuses first and foremost on crucial subsystems and their associated properties.

Our current work on DSLs has two complementary goals: (i) the design of a unified framework for developing and composing DSLs, following our work on Filet-o-Fish, and (ii) the design of domain-specific languages for domains where there is a critical need for code correctness, and corresponding methodologies for proving properties of the run-time behavior of the system.

## WIDE Project-Team

### 3. Research Program

#### 3.1. Overview

In order to progress in the four fields described above, the WIDE team is developing a research program which aims to **help developers control and master the inherent uncertainties and performance challenges brought by scale and distribution**.

More specifically, our program revolves around four key challenges.

- Challenge 1: Designing Hybrid Scalable Architectures,
- Challenge 2: Constructing Personalizable Privacy-Aware Distributed Systems,
- Challenge 3: Understanding Controllable Network Diffusion Processes,
- Challenge 4: Systemizing Modular Distributed Computability and Efficiency.

These four challenges have in common **the inherent tension between coordination and scalability in large-scale distributed systems**: strong coordination mechanisms can deliver strong guarantees (in terms of consistency, agreement, fault-tolerance, and privacy protection), but are generally extremely costly and inherently non-scalable if applied indiscriminately. By contrast, highly scalable coordination approaches (such as epidemic protocols, eventual consistency, or self-organizing overlays) perform much better when the size of a system increases, but do not, in most cases, provide any strong guarantees in terms of consistency or agreement.

The above four challenges explore these tensions from *four complementary angles*: from an architectural perspective (Challenge 1), from the point of view of a fundamental system-wide guarantee (privacy protection, Challenge 2), looking at one universal scalable mechanism (network diffusion, Challenge 3), and considering the interplay between modularity and computability in large-scale systems (Challenge 4). These four challenges range from practical concerns (Challenges 1 and 2) to more theoretical questions (Challenges 3 and 4), yet present *strong synergies* and *fertile interaction points*. E.g. better understanding network diffusion (Challenge 3) is a key enabler to develop more private decentralized systems (Challenge 2), while the development of a theoretically sound modular computability hierarchy (Challenge 4) has a direct impact on our work on hybrid architectures (Challenge 1).

#### 3.2. Hybrid Scalable Architectures

The rise of planetary-scale distributed systems calls for novel software and system architectures that can support user-facing applications while scaling to large numbers of devices, and leveraging established and emerging technologies. The members of WIDE are particularly well positioned to explore this avenue of research thanks to their experience on de-concentrated architectures combining principles from both decentralized peer-to-peer [48], [58] systems and hybrid infrastructures (i.e. architectures that combines centralized or hierarchical elements, often hosted in well-provisioned data-centers, and a decentralized part, often hosted in a peer-to-peer overlay) [52]. In the short term, we aim to explore two axes in this direction: browser-based communication, and micro services.

##### 3.2.1. Browser-based fog computing

The dramatic increase in the amount of data being produced and processed by connected devices has led to paradigms that seek to decentralize the traditional cloud model. In 2011 Cisco [49] introduced the vision of *fog computing* that combines the cloud with resources located at the edge of the network and in between. More generally, the term *edge computing* has been associated with the idea of adding edge-of-the-network storage and computation to traditional cloud infrastructures [44].

A number of efforts in this directions focus on specific hardware, e.g. fog nodes that are responsible for connected IoT devices [50]. However, many of today's applications run within web browsers or mobile phones. In this context, the recent introduction of the WebRTC API, makes it possible for browsers and smartphones to exchange directly between each other, enabling mobile, or browser-based decentralized applications. Maygh [72], for example, uses the WebRTC API to build a decentralized Content Delivery Network that runs solely on web browsers. The fact that the application is hosted completely on a web server and downloaded with enabled websites means that webmasters can adopt the Content Delivery Network (CDN) without requiring users to install any specific software.

For us, the ability of browsers to communicate with each other using the WebRTC paradigm provides a novel playground for new programming models, and for a *browser-based fog architecture* combining both a centralized, cloud-based part, and a decentralized, browser-supported part.

This model offers tremendous potential by making edge-of-the-network resources available through the interconnection of web-browsers, and offers new opportunities for the protection of the personal data of end users. But consistently engineering browser-based components requires novel tools and methodologies.

In particular, WebRTC was primarily designed for exchanging media and data between two browsers in the presence of a coordinating server. Its complex mechanisms for connection establishment make many of the existing peer-to-peer protocols inefficient. To address this challenge, we plan to consider two angles of attack. First, we plan to design novel protocols that take into account the specific requirements set by this new technology. Second, we envisage to investigate variants of the current WebRTC model with cheaper connection-establishment protocols, in order to provide lower delays and bandwidth consumption in large-scale browser-based applications.

We also plan to address the trade-offs associated with hybrid browser-cloud models. For example, when should computation be delegated to browsers and when should it be executed on the cloud in order to maximize the quality of service? Or, how can a decentralized analytics algorithms operating on browser-based data complement or exploit the knowledge built by cloud-based data analytics solutions?

### 3.2.2. *Emergent micro-service deployment and management*

Micro-services tend to produce fine-grained applications in which many small services interact in a loosely coupled manner to produce a wide range of services within an organization. Individual services need to evolve independently of each other over time without compromising the availability of the overall application. Lightweight isolation solutions such as containers (Docker, ...), and their associated tooling ecosystem (e.g. Google's Borg [71], Kubernetes [47]) have emerged to facilitate the deployment of large-scale micro-service-based applications, but only provide preliminary solutions for key concerns in these systems, which we would like to investigate and extend.

Most of today's on-line computer systems are now too large to evolve in monolithic, entirely pre-planned ways. This applies to very large data centres, for example, where the placement of virtual machines to reduce heating and power consumption can no longer be treated using top-down exhaustive optimisation approaches beyond a critical size. This is also true of social networking applications, where different mechanisms—e.g. to spread news notifications, or to recommend new contacts—must be adapted to the different sub-communities present in the system.

To cope with the inherent complexity of building complex loosely-coupled distributed systems while fostering and increasing efficiency, maintainability, and scalability, we plan to study how novel programming techniques based on declarative programming, components and epidemic protocols can help design, deploy, and maintain self-adaptive structures (e.g. placement of VM) and mechanisms (e.g. contact recommendations) that are optimized to the local context of very large distributed systems. To fulfill this vision, we plan to explore a three-pronged strategy to raise the level of programming abstraction offered to developers.

- First, we plan to explore the use of high-level domain-specific languages (DSL) to declare how large-scale topologies should be achieved, deployed, and maintained. Our vision is a declarative approach to describe how to combine, deploy and orchestrate micro-services in an abstract manner

thus abstracting away developers from the underlying cloud infrastructures, and from the intricacies involved in writing low-level code to build a large-scale distributed application that scales. With this effort, we plan notably to directly support the twin properties of *emergence* (the adaptation “from within”) and *differentiation* (the possibility from parts of the system to diverge while still forming a whole). Our central objective is to search for principled programming constructs to support these two capabilities using a modular and incremental software development approach.

- On a second strand of work, we plan to investigate how unikernels enable smaller footprints, more optimization options, and faster boot times for micro-services. Isolating micro-services into VMs is not the most adequate approach as it requires the use of hypervisors, or virtual machine monitors (VMMs), to virtualize hardware resources. VMMs are well known to be heavyweight with both boot and run time overheads that may have a strong impact on performances. Unikernels seem to offer the right balance between performance and flexibility to address this challenge. One of the key underlying challenges is to compile directly the aforementioned provided DSL to a dedicated and customized machine image, ready to be deployed directly on top of a large set of bare metal servers.
- Depending on the workload it is subjected to, and the state of its execution environment (network, VMs), a large-scale distributed application may present erratic or degraded performance that is hard to anticipate and plan for. There is therefore a strong need to adapt dynamically the way resources are allocated to a running application. We would like to study how the DSL approach we envisage can be extended to enable developers to express orchestration algorithms based on machine learning algorithms.

### 3.3. Personalizable Privacy-Aware Distributed Systems

On-line services are increasingly moving towards an in-depth analysis of user data, with the objective of providing ever better personalization. But in doing so, personalized on-line services inevitably pose risks to the privacy of users. Eliminating, or even reducing these risks raises important challenges caused by the inherent trade-off between the level of personalization users wish to achieve, and the amount of information they are willing to reveal about themselves (explicitly or through the many implicit sources of digital information such as smart homes, smart cars, and IoT environments).

At a general level, we would like to address these challenges through protocols that can provide access to unprecedented amounts of data coming from sensors, users, and documents published by users, while protecting the privacy of individuals and data sources. To this end, we plan to rely on our experience in the context of distributed systems, recommender systems, and privacy, as well as in our collaborations with experts in neighboring fields such as machine learning, and security. In particular, we aim to explore different privacy-utility tradeoffs that make it possible to provide differentiated levels of privacy guarantees depending on the context associated with data, on the users that provide the data, and on those that access it. Our research targets the general goal of privacy-preserving decentralized learning, with applications in different contexts such as user-oriented applications, and the Internet-of-Things (IoT).

#### 3.3.1. Privacy-preserving decentralized learning

Personalization and recommendation can be seen as a specific case of general machine learning. Production-grade recommenders and personalizers typically centralize and process the available data in one location (a data-center, a cloud service). This is highly problematic, as it endangers the privacy of users, while hampering the analysis of datasets subject to privacy constraints that are held by multiple independent organizations (such as health records). A decentralized approach to machine learning appears as a promising candidate to overcome these weaknesses: if each user or participating organization keeps its data, while only exchanging gradient or model information, privacy leaks seem less likely to occur.

In some cases, decentralized learning may be achieved through relatively simple adaptations of existing centralized models, for instance by defining alternative learning models that may be more easily decentralized. But in all cases, processing growing amounts of information calls for high-performance algorithms and middleware that can handle diverse storage and computation resources, in the presence of dynamic and

privacy-sensitive data. To reach this objective, we will therefore leverage our work in distributed and privacy-preserving algorithms and middleware [51], [53], [54] as well as the results of our work on large-scale hybrid architectures in Objective 1.

### 3.3.2. Personalization in user-oriented applications

As a first application perspective, we plan to design tools that exploit decentralized analytics to enhance user-centric personalized applications. As we observed above, such applications exhibit an inherent trade-off between personalization quality and privacy preservation. The most obvious goal in this direction consists in designing algorithms that can achieve high levels of personalization while protecting sensitive user information. But an equally important one consists in personalizing the trade-off itself by adapting the quality of the personalization provided to a user to his/her willingness to expose information. This, like other desirable behaviors, appears at odds with the way current systems work. For example, a user of a recommender system that does not reveal his/her profile information penalizes other users causing them to receive less accurate recommendations. We would like to mitigate this situation by means of protocols that reward users for sharing information. On the one hand, we plan to take inspiration from protocols for free-riding avoidance in peer-to-peer systems [55], [60]. On the other hand, we will consider blockchains as a tool for tracking and rewarding data contributions. Ultimately, we aim at enabling users to configure the level of privacy and personalization they wish to experience.

### 3.3.3. Privacy preserving decentralized aggregation

As a second setting we would like to consider target applications running on constrained devices like in the Internet-of-Things (IoT). This setting makes it particularly important to operate on decentralized data in a light-weight privacy-preserving manner, and further highlights the synergy between this objective and Objective 1. For example, we plan to provide data subjects with the possibility to store and manage their data locally on their own devices, without having to rely on third-party managers or aggregators, but possibly storing less private information or results in the cloud. Using this strategy, we intend to design protocols that enable users themselves, or third-party companies to query distributed data in aggregate form, or to run data analytics processes on a distributed set of data repositories, thereby gathering knowledge without violating the privacy of other users. For example, we have started working on the problem of computing an aggregate function over a subset of the data in a distributed setting. This involves two major steps: selection and aggregation. With respect to selection, we envision defining a decentralized data-selection operation that can apply a selection predicate without violating privacy constraints. With respect to aggregation, we will continue our investigation of lightweight protocols that can provide privacy with limited computational complexity [45].

## 3.4. Network Diffusion Processes

Social, biological, and technological networks can serve as conduits for the spread of ideas, trends, diseases, or viruses. In social networks, rumors, trends and behaviors, or the adoption of new products, spread from person to person. In biological networks, diseases spread through contact between individuals, and mutations spread from an individual to its offsprings. In technological networks, such as the Internet and the power grid, viruses and worms spread from computer to computer, and power failures often lead to cascading failures. The common theme in all the examples above is that the rumor, disease, or failure starts out with a single or a few individual nodes, and propagates through the network, from node to node, to reach a potentially much larger number of nodes.

These types of *network diffusion processes* have long been a topic of study in various disciplines, including sociology, biology, physics, mathematics, and more recently, computer science. A main goal has been to devise mathematical models for these processes, describing how the state of an individual node can change as a function of the state of its neighbors in the network, and then analyse the role of the network structure in the outcome of the process. Based on our previous work, we would like to study to what extent one can affect the outcome of the diffusion process by controlling a small, possibly carefully selected fraction of the network.

For example, we plan to explore how we may increase the spread or speed of diffusion by choosing an appropriate set of seed nodes (a standard goal in viral marketing by word-of-mouth), or achieve the opposite effect either by choosing a small set of nodes to remove (a goal in immunization against diseases), or by seeding a competing diffusion (e.g., to limit the spread of misinformation in a social network).

Our goal is to provide a framework for a systematic and rigorous study of these problems. We will consider several standard diffusion models and extensions of them, including models from mathematical sociology, mathematical epidemiology, and interacting particle systems. We will consider existing and new variants of spread maximization/limitation problems, and will provide (approximation) algorithms or show negative (inapproximability) results. In case of negative results, we will investigate general conditions that make the problem tractable. We will consider both general network topologies and specific network models, and will relate the efficiency of solutions to structural properties of the topology. Finally, we will use these insights to engineer new network diffusion processes for efficient data dissemination.

### 3.4.1. Spread maximization

Our goal is in particular to study spread maximization in a broader class of diffusion processes than the basic independent cascade (IC) and linear threshold (LT) models of influence [64], [65], [66] that have been studied in this context so far. This includes the *randomized rumor spreading (RS)* model for information dissemination [57], *biased versions of the voter model* [61] modelling influence, and the (graph-based) *Moran processes* [68] modelling the spread of mutations. We would like to consider several natural versions of the spread maximization problem, and the relationships between them. For these problems we will use the greedy algorithm and the submodularity-based analytical framework of [64], and will also explore new approaches.

### 3.4.2. Immunization optimization

Conversely we would also like to explore immunization optimization problems. Existing works on these types of problem assume a *perfect-contagion* model, i.e., once a node gets infected, it deterministically infects all its non-immunized neighbors. We plan to consider various diffusion processes, including the standard *susceptible–infected (SI)*, *susceptible–infected–recovered (SIR)* and *susceptible–infected–susceptible (SIS)* epidemic models, and explore the extent to which results and techniques for the perfect-contagion model carry over to these probabilistic models. We will also investigate whether techniques for spread maximization could be applied to immunization problems.

Some immunization problems are known to be hard to approximate in general graphs, even for the perfect-contagion model, e.g., the fixed-budget version of the fire-fighter problem cannot be approximated to any  $n^{1-\epsilon}$  factor [46]. This strand of work will consider restricted graph families, such as trees or graphs of small treewidth, for such problems. In addition, for some immunization problems, there is a large gap between the best known approximation algorithm and the best known inapproximability result, and we would like to make progress in reducing these gaps.

## 3.5. Systemizing Modular Distributed Computability and Efficiency

The applications and services envisaged in Objectives 1 and 2 will lead to increasingly complex and multifaceted systems. Constructing these novel hybrid and decentralized systems will naturally push our need to understand distributed computing beyond the current state of the art. These trends therefore demand research efforts in establishing sound theoretical foundations to allow everyday developers to master the design, properties and implementation of these systems. We plan to investigate these foundations along two directions: first by studying novel approaches to some fundamental problems of *mutual exclusion and distributed coordination*, and second by exploring how we can build a *comprehensive and modular framework* capturing the foundations of *distributed computation*.



### 3.5.1. Randomized algorithm for mutual exclusion and coordination

To exploit the power of massive distributed applications and systems (such as those envisaged in Objectives 1 and 2) or multiple processors, algorithms must cope with the scale and asynchrony of these systems, and their inherent instability, e.g., due to node, link, or processor failures. Our goal is to explore the power and limits of randomized algorithms for large-scale networks of distributed systems, and for shared memory multi-processor systems, in effect providing fundamental building blocks to the work envisioned in Objectives 1 and 2.

For shared memory systems, randomized algorithms have notably proved extremely useful to deal with asynchrony and failures. Sometimes probabilistic algorithms provide the only solution to a problem; sometimes they are more efficient; sometimes they are simply easier to implement. We plan to devise efficient algorithms for some of the fundamental problems of shared memory computing, such as mutual exclusion, renaming, and consensus.

In particular, looking at the problem of *mutual exclusion*, it is desirable that mutual exclusion algorithms be *abortable*. This means that a process that is trying to lock the resource can abort its attempt in case it has to wait too long. Abortability is difficult to achieve for mutual exclusion algorithms. We will try to extend our algorithms for the *cache-coherent* (CC) and the *distributed shared memory* (DSM) model in order to make them abortable, while maintaining expected constant *Remote Memory References* (RMRs) complexity, under optimistic system assumptions. In order to achieve this, the algorithm will use strong synchronization primitives, called compare-and-swap objects. As part of our collaboration with the University of Calgary, we will work on implementing those objects from registers in such a way that they also allow aborts. Our goal is to build on existing non-abortable implementations [59]. We plan then later to use these objects as building blocks in our mutual exclusion algorithm, in order to make them work even if the system does not readily provide such primitives.

We have also started working on blockchains, as these represent a new and interesting trade-off between probabilistic guarantees, scalability, and system dynamics, while revisiting some of the fundamental questions and limitations of consensus in fault-prone asynchronous systems.

### 3.5.2. Modular theory of distributed computing

Practitioners and engineers have proposed a number of reusable frameworks and services to implement specific distributed services (from Remote Procedure Calls with Java RMI or SOAP-RPC, to JGroups for group communication, and Apache Zookeeper for state machine replication). In spite of the high conceptual and practical interest of such frameworks, many of these efforts lack a sound grounding in distributed computation theory (with the notable exceptions of JGroups and Zookeeper), and often provide punctual and partial solutions for a narrow range of services. We argue that this is because we still lack a generic framework that unifies the large body of fundamental knowledge on distributed computation that has been acquired over the last 40 years.

To overcome this gap we would like to develop a systematic model of distributed computation that organizes the functionalities of a distributed computing system into reusable modular constructs assembled via well-defined mechanisms that maintain sound theoretical guarantees on the resulting system. This research vision arises from the strong belief that distributed computing is now mature enough to resolve the tension between the social needs for distributed computing systems, and the lack of a fundamentally sound and systematic way to realize these systems.

To progress on this vision, we plan in the near future to investigate, from a distributed software point of view, the impact due to failures and asynchrony on the layered architecture of distributed computing systems. A first step in this direction will address the notions of *message adversaries* (introduced a long time ago in [70]) and *process adversaries* (investigated in several papers, e.g. [69], [56], [62], [63], [67]). The aim of these notions is to consider failures, not as “bad events”, but as part of the normal behavior of a system. As an example, when considering round-based algorithms, a message adversary is a daemon which, at every round, is allowed to suppress some messages. The aim is then, given a problem  $P$ , to find the strongest adversary under which  $P$

can be solved (“strongest” means here that giving more power to the adversary makes the problem impossible to solve). This work will allow us to progress in terms of general *layered* theory of distributed computing, and allow us to better *map* distributed computing models and their relations, in the steps of noticeable early efforts in this direction [69], [43].



## ALPINES Project-Team

### 3. Research Program

#### 3.1. Overview

The research described here is directly relevant to several steps of the numerical simulation chain. Given a numerical simulation that was expressed as a set of differential equations, our research focuses on mesh generation methods for parallel computation, novel numerical algorithms for linear algebra, as well as algorithms and tools for their efficient and scalable implementation on high performance computers. The validation and the exploitation of the results is performed with collaborators from applications and is based on the usage of existing tools. In summary, the topics studied in our group are the following:

- Numerical methods and algorithms
  - Mesh generation for parallel computation
  - Solvers for numerical linear algebra
    - \* Domain decomposition methods
    - \* Preconditioning for iterative methods
  - Computational kernels for numerical linear algebra
  - Tensor computations
- Validation on numerical simulations and other numerical applications

#### 3.2. Domain specific language - parallel FreeFem++

In the engineering, researchers, and teachers communities, there is a strong demand for simulation frameworks that are simple to install and use, efficient, sustainable, and that solve efficiently and accurately complex problems for which there are no dedicated tools or codes available. In our group we develop FreeFem++ (see <https://www.freefem.org/>), a user dedicated language for solving PDEs. The goal of FreeFem++ is not to be a substitute for complex numerical codes, but rather to provide an efficient and relatively generic tool for:

- getting a quick answer to a specific problem,
- prototyping the resolution of a new complex problem.

The current users of FreeFem++ are mathematicians, engineers, university professors, and students. In general for these users the installation of public libraries as MPI, MUMPS, Ipopt, Blas, lapack, OpenGL, fftw, scotch, PETSc, SLEPc is a very difficult problem. For this reason, the authors of FreeFem++ have created a user friendly language, and over years have enriched its capabilities and provided tools for compiling FreeFem++ such that the users do not need to have special knowledge of computer science. This leads to an important work on porting the software on different emerging architectures.

Today, the main components of parallel FreeFem++ are:

1. definition of a coarse grid,
2. splitting of the coarse grid,
3. mesh generation of all subdomains of the coarse grid, and construction of parallel data structures for vectors and sparse matrices from the mesh of the subdomain,
4. call to a linear solver,
5. analysis of the result.

All these components are parallel, except for point (5) which is not in the focus of our research. However for the moment, the parallel mesh generation algorithm is very simple and not sufficient, for example it addresses only polygonal geometries. Having a better parallel mesh generation algorithm is one of the goals of our project. In addition, in the current version of FreeFem++, the parallelism is not hidden from the user, it is done through direct calls to MPI. Our goal is also to hide all the MPI calls in the specific language part of FreeFem++. In addition to these in-house domain decomposition methods, FreeFem++ is also linked to PETSc solvers which enables an easy use of third parties parallel multigrid methods.

### **3.3. Solvers for numerical linear algebra**

Iterative methods are widely used in industrial applications, and preconditioning is the most important research subject here. Our research considers domain decomposition methods and iterative methods and its goal is to develop solvers that are suitable for parallelism and that exploit the fact that the matrices are arising from the discretization of a system of PDEs on unstructured grids.

One of the main challenges that we address is the lack of robustness and scalability of existing methods as incomplete LU factorizations or Schwarz-based approaches, for which the number of iterations increases significantly with the problem size or with the number of processors. This is often due to the presence of several low frequency modes that hinder the convergence of the iterative method. To address this problem, we study different approaches for dealing with the low frequency modes as coarse space correction in domain decomposition or deflation techniques.

We also focus on developing boundary integral equation methods that would be adapted to the simulation of wave propagation in complex physical situations, and that would lend themselves to the use of parallel architectures. The final objective is to bring the state of the art on boundary integral equations closer to contemporary industrial needs. From this perspective, we investigate domain decomposition strategies in conjunction with boundary element method as well as acceleration techniques (H-matrices, FMM and the like) that would appear relevant in multi-material and/or multi-domain configurations. Our work on this topic also includes numerical implementation on large scale problems, which appears as a challenge due to the peculiarities of boundary integral equations.

### **3.4. Computational kernels for numerical linear and multilinear algebra**

The design of new numerical methods that are robust and that have well proven convergence properties is one of the challenges addressed in Alpines. Another important challenge is the design of parallel algorithms for the novel numerical methods and the underlying building blocks from numerical linear algebra. The goal is to enable their efficient execution on a diverse set of node architectures and their scaling to emerging high-performance clusters with an increasing number of nodes.

Increased communication cost is one of the main challenges in high performance computing that we address in our research by investigating algorithms that minimize communication, as communication avoiding algorithms. We propose to integrate the minimization of communication into the algorithmic design of numerical linear algebra problems. This is different from previous approaches where the communication problem was addressed as a scheduling or as a tuning problem. The communication avoiding algorithmic design is an approach originally developed in our group since 2007 (initially in collaboration with researchers from UC Berkeley and CU Denver). While at mid term we focus on reducing communication in numerical linear algebra, at long term we aim at considering the communication problem one level higher, during the parallel mesh generation tool described earlier.

Our research also focuses on solving problems of large size that feature high dimensions as in molecular simulations. The data in this case is represented by objects called tensors, or multilinear arrays. The goal is to design novel tensor techniques to allow their effective compression, i.e. their representation by simpler objects in small dimensions, while controlling the loss of information. The algorithms are aiming to be highly parallel to allow to deal with the large number of dimensions and large data sets, while preserving the required information for obtaining the solution of the problem.

## **AVALON Project-Team**

### **3. Research Program**

#### **3.1. Energy Application Profiling and Modeling**

Despite recent improvements, there is still a long road to follow in order to obtain energy efficient, energy proportional and eco-responsible exascale systems by 2022. Energy efficiency is therefore a major challenge for building next generation large-scale platforms. The targeted platforms will gather hundreds of millions of cores, low power servers, or CPUs. Besides being very important, their power consumption will be dynamic and irregular.

Thus, to consume energy efficiently, we aim at investigating two research directions. First, we need to improve measurement, understanding, and analysis on how large-scale platforms consume energy. Unlike some approaches [24] that mix the usage of internal and external wattmeters on a small set of resources, we target high frequency and precise internal and external energy measurements of each physical and virtual resource on large-scale distributed systems.

Secondly, we need to find new mechanisms that consume less and better on such platforms. Combined with hardware optimizations, several works based on shutdown or slowdown approaches aim at reducing energy consumption of distributed platforms and applications. To consume less, we first plan to explore the provision of accurate estimation of the energy consumed by applications without pre-executing and knowing them while most of the works try to do it based on in-depth application knowledge (code instrumentation [27], phase detection for specific HPC applications [31], *etc.* ). As a second step, we aim at designing a framework model that allows interaction, dialogue and decisions taken in cooperation among the user/application, the administrator, the resource manager, and the energy supplier. While smart grid is one of the last killer scenarios for networks, electrical provisioning of next generation large IT infrastructures remains a challenge.

#### **3.2. Data-intensive Application Profiling, Modeling, and Management**

Recently, the term “Big Data” has emerged to design data sets or collections so large that they become intractable for classical tools. This term is most time implicitly linked to “analytics” to refer to issues such as data curation, storage, search, sharing, analysis, and visualization. However, the Big Data challenge is not limited to data-analytics, a field that is well covered by programming languages and run-time systems such as Map-Reduce. It also encompasses data-intensive applications. These applications can be sorted into two categories. In High Performance Computing (HPC), data-intensive applications leverage post-petascale infrastructures to perform highly parallel computations on large amount of data, while in High Throughput Computing (HTC), a large amount of independent and sequential computations are performed on huge data collections.

These two types of data-intensive applications (HTC and HPC) raise challenges related to profiling and modeling that the AVALON team proposes to address. While the characteristics of data-intensive applications are very different, our work will remain coherent and focused. Indeed, a common goal will be to acquire a better understanding of both the applications and the underlying infrastructures running them to propose the best match between application requirements and infrastructure capacities. To achieve this objective, we will extensively rely on logging and profiling in order to design sound, accurate, and validated models. Then, the proposed models will be integrated and consolidated within a single simulation framework (SIMGRID). This will allow us to explore various potential “what-if?” scenarios and offer objective indicators to select interesting infrastructure configurations that match application specificities.

Another challenge is the ability to mix several heterogeneous infrastructures that scientists have at their disposal (*e.g.*, Grids, Clouds, and Desktop Grids) to execute data-intensive applications. Leveraging the aforementioned results, we will design strategies for efficient data management service for hybrid computing infrastructures.

### 3.3. Resource-Agnostic Application Description Model

With parallel programming, users expect to obtain performance improvement, regardless its cost. For long, parallel machines have been simple enough to let a user program use them given a minimal abstraction of their hardware. For example, MPI [26] exposes the number of nodes but hides the complexity of network topology behind a set of collective operations; OpenMP [30] simplifies the management of threads on top of a shared memory machine while OpenACC [29] aims at simplifying the use of GPGPU.

However, machines and applications are getting more and more complex so that the cost of manually handling an application is becoming very high [25]. Hardware complexity also stems from the unclear path towards next generations of hardware coming from the frequency wall: multi-core CPU, many-core CPU, GPGPUs, deep memory hierarchy, *etc.* have a strong impact on parallel algorithms. Parallel languages (UPC, Fortress, X10, *etc.*) is a first piece of the solution. However, they will still face the challenge of supporting distinct codes corresponding to different algorithms corresponding to distinct hardware capacities.

Therefore, the challenge we aim to address is to define a model, for describing the structure of parallel and distributed applications that enables code variations but also efficient executions on parallel and distributed infrastructures. Indeed, this issue appears for HPC applications but also for cloud oriented applications. The challenge is to adapt an application to user constraints such as performance, energy, security, *etc.*

Our approach is to consider component based models [32] as they offer the ability to manipulate the software architecture of an application. To achieve our goal, we consider a “compilation” approach that transforms a resource-agnostic application description into a resource-specific description. The challenge is thus to determine a component based model that enables to efficiently compute application mapping while being tractable. In particular, it has to provide an efficient support with respect to application and resource elasticity, energy consumption and data management. OpenMP runtime is a specific use case that we target.

### 3.4. Application Mapping and Scheduling

This research axis is at the crossroad of the AVALON team. In particular, it gathers results of the three other research axis. We plan to consider application mapping and scheduling addressing the following three issues.

#### 3.4.1. Application Mapping and Software Deployment

Application mapping and software deployment consist in the process of assigning distributed pieces of software to a set of resources. Resources can be selected according to different criteria such as performance, cost, energy consumption, security management, *etc.* A first issue is to select resources at application launch time. With the wide adoption of elastic platforms, *i.e.*, platforms that let the number of resources allocated to an application to be increased or decreased during its execution, the issue is also to handle resource selection at runtime.

The challenge in this context corresponds to the mapping of applications onto distributed resources. It will consist in designing algorithms that in particular take into consideration application profiling, modeling, and description.

A particular facet of this challenge is to propose scheduling algorithms for dynamic and elastic platforms. As the number of elements can vary, some kind of control of the platforms must be used accordingly to the scheduling.

#### 3.4.2. Non-Deterministic Workflow Scheduling

Many scientific applications are described through workflow structures. Due to the increasing level of parallelism offered by modern computing infrastructures, workflow applications now have to be composed not only of sequential programs, but also of parallel ones. New applications are now built upon workflows with conditionals and loops (also called non-deterministic workflows).

These workflows cannot be scheduled beforehand. Moreover cloud platforms bring on-demand resource provisioning and pay-as-you-go billing models. Therefore, there is a problem of resource allocation for non-deterministic workflows under budget constraints and using such an elastic management of resources.

Another important issue is data management. We need to schedule the data movements and replications while taking job scheduling into account. If possible, data management and job scheduling should be done at the same time in a closely coupled interaction.

### ***3.4.3. Software Asset Management***

The use of software is generally regulated by licenses, whether they are free or paid and with or without access to their sources. The world of licenses is very vast and unknown (especially in the industrial world). Often only the general public version is known (a software purchase corresponds to a license). For enterprises, the reality is much more complex, especially for main publishers. We work on the OpTISAM software, a software offering tools to perform Software Asset Management (SAM) much more efficiently in order to be able to ensure the full compliance with all contracts from each software and a new type of deployment taking into account these aspects and other additional parameters like energy and performance. This work is built on an Orange™ collaboration.

### ***3.4.4. Cloud deployment and reproducibility***

As part of the scientific method, any researcher should be able to reproduce the experimentation in order to not only verify the result but also evaluate and compare this experimentation with other approaches. The need of a standard tool allowing researchers to easily generate, share and reproduce experiments set-up arises. In our research, through a Nokia collaboration, we created SeeDep [10], a framework aiming at being such a standard tool. By associating a generation key to a network experiment set-up, SeeDep allows for reproducing network experiments independently from the used infrastructure.

## DATAMOVE Project-Team

### 3. Research Program

#### 3.1. Motivation

Today's largest supercomputers<sup>0</sup> are composed of few millions of cores, with performances almost reaching 100 PetaFlops<sup>0</sup> for the largest machine. Moving data in such large supercomputers is becoming a major performance bottleneck, and the situation is expected to worsen even more at exascale and beyond. The data transfer capabilities are growing at a slower rate than processing power ones. The profusion of available flops will very likely be underused due to constrained communication capabilities. It is commonly admitted that data movements account for 50% to 70% of the global power consumption<sup>0</sup>. Thus, data movements are potentially one of the most important source of savings for enabling supercomputers to stay in the commonly adopted energy barrier of 20 MegaWatts. In the mid to long term, non volatile memory (NVRAM) is expected to deeply change the machine I/Os. Data distribution will shift from disk arrays with an access time often considered as uniform, towards permanent storage capabilities at each node of the machine, making data locality an even more prevalent paradigm.

The proposed DataMove team will work on **optimizing data movements for large scale computing** mainly at two related levels:

- Resource allocation
- Integration of numerical simulation and data analysis

The resource and job management system (also called batch scheduler or RJMS) is in charge of allocating resources upon user requests for executing their parallel applications. The growing cost of data movements requires adapted scheduling policies able to take into account the influence of intra-application communications, I/Os as well as contention caused by data traffic generated by other concurrent applications. Modelling the application behavior to anticipate its actual resource usage on such architecture is known to be challenging, but it becomes critical for improving performances (execution time, energy, or any other relevant objective). The job management system also needs to handle new types of workloads: high performance platforms now need to execute more and more often data intensive processing tasks like data analysis in addition to traditional computation intensive numerical simulations. In particular, the ever growing amount of data generated by numerical simulation calls for a tighter integration between the simulation and the data analysis. The challenge here is to reduce data traffic and to speed-up result analysis by performing result processing (compression, indexation, analysis, visualization, etc.) as closely as possible to the locus and time of data generation. This emerging trend called *in-situ analytics* requires to revisit the traditional workflow (loop of batch processing followed by postmortem analysis). The application becomes a whole including the simulation, in-situ processing and I/Os. This motivates the development of new well-adapted resource sharing strategies, data structures and parallel analytics schemes to efficiently interleave the different components of the application and globally improve the performance.

#### 3.2. Strategy

DataMove targets HPC (High Performance Computing) at Exascale. But such machines and the associated applications are expected to be available only in 5 to 10 years. Meanwhile, we expect to see a growing number of petaflop machines to answer the needs for advanced numerical simulations. A sustainable exploitation of these petaflop machines is a real and hard challenge that we will address. We may also see in the coming years a convergence between HPC and Big Data, HPC platforms becoming more elastic and supporting Big

<sup>0</sup>Top500 Ranking, <http://www.top500.org>

<sup>0</sup>10<sup>15</sup> floating point operations per second

<sup>0</sup>SciDAC Review, 2010, <http://scidacreview.org/1001/pdf/hardware.pdf>

Data jobs, or HPC applications being more commonly executed on cloud like architectures. This is the second top objective of the 2015 US Strategic Computing Initiative <sup>0</sup>: *Increasing coherence between the technology base used for modelling and simulation and that used for data analytic computing*. We will contribute to that convergence at our level, considering more dynamic and versatile target platforms and types of workloads.

Our approaches should entail minimal modifications on the code of numerical simulations. Often large scale numerical simulations are complex domain specific codes with a long life span. We assume these codes as being sufficiently optimized. We will influence the behavior of numerical simulations through resource allocation at the job management system level or when interleaving them with analytics code.

To tackle these issues, we propose to intertwine theoretical research and practical developments in an agile mode. Algorithms with performance guarantees will be designed and experimented on large scale platforms with realistic usage scenarios developed with partner scientists or based on logs of the biggest available computing platforms (national supercomputers like Curie, or the BlueWaters machine accessible through our collaboration with Argonne National Lab). Conversely, a strong experimental expertise will enable to feed theoretical models with sound hypotheses, to twist proven algorithms with practical heuristics that could be further retro-fed into adequate theoretical models.

A central scientific question is to make the relevant choices for optimizing performance (in a broad sense) in a reasonable time. HPC architectures and applications are increasingly complex systems (heterogeneity, dynamicity, uncertainties), which leads to consider the **optimization of resource allocation based on multiple objectives**, often contradictory (like energy and run-time for instance). Focusing on the optimization of one particular objective usually leads to worsen the others. The historical positioning of some members of the team who are specialists in multi-objective optimization is to generate a (limited) set of trade-off configurations, called *Pareto points*, and choose when required the most suitable trade-off between all the objectives. This methodology differs from the classical approaches, which simplify the problem into a single objective one (focus on a particular objective, combining the various objectives or agglomerate them). The real challenge is thus to combine algorithmic techniques to account for this diversity while guaranteeing a target efficiency for all the various objectives.

The DataMove team aims to elaborate generic and effective solutions of practical interest. We will make our new algorithms accessible through the team flagship software tools, **the OAR batch scheduler and the in-situ processing framework FlowVR**. We will maintain and enforce strong links with teams closely connected with large architecture design and operation (CEA DAM, BULL, Argonne National Lab), as well as scientists of other disciplines, in particular computational biologists, with whom we will elaborate and validate new usage scenarios (IBPC, CEA DAM, EDF).

### 3.3. Research Directions

DataMove research activity is organised around three directions. When a parallel job executes on a machine, it triggers data movements through the input data it needs to read, the results it produces (simulation results as well as traces) that need to be stored in the file system, as well as internal communications and temporary storage (for fault tolerance related data for instance). Modeling in details the simulation and the target machines to analyze scheduling policies is not feasible at large scales. We propose to investigate alternative approaches, including learning approaches, to capture and model the influence of data movements on the performance metrics of each job execution to develop **Data Aware Batch Scheduling** models and algorithms (Sec. 4.1 ). Experimenting new scheduling policies on real platforms at scale is unfeasible. Theoretical performance guarantees are not sufficient to ensure a new algorithm will actually perform as expected on a real platform. An intermediate evaluation level is required to probe novel scheduling policies. The second research axe focuses on the **Empirical Studies of Large Scale Platforms** (Sec. 4.2 ). The goal is to investigate how we could extract from actual computing centers traces information to replay the job allocations and executions on a simulated or emulated platform with new scheduling policies. Schedulers need information about jobs behavior on target machines to actually be able to make efficient allocation decisions. Asking users

<sup>0</sup><https://www.whitehouse.gov/the-press-office/2015/07/29/executive-order-creating-national-strategic-computing-initiative>

to characterize jobs often does not lead to reliable information. The third research direction **Integration of High Performance Computing and Data Analytics** (Sec. 4.3 ) addresses the data movement issue from a different perspective. New data analysis techniques on the HPC platform introduce new type of workloads, potentially more data than compute intensive, but could also enable to reduce data movements by directly enabling to pipe-line simulation execution with a live analysis of the produced results. Our goal is here to investigate how to program and schedule such analysis workflows in the HPC context.



## HIEPACS Project-Team

### 3. Research Program

#### 3.1. Introduction

The methodological component of **HIEPACS** 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 **HIEPACS** 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 preconditioned 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 developed 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 breakdown detection, deflation at restart and subspace recycling. Beyond new numerical investigations, a software implementation to be included in our linear solver library **Fabulous** originally 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 through 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 conducted 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 conducted on close connection with the activity described in Section 3.3. Following our earlier experience, we plan to first consider approximate inverse preconditioners 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 use 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.

## KERDATA Project-Team

### 3. Research Program

#### 3.1. Research axis 1: Convergence of HPC and Big Data

The tools and cultures of High Performance Computing and Big Data Analytics have evolved in divergent ways. This is to the detriment of both. However, big computations still generate and are needed to analyze Big Data. As scientific research increasingly depends on both high-speed computing and data analytics, the potential interoperability and scaling convergence of these two ecosystems is crucial to the future.

Our objective is premised on the idea that we must explore the ways in which the major challenges associated with Big Data analytics intersect with, impact, and potentially change the directions now in progress for achieving Exascale computing.

In particular, a key milestone will be to achieve convergence through common abstractions and techniques for data storage and processing in support of complex workflows combining simulations and analytics. Such application workflows will need such a convergence to run on hybrid infrastructures combining HPC systems and clouds (potentially in extension to edge devices, in a complete digital continuum).

**Collaboration.** *This axis is addressed in close collaboration with [María Pérez](#) (UPM), [Rob Ross](#) (ANL), [Toni Cortes](#) (BSC), Several groups at Argonne National Laboratory and NCSA ([Franck Cappello](#), [Rob Ross](#), [Bill Kramer](#), [Tom Peterka](#)).*

*Relevant groups with similar interests are the following ones.*

- *The group of [Jack Dongarra](#), Innovative Computing Laboratory at University of Tennessee, who is leading international efforts for the convergence of Exascale Computing and Big Data.*
- *The group of [Satoshi Matsuoka](#), RIKEN, working on system software for clouds and HPC.*
- *The group of [Ian Foster](#), Argonne National Laboratory, working on on-demand data analytics and storage for extreme-scale simulations and experiments.*

##### 3.1.1. High-performance storage for concurrent Big Data applications

Storage is a plausible pathway to convergence. In this context, we plan to focus on the needs of concurrent Big Data applications that require high-performance storage, as well as transaction support. Although blobs (binary large objects) are an increasingly popular storage model for such applications, state-of-the-art blob storage systems offer no transaction semantics. This demands users to coordinate data access carefully in order to avoid race conditions, inconsistent writes, overwrites and other problems that cause erratic behavior.

There is a gap between existing storage solutions and application requirements, which limits the design of transaction-oriented applications. In this context, one idea on which we plan to focus our efforts is exploring how blob storage systems could provide built-in, multiblob transactions, while retaining sequential consistency and high throughput under heavy access concurrency.

The early principles of this research direction have already raised interest from our partners at ANL (Rob Ross) and UPM (María Pérez) for potential collaborations. In this direction, the acceptance of our paper on the Týr transactional blob storage system as a Best Student Paper Award Finalist at the SC16 conference [10] is a very encouraging step.

### 3.1.2. Towards unified data processing techniques for Extreme Computing and Big Data applications

In the high-performance computing area (HPC), the need to get fast and relevant insights from massive amounts of data generated by extreme-scale computations led to the emergence of *in situ processing*. It allows data to be visualized and processed in real-time on the supercomputer generating them, in an interactive way, as they are produced, as opposed to the traditional approach consisting of transferring data off-site after the end of the computation, for offline analysis. As such processing runs on the same resources executing the simulation, if it consumes too many resources, there is a risk to "disturb" the simulation.

Consequently, an alternative approach was proposed (*in transit processing*), as a means to reduce this impact: data are transferred to some temporary processing resources (with high memory and processing capacities). After this real-time processing, they are moved to persistent storage.

In the Big Data area, the search for real-time, fast analysis was materialized through a different approach: stream-based processing. Such an approach is based on a different abstraction for data, that are seen as a dynamic flow of items to be processed. Stream-based processing and *in situ/in transit processing* have been developed separately and implemented in different tools in the BDA and HPC areas respectively.

A major challenge from the perspective of the HPC-BDA convergence is their joint use in a unified data processing architecture. This is one of the future research challenges that I plan to address in the near future, by combining ongoing approaches currently active in my team: Damaris and KerA. We started preliminary work within the "*Frameworks*" work package of the HPC-Big Data IPL. Further exploring this convergence is a core direction of our current efforts to build collaborative European projects.

## 3.2. Research axis 2: Cloud and Edge processing

The recent evolutions in the area of Big Data processing have pointed out some limitations of the initial Map-Reduce model. It is well suited for batch data processing, but less suited for real-time processing of dynamic data streams. New types of data-intensive applications emerge, e.g., for enterprises who need to perform analysis on their stream data in ways that can give fast results (i.e., in real time) at scale (e.g., click-stream analysis and network-monitoring log analysis). Similarly, scientists require fast and accurate data processing techniques in order to analyze their experimental data correctly at scale (e.g., collectively analysis of large data sets distributed in multiple geographically distributed locations).

Our plan is to revisit current data storage and processing techniques to cope with the volatile requirements of data-intensive applications on large-scale dynamic clouds in a cost-efficient way, with a particular focus on streaming. More recently, the strong emergence of edge/fog-based infrastructures leads to additional challenges for new scenarios involving hybrid cloud/fog/edge systems.

**Collaboration.** This axis is addressed in close collaboration with *María Pérez (UPM)*, *Kate Keahey (ANL)*

Relevant groups with similar interests include the following ones.

- The group of *Geoffrey Fox*, Indiana University, working on data analytics, cloud data processing, stream processing.
- The group at RISE Lab, UC Berkeley, working on real-time stream-based processing and analytics.
- The group of *Ewa Deelman*, USC Information Sciences Institute, working on resource management for workflows in clouds.

### 3.2.1. Stream-oriented, Big Data processing on clouds

The state-of-the-art Hadoop Map-Reduce framework cannot deal with stream data applications, as it requires the data to be initially stored in a distributed file system in order to process them. To better cope with the above-mentioned requirements, several systems have been introduced for stream data processing such as Flink [27], Spark [32], Storm [33], and Google MillWheel [34]. These systems keep computation in memory to decrease

latency, and preserve scalability by using data-partitioning or dividing the streams into a set of deterministic batch computations.

However, they are designed to work in dedicated environments and they do not consider the performance variability (i.e., network, I/O, etc.) caused by resource contention in the cloud. This variability may in turn cause high and unpredictable latency when output streams are transmitted to further analysis. Moreover, they overlook the dynamic nature of data streams and the volatility in their computation requirements. Finally, they still address failures in a best-effort manner.

Our objective is to investigate new approaches for reliable, stream Big Data processing on clouds.

### ***3.2.2. Efficient Edge, Cloud and hybrid Edge/Cloud data processing***

Today, we are approaching an important technological milestone: applications are generating huge amounts of data and are demanding low-latency responses to their requests. Mobile computing and Internet of Things (IoT) applications are good illustrations of such scenarios. Using only Cloud computing for such scenarios is challenging. Firstly, Cloud resources are most of the time accessed through Internet, hence, data are sent across high-latency wide area networks, which may degrade the performance of applications. Secondly, it may be impossible to send data to the Cloud due to data regulations, national security laws or simply because an Internet connection is not available. Finally, data transmission costs (e.g., Cloud provider fees, carrier costs) could make a business solution impractical.

Edge computing is a new paradigm which aims to address some of these issues. The key idea is to leverage computing and storage resources at the "edge" of the network, i.e., on processing units located close to the data sources. This allows applications to outsource task execution from the main (Cloud) processing data centers to the edge. The development of Edge computing was accelerated by the recent emergence of stream processing, a new model for handling continuous flows of data in real-time, as opposed to batch processing, which typically processes bounded datasets offline.

However, Edge computing is not a silver bullet. Besides being a new concept not fully established in the community, issues like node volatility, limited processing power, high latency between nodes, fault tolerance and data degradation may impact applications depending on the characteristics of the infrastructure.

Some relevant research questions are: How much can one improve (or degrade) the performance of an application by performing data processing closer to the data sources rather than performing it in the cloud? How to progress towards a seamless scheduling and execution of a data analytics workflow and break the limitation the current dual approaches used in preliminary efforts in this area, that rely on manual and empirical deployment of the corresponding dataflow operator graphs, using separate analytics engines for centralized clouds and for edge systems respectively?

Our objective is to try to answer precisely such questions. We are interested in understanding the conditions that enable the usage of Edge or Cloud computing to reduce the time to results and the associated costs. While some state-of-the-art approaches advocate either "100% Cloud" or "100% Edge" solutions, the relative efficiency of a method over the other may vary. Intuitively, it depends on many parameters, including network technology, hardware characteristics, volume of data or computing power, processing framework configuration and application requirements, to cite a few. We plan to study their impact on the overall application performance.

## **3.3. Research axis 3: Supporting AI across the digital continuum**

Integrating and processing high-frequency data streams from multiple sensors scattered over a large territory in a timely manner requires high-performance computing techniques and equipments. For instance, a machine learning earthquake detection solution has to be designed jointly with experts in distributed computing and cyber-infrastructure to enable real-time alerts. Because of the large number of sensors and their high sampling rate, a traditional centralized approach which transfers all data to a single point may be impractical. Our goal is to investigate innovative solutions for the design of efficient data processing infrastructures for a distributed machine learning-based approach.

In particular, building on our previous results in the area of efficient stream processing systems, we aim to explore approaches for unified data storage, processing and machine-learning based analytics across the whole digital continuum (i.e., for highly distributed applications deployed on hybrid edge/cloud/HPC infrastructures). Our ZettaFlow project is targeting a startup creation precisely this area.

**Collaboration.** *This recently started axis is worked out in close collaboration with the group of [Manish Parashar](#), Rutgers University, and with the [LACODAM](#) team at Inria, focused on large-scale collaborative data mining.*

## POLARIS Project-Team

### 3. Research Program

#### 3.1. Sound and Reproducible Experimental Methodology

**Participants:** Vincent Danjean, Nicolas Gast, Guillaume Huard, Arnaud Legrand, Patrick Loiseau, Jean-Marc Vincent.

Experiments in large scale distributed systems are costly, difficult to control and therefore difficult to reproduce. Although many of these digital systems have been built by men, they have reached such a complexity level that we are no longer able to study them like artificial systems and have to deal with the same kind of experimental issues as natural sciences. The development of a sound experimental methodology for the evaluation of resource management solutions is among the most important ways to cope with the growing complexity of computing environments. Although computing environments come with their own specific challenges, we believe such general observation problems should be addressed by borrowing good practices and techniques developed in many other domains of science.

This research theme builds on a transverse activity on *Open science and reproducible research* and is organized into the following two directions: (1) *Experimental design* (2) *Smart monitoring and tracing*. As we will explain in more detail hereafter, these transverse activity and research directions span several research areas and our goal within the POLARIS project is foremost to transfer original ideas from other domains of science to the distributed and high performance computing community.

#### 3.2. Multi-Scale Analysis and Visualization

**Participants:** Vincent Danjean, Guillaume Huard, Arnaud Legrand, Jean-Marc Vincent, Panayotis Mertikopoulos.

As explained in the previous section, the first difficulty encountered when modeling large scale computer systems is to observe these systems and extract information on the behavior of both the architecture, the middleware, the applications, and the users. The second difficulty is to *visualize* and *analyze* such *multi-level traces to understand how the performance of the application can be improved*. While a lot of efforts are put into visualizing scientific data, in comparison little effort have gone into to developing techniques specifically tailored for understanding the behavior of distributed systems. Many visualization tools have been developed by renowned HPC groups since decades (e.g., BSC [91], Jülich and TU Dresden [90], [61], UIUC [79], [94], [82] and ANL [107], Inria Bordeaux [67] and Grenoble [109], ...) but most of these tools build on the classical information visualization mantra [99] that consists in always first presenting an overview of the data, possibly by plotting everything if computing power allows, and then to allow users to zoom and filter, providing details on demand. However in our context, the amount of data comprised in such traces is several orders of magnitude larger than the number of pixels on a screen and displaying even a small fraction of the trace leads to harmful visualization artifacts [86]. Such traces are typically made of events that occur at very different time and space scales, which unfortunately hinders classical approaches. Such visualization tools have focused on easing interaction and navigation in the trace (through gantcharts, intuitive filters, pie charts and kiviats) but they are very difficult to maintain and evolve and they require some significant experience to identify performance bottlenecks.

Therefore many groups have more recently proposed in combination to these tools some techniques to help identifying the structure of the application or regions (applicative, spatial or temporal) of interest. For example, researchers from the SDSC [89] propose some segment matching techniques based on clustering (Euclidean or Manhattan distance) of start and end dates of the segments that enables to reduce the amount of information to display. Researchers from the BSC use clustering, linear regression and Kriging techniques [98], [85], [78] to identify and characterize (in term of performance and resource usage) application phases and present aggregated representations of the trace [97]. Researchers from Jülich and TU Darmstadt have proposed techniques to identify specific communication patterns that incur wait states [104], [54]

### 3.3. Fast and Faithful Performance Prediction of Very Large Systems

**Participants:** Jonatha Anselmi, Vincent Danjean, Bruno Gaujal, Arnaud Legrand, Florence Perronnin, Jean-Marc Vincent.

Evaluating the scalability, robustness, energy consumption and performance of large infrastructures such as exascale platforms and clouds raises severe methodological challenges. The complexity of such platforms mandates empirical evaluation but direct experimentation via an application deployment on a real-world testbed is often limited by the few platforms available at hand and is even sometimes impossible (cost, access, early stages of the infrastructure design, ...). Unlike direct experimentation via an application deployment on a real-world testbed, simulation enables fully repeatable and configurable experiments that can often be conducted quickly for arbitrary hypothetical scenarios. In spite of these promises, current simulation practice is often not conducive to obtaining scientifically sound results. To date, most simulation results in the parallel and distributed computing literature are obtained with simulators that are ad hoc, unavailable, undocumented, and/or no longer maintained. For instance, Naicken et al. [53] point out that out of 125 recent papers they surveyed that study peer-to-peer systems, 52% use simulation and mention a simulator, but 72% of them use a custom simulator. As a result, most published simulation results build on throw-away (short-lived and non validated) simulators that are specifically designed for a particular study, which prevents other researchers from building upon it. There is thus a strong need for recognized simulation frameworks by which simulation results can be reproduced, further analyzed and improved.

The *SimGrid* simulation toolkit [65], whose development is partially supported by POLARIS, is specifically designed for studying large scale distributed computing systems. It has already been successfully used for simulation of grid, volunteer computing, HPC, cloud infrastructures and we have constantly invested on the software quality, the scalability [57] and the validity of the underlying network models [55], [102]. Many simulators of MPI applications have been developed by renowned HPC groups (e.g., at SDSC [100], BSC [51], UIUC [108], Sandia Nat. Lab. [103], ORNL [64] or ETH Zürich [80] for the most prominent ones). Yet, to scale most of them build on restrictive network and application modeling assumptions that make them difficult to extend to more complex architectures and to applications that do not solely build on the MPI API. Furthermore, simplistic modeling assumptions generally prevent to faithfully predict execution times, which limits the use of simulation to indication of gross trends at best. Our goal is to improve the quality of SimGrid to the point where it can be used effectively on a daily basis by practitioners to *reproduce the dynamic of real HPC systems*.

We also develop another simulation software, *PSI* (Perfect SIMulator) [69], [62], dedicated to the simulation of very large systems that can be modeled as Markov chains. PSI provides a set of simulation kernels for Markov chains specified by events. It allows one to sample stationary distributions through the Perfect Sampling method (pioneered by Propp and Wilson [92]) or simply to generate trajectories with a forward Monte-Carlo simulation leveraging time parallel simulation (pioneered by Fujimoto [73], Lin and Lazowska [84]). One of the strength of the PSI framework is its expressiveness that allows us to easily study networks with finite and infinite capacity queues [63]. Although PSI already allows to simulate very large and complex systems, our main objective is to push its scalability even further and *improve its capabilities by one or several orders of magnitude*.

### 3.4. Local Interactions and Transient Analysis in Adaptive Dynamic Systems

**Participants:** Jonatha Anselmi, Nicolas Gast, Bruno Gaujal, Florence Perronnin, Jean-Marc Vincent, Panayotis Mertikopoulos.

Many systems can be effectively described by stochastic population models. These systems are composed of a set of  $n$  entities interacting together and the resulting stochastic process can be seen as a continuous-time Markov chain with a finite state space. Many numerical techniques exist to study the behavior of Markov chains, to solve stochastic optimal control problems [93] or to perform model-checking [52]. These techniques, however, are limited in their applicability, as they suffer from the *curse of dimensionality*: the state-space grows exponentially with  $n$ .



This results in the need for approximation techniques. Mean field analysis offers a viable, and often very accurate, solution for large  $n$ . The basic idea of the mean field approximation is to count the number of entities that are in a given state. Hence, the fluctuations due to stochasticity become negligible as the number of entities grows. For large  $n$ , the system becomes essentially deterministic. This approximation has been originally developed in statistical mechanics for vary large systems composed of more than  $10^{20}$  particles (called entities here). More recently, it has been claimed that, under some conditions, this approximation can be successfully used for stochastic systems composed of a few tens of entities. The claim is supported by various convergence results [74], [83], [106], and has been successfully applied in various domains: wireless networks [56], computer-based systems [77], [88], [101], epidemic or rumour propagation [66], [81] and bike-sharing systems [70]. It is also used to develop distributed control strategies [105], [87] or to construct approximate solutions of stochastic model checking problems [58], [59], [60].

Within the POLARIS project, we will continue developing both the theory behind these approximation techniques and their applications. Typically, these techniques require a homogeneous population of objects where the dynamics of the entities depend only on their state (the state space of each object must not scale with  $n$  the number of objects) but neither on their identity nor on their spatial location. Continuing our work in [74], we would like to be able to handle heterogeneous or uncertain dynamics. Typical applications are caching mechanisms [77] or bike-sharing systems [71]. A second point of interest is the use of mean field or large deviation asymptotics to compute the time between two regimes [96] or to reach an equilibrium state. Last, mean-field methods are mostly descriptive and are used to analyse the performance of a given system. We wish to extend their use to solve optimal control problems. In particular, we would like to implement numerical algorithms that use the framework that we developed in [75] to build distributed control algorithms [68] and optimal pricing mechanisms [76].

### 3.5. Distributed Learning in Games and Online Optimization

**Participants:** Nicolas Gast, Bruno Gaujal, Arnaud Legrand, Patrick Loiseau, Panayotis Mertikopoulos, Bary Pradelski.

Game theory is a thriving interdisciplinary field that studies the interactions between competing optimizing agents, be they humans, firms, bacteria, or computers. As such, game-theoretic models have met with remarkable success when applied to complex systems consisting of interdependent components with vastly different (and often conflicting) objectives – ranging from latency minimization in packet-switched networks to throughput maximization and power control in mobile wireless networks.

In the context of large-scale, decentralized systems (the core focus of the POLARIS project), it is more relevant to take an inductive, “bottom-up” approach to game theory, because the components of a large system cannot be assumed to perform the numerical calculations required to solve a very-large-scale optimization problem. In view of this, POLARIS’ overarching objective in this area is to *develop novel algorithmic frameworks that offer robust performance guarantees when employed by all interacting decision-makers.*

A key challenge here is that most of the literature on learning in games has focused on *static* games with a *finite number of actions* per player [72], [95]. While relatively tractable, such games are ill-suited to practical applications where players pick an action from a continuous space or when their payoff functions evolve over time – this being typically the case in our target applications (e.g., routing in packet-switched networks or energy-efficient throughput maximization in wireless). On the other hand, the framework of online convex optimization typically provides worst-case performance bounds on the learner’s *regret* that the agents can attain irrespectively of how their environment varies over time. However, if the agents’ environment is determined chiefly by their interactions these bounds are fairly loose, so more sophisticated convergence criteria should be applied.

From an algorithmic standpoint, a further challenge occurs when players can only observe their own payoffs (or a perturbed version thereof). In this bandit-like setting regret-matching or trial-and-error procedures guarantee convergence to an equilibrium in a weak sense in certain classes of games. However, these results apply exclusively to static, finite games: learning in games with continuous action spaces and/or nonlinear

payoff functions cannot be studied within this framework. Furthermore, even in the case of finite games, the complexity of the algorithms described above is not known, so it is impossible to decide a priori which algorithmic scheme can be applied to which application.

## ROMA Project-Team

### 3. Research Program

#### 3.1. Algorithms for probabilistic environments

There are two main research directions under this research theme. In the first one, we consider the problem of the efficient execution of applications in a failure-prone environment. Here, probability distributions are used to describe the potential behavior of computing platforms, namely when hardware components are subject to faults. In the second research direction, probability distributions are used to describe the characteristics and behavior of applications.

##### 3.1.1. Application resilience

An application is resilient if it can successfully produce a correct result in spite of potential faults in the underlying system. Application resilience can involve a broad range of techniques, including fault prediction, error detection, error containment, error correction, checkpointing, replication, migration, recovery, etc. Faults are quite frequent in the most powerful existing supercomputers. The Jaguar platform, which ranked third in the TOP 500 list in November 2011 [44], had an average of 2.33 faults per day during the period from August 2008 to February 2010 [68]. The mean-time between faults of a platform is inversely proportional to its number of components. Progresses will certainly be made in the coming years with respect to the reliability of individual components. However, designing and building high-reliability hardware components is far more expensive than using lower reliability top-of-the-shelf components. Furthermore, low-power components may not be available with high-reliability. Therefore, it is feared that the progresses in reliability will far from compensate the steady projected increase of the number of components in the largest supercomputers. Already, application failures have a huge computational cost. In 2008, the DARPA white paper on “System resilience at extreme scale” [43] stated that high-end systems wasted 20% of their computing capacity on application failure and recovery.

In such a context, any application using a significant fraction of a supercomputer and running for a significant amount of time will have to use some fault-tolerance solution. It would indeed be unacceptable for an application failure to destroy centuries of CPU-time (some of the simulations run on the Blue Waters platform consumed more than 2,700 years of core computing time [39] and lasted over 60 hours; the most time-consuming simulations of the US Department of Energy (DoE) run for weeks to months on the most powerful existing platforms [42]).

Our research on resilience follows two different directions. On the one hand we design new resilience solutions, either generic fault-tolerance solutions or algorithm-based solutions. On the other hand we model and theoretically analyze the performance of existing and future solutions, in order to tune their usage and help determine which solution to use in which context.

##### 3.1.2. Scheduling strategies for applications with a probabilistic behavior

Static scheduling algorithms are algorithms where all decisions are taken before the start of the application execution. On the contrary, in non-static algorithms, decisions may depend on events that happen during the execution. Static scheduling algorithms are known to be superior to dynamic and system-oriented approaches in stable frameworks [50], [56], [57], [67], that is, when all characteristics of platforms and applications are perfectly known, known a priori, and do not evolve during the application execution. In practice, the prediction of application characteristics may be approximative or completely infeasible. For instance, the amount of computations and of communications required to solve a given problem in parallel may strongly depend on some input data that are hard to analyze (this is for instance the case when solving linear systems using full pivoting).

We plan to consider applications whose characteristics change dynamically and are subject to uncertainties. In order to benefit nonetheless from the power of static approaches, we plan to model application uncertainties and variations through probabilistic models, and to design for these applications scheduling strategies that are either static, or partially static and partially dynamic.

## 3.2. Platform-aware scheduling strategies

In this theme, we study and design scheduling strategies, focusing either on energy consumption or on memory behavior. In other words, when designing and evaluating these strategies, we do not limit our view to the most classical platform characteristics, that is, the computing speed of cores and accelerators, and the bandwidth of communication links.

In most existing studies, a single optimization objective is considered, and the target is some sort of absolute performance. For instance, most optimization problems aim at the minimization of the overall execution time of the application considered. Such an approach can lead to a very significant waste of resources, because it does not take into account any notion of efficiency nor of yield. For instance, it may not be meaningful to use twice as many resources just to decrease by 10% the execution time. In all our work, we plan to look only for algorithmic solutions that make a “clever” usage of resources. However, looking for the solution that optimizes a metric such as the efficiency, the energy consumption, or the memory-peak minimization, is doomed for the type of applications we consider. Indeed, in most cases, any optimal solution for such a metric is a sequential solution, and sequential solutions have prohibitive execution times. Therefore, it becomes mandatory to consider multi-criteria approaches where one looks for trade-offs between some user-oriented metrics that are typically related to notions of Quality of Service—execution time, response time, stretch, throughput, latency, reliability, etc.—and some system-oriented metrics that guarantee that resources are not wasted. In general, we will not look for the Pareto curve, that is, the set of all dominating solutions for the considered metrics. Instead, we will rather look for solutions that minimize some given objective while satisfying some bounds, or “budgets”, on all the other objectives.

### 3.2.1. Energy-aware algorithms

Energy-aware scheduling has proven an important issue in the past decade, both for economical and environmental reasons. Energy issues are obvious for battery-powered systems. They are now also important for traditional computer systems. Indeed, the design specifications of any new computing platform now always include an upper bound on energy consumption. Furthermore, the energy bill of a supercomputer may represent a significant share of its cost over its lifespan.

Technically, a processor running at speed  $s$  dissipates  $s^\alpha$  watts per unit of time with  $2 \leq \alpha \leq 3$  [48], [49], [54]; hence, it consumes  $s^\alpha \times d$  joules when operated during  $d$  units of time. Therefore, energy consumption can be reduced by using speed scaling techniques. However it was shown in [69] that reducing the speed of a processor increases the rate of transient faults in the system. The probability of faults increases exponentially, and this probability cannot be neglected in large-scale computing [65]. In order to make up for the loss in *reliability* due to the energy efficiency, different models have been proposed for fault tolerance: (i) *re-execution* consists in re-executing a task that does not meet the reliability constraint [69]; (ii) *replication* consists in executing the same task on several processors simultaneously, in order to meet the reliability constraints [47]; and (iii) *checkpointing* consists in “saving” the work done at some certain instants, hence reducing the amount of work lost when a failure occurs [64].

Energy issues must be taken into account at all levels, including the algorithm-design level. We plan to both evaluate the energy consumption of existing algorithms and to design new algorithms that minimize energy consumption using tools such as resource selection, dynamic frequency and voltage scaling, or powering-down of hardware components.

### 3.2.2. Memory-aware algorithms

For many years, the bandwidth between memories and processors has increased more slowly than the computing power of processors, and the latency of memory accesses has been improved at an even slower

pace. Therefore, in the time needed for a processor to perform a floating point operation, the amount of data transferred between the memory and the processor has been decreasing with each passing year. The risk is for an application to reach a point where the time needed to solve a problem is no longer dictated by the processor computing power but by the memory characteristics, comparable to the *memory wall* that limits CPU performance. In such a case, processors would be greatly under-utilized, and a large part of the computing power of the platform would be wasted. Moreover, with the advent of multicore processors, the amount of memory per core has started to stagnate, if not to decrease. This is especially harmful to memory intensive applications. The problems related to the sizes and the bandwidths of memories are further exacerbated on modern computing platforms because of their deep and highly heterogeneous hierarchies. Such a hierarchy can extend from core private caches to shared memory within a CPU, to disk storage and even tape-based storage systems, like in the Blue Waters supercomputer [40]. It may also be the case that heterogeneous cores are used (such as hybrid CPU and GPU computing), and that each of them has a limited memory.

Because of these trends, it is becoming more and more important to precisely take memory constraints into account when designing algorithms. One must not only take care of the amount of memory required to run an algorithm, but also of the way this memory is accessed. Indeed, in some cases, rather than to minimize the amount of memory required to solve the given problem, one will have to maximize data reuse and, especially, to minimize the amount of data transferred between the different levels of the memory hierarchy (minimization of the volume of memory inputs-outputs). This is, for instance, the case when a problem cannot be solved by just using the in-core memory and that any solution must be out-of-core, that is, must use disks as storage for temporary data.

It is worth noting that the cost of moving data has led to the development of so called “communication-avoiding algorithms” [60]. Our approach is orthogonal to these efforts: in communication-avoiding algorithms, the application is modified, in particular some redundant work is done, in order to get rid of some communication operations, whereas in our approach, we do not modify the application, which is provided as a task graph, but we minimize the needed memory peak only by carefully scheduling tasks.

### 3.3. High-performance computing and linear algebra

Our work on high-performance computing and linear algebra is organized along three research directions. The first direction is devoted to direct solvers of sparse linear systems. The second direction is devoted to combinatorial scientific computing, that is, the design of combinatorial algorithms and tools that solve problems encountered in some of the other research themes, like the problems faced in the preprocessing phases of sparse direct solvers. The last direction deals with the adaptation of classical dense linear algebra kernels to the architecture of future computing platforms.

#### 3.3.1. Direct solvers for sparse linear systems

The solution of sparse systems of linear equations (symmetric or unsymmetric, often with an irregular structure, from a few hundred thousand to a few hundred million equations) is at the heart of many scientific applications arising in domains such as geophysics, structural mechanics, chemistry, electromagnetism, numerical optimization, or computational fluid dynamics, to cite a few. The importance and diversity of applications are a main motivation to pursue research on sparse linear solvers. Because of this wide range of applications, any significant progress on solvers will have a significant impact in the world of simulation. Research on sparse direct solvers in general is very active for the following main reasons:

- many applications fields require large-scale simulations that are still too big or too complicated with respect to today’s solution methods;
- the current evolution of architectures with massive, hierarchical, multicore parallelism imposes to overhaul all existing solutions, which represents a major challenge for algorithm and software development;
- the evolution of numerical needs and types of simulations increase the importance, frequency, and size of certain classes of matrices, which may benefit from a specialized processing (rather than resort to a generic one).

Our research in the field is strongly related to the software package MUMPS, which is both an experimental platform for academics in the field of sparse linear algebra, and a software package that is widely used in both academia and industry. The software package MUMPS enables us to (i) confront our research to the real world, (ii) develop contacts and collaborations, and (iii) receive continuous feedback from real-life applications, which is extremely critical to validate our research work. The feedback from a large user community also enables us to direct our long-term objectives towards meaningful directions.

In this context, we aim at designing parallel sparse direct methods that will scale to large modern platforms, and that are able to answer new challenges arising from applications, both efficiently—from a resource consumption point of view—and accurately—from a numerical point of view. For that, and even with increasing parallelism, we do not want to sacrifice in any manner numerical stability, based on threshold partial pivoting, one of the main originalities of our approach (our “trademark”) in the context of direct solvers for distributed-memory computers; although this makes the parallelization more complicated, applying the same pivoting strategy as in the serial case ensures numerical robustness of our approach, which we generally measure in terms of sparse backward error. In order to solve the hard problems resulting from the always-increasing demands in simulations, special attention must also necessarily be paid to memory usage (and not only execution time). This requires specific algorithmic choices and scheduling techniques. From a complementary point of view, it is also necessary to be aware of the functionality requirements from the applications and from the users, so that robust solutions can be proposed for a wide range of applications.

Among direct methods, we rely on the multifrontal method [58], [59], [63]. This method usually exhibits a good data locality and hence is efficient in cache-based systems. The task graph associated with the multifrontal method is in the form of a tree whose characteristics should be exploited in a parallel implementation.

Our work is organized along two main research directions. In the first one we aim at efficiently addressing new architectures that include massive, hierarchical parallelism. In the second one, we aim at reducing the running time complexity and the memory requirements of direct solvers, while controlling accuracy.

### **3.3.2. Combinatorial scientific computing**

Combinatorial scientific computing (CSC) is a recently coined term (circa 2002) for interdisciplinary research at the intersection of discrete mathematics, computer science, and scientific computing. In particular, it refers to the development, application, and analysis of combinatorial algorithms to enable scientific computing applications. CSC’s deepest roots are in the realm of direct methods for solving sparse linear systems of equations where graph theoretical models have been central to the exploitation of sparsity, since the 1960s. The general approach is to identify performance issues in a scientific computing problem, such as memory use, parallel speed up, and/or the rate of convergence of a method, and to develop combinatorial algorithms and models to tackle those issues.

Our target scientific computing applications are (i) the preprocessing phases of direct methods (in particular MUMPS), iterative methods, and hybrid methods for solving linear systems of equations, and general sparse matrix and tensor computations; and (ii) the mapping of tasks (mostly the sub-tasks of the mentioned solvers) onto modern computing platforms. We focus on the development and the use of graph and hypergraph models, and related tools such as hypergraph partitioning algorithms, to solve problems of load balancing and task mapping. We also focus on bipartite graph matching and vertex ordering methods for reducing the memory overhead and computational requirements of solvers. Although we direct our attention on these models and algorithms through the lens of linear system solvers, our solutions are general enough to be applied to some other resource optimization problems.

### **3.3.3. Dense linear algebra on post-petascale multicore platforms**

The quest for efficient, yet portable, implementations of dense linear algebra kernels (QR, LU, Cholesky) has never stopped, fueled in part by each new technological evolution. First, the LAPACK library [52] relied on BLAS level 3 kernels (Basic Linear Algebra Subroutines) that enable to fully harness the computing power of a single CPU. Then the SCALAPACK library [51] built upon LAPACK to provide a coarse-grain parallel version, where processors operate on large block-column panels. Inter-processor communications

occur through highly tuned MPI send and receive primitives. The advent of multi-core processors has led to a major modification in these algorithms [53], [66], [61]. Each processor runs several threads in parallel to keep all cores within that processor busy. Tiled versions of the algorithms have thus been designed: dividing large block-column panels into several tiles allows for a decrease in the granularity down to a level where many smaller-size tasks are spawned. In the current panel, the diagonal tile is used to eliminate all the lower tiles in the panel. Because the factorization of the whole panel is now broken into the elimination of several tiles, the update operations can also be partitioned at the tile level, which generates many tasks to feed all cores.

The number of cores per processor will keep increasing in the following years. It is projected that high-end processors will include at least a few hundreds of cores. This evolution will require to design new versions of libraries. Indeed, existing libraries rely on a static distribution of the work: before the beginning of the execution of a kernel, the location and time of the execution of all of its component is decided. In theory, static solutions enable to precisely optimize executions, by taking parameters like data locality into account. At run time, these solutions proceed at the pace of the slowest of the cores, and they thus require a perfect load-balancing. With a few hundreds, if not a thousand, cores per processor, some tiny differences between the computing times on the different cores (“jitter”) are unavoidable and irremediably condemn purely static solutions. Moreover, the increase in the number of cores per processor once again mandates to increase the number of tasks that can be executed in parallel.

We study solutions that are part-static part-dynamic, because such solutions have been shown to outperform purely dynamic ones [55]. On the one hand, the distribution of work among the different nodes will still be statically defined. On the other hand, the mapping and the scheduling of tasks inside a processor will be dynamically defined. The main difficulty when building such a solution will be to design lightweight dynamic schedulers that are able to guarantee both an excellent load-balancing and a very efficient use of data locality.



## STORM Project-Team

### 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 adaptive 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<sup>0</sup>. However, very few applications use the whole machine. In general, an application runs only on a subset of the nodes<sup>0</sup>. 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.

---

<sup>0</sup>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.

<sup>0</sup>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.

## DIVERSE Project-Team

### 3. Research Program

#### 3.1. Scientific background

##### 3.1.1. Model-Driven Engineering

Model-Driven Engineering (MDE) aims at reducing the accidental complexity associated with developing complex software-intensive systems (e.g., use of abstractions of the problem space rather than abstractions of the solution space) [120]. It provides DIVERSE with solid foundations to specify, analyze and reason about the different forms of diversity that occur through the development lifecycle. A primary source of accidental complexity is the wide gap between the concepts used by domain experts and the low-level abstractions provided by general-purpose programming languages [91]. MDE approaches address this problem through modeling techniques that support separation of concerns and automated generation of major system artifacts from models (e.g., test cases, implementations, deployment and configuration scripts). In MDE, a model describes an aspect of a system and is typically created or derived for specific development purposes [73]. Separation of concerns is supported through the use of different modeling languages, each providing constructs based on abstractions that are specific to an aspect of a system. MDE technologies also provide support for manipulating models, for example, support for querying, slicing, transforming, merging, and analyzing (including executing) models. Modeling languages are thus at the core of MDE, which participates in the development of a sound *Software Language Engineering*<sup>0</sup>, including a unified typing theory that integrate models as first class entities [123].

Incorporating domain-specific concepts and high-quality development experience into MDE technologies can significantly improve developer productivity and system quality. Since the late nineties, this realization has led to work on MDE language workbenches that support the development of domain-specific modeling languages (DSMLs) and associated tools (e.g., model editors and code generators). A DSML provides a bridge between the field in which domain experts work and the implementation (programming) field. Domains in which DSMLs have been developed and used include, among others, automotive, avionics, and the emerging cyber-physical systems. A study performed by Hutchinson et al. [97] indicates that DSMLs can pave the way for wider industrial adoption of MDE.

More recently, the emergence of new classes of systems that are complex and operate in heterogeneous and rapidly changing environments raises new challenges for the software engineering community. These systems must be adaptable, flexible, reconfigurable and, increasingly, self-managing. Such characteristics make systems more prone to failure when running and thus development and study of appropriate mechanisms for continuous design and runtime validation and monitoring are needed. In the MDE community, research is focused primarily on using models at design, implementation, and deployment stages of development. This work has been highly productive, with several techniques now entering a commercialization phase. As software systems are becoming more and more dynamic, the use of model-driven techniques for validating and monitoring runtime behavior is extremely promising [105].

##### 3.1.2. Variability modeling

While the basic vision underlying *Software Product Lines* (SPL) can probably be traced back to David Parnas' seminal article [113] on the Design and Development of Program Families, it is only quite recently that SPLs are emerging as a paradigm shift towards modeling and developing software system families rather than individual systems [111]. SPL engineering embraces the ideas of mass customization and software reuse. It focuses on the means of efficiently producing and maintaining multiple related software products, exploiting what they have in common and managing what varies among them.

---

<sup>0</sup>See <http://planet-sl.org>

Several definitions of the *software product line* concept can be found in the research literature. Clements *et al.* define it as a *set of software-intensive systems sharing a common, managed set of features that satisfy the specific needs of a particular market segment or mission and are developed from a common set of core assets in a prescribed way* [110]. Bosch provides a different definition [79]: *A SPL consists of a product line architecture and a set of reusable components designed for incorporation into the product line architecture. In addition, the PL consists of the software products developed using the mentioned reusable assets.* In spite of the similarities, these definitions provide different perspectives of the concept: *market-driven*, as seen by Clements *et al.*, and *technology-oriented* for Bosch.

SPL engineering is a process focusing on capturing the *commonalities* (assumptions true for each family member) and *variability* (assumptions about how individual family members differ) between several software products [85]. Instead of describing a single software system, a SPL model describes a set of products in the same domain. This is accomplished by distinguishing between elements common to all SPL members, and those that may vary from one product to another. Reuse of core assets, which form the basis of the product line, is key to productivity and quality gains. These core assets extend beyond simple code reuse and may include the architecture, software components, domain models, requirements statements, documentation, test plans or test cases.

The SPL engineering process consists of two major steps:

1. **Domain Engineering**, or *development for reuse*, focuses on core assets development.
2. **Application Engineering**, or *development with reuse*, addresses the development of the final products using core assets and following customer requirements.

Central to both processes is the management of **variability** across the product line [93]. In common language use, the term *variability* refers to *the ability or the tendency to change*. Variability management is thus seen as the key feature that distinguishes SPL engineering from other software development approaches [80]. Variability management is thus growingly seen as the cornerstone of SPL development, covering the entire development life cycle, from requirements elicitation [125] to product derivation [130] to product testing [109], [108].

Halmans *et al.* [93] distinguish between *essential* and *technical* variability, especially at requirements level. Essential variability corresponds to the customer's viewpoint, defining what to implement, while technical variability relates to product family engineering, defining how to implement it. A classification based on the dimensions of variability is proposed by Pohl *et al.* [115]: beyond **variability in time** (existence of different versions of an artifact that are valid at different times) and **variability in space** (existence of an artifact in different shapes at the same time) Pohl *et al.* claim that variability is important to different stakeholders and thus has different levels of visibility: **external variability** is visible to the customers while **internal variability**, that of domain artifacts, is hidden from them. Other classification proposals come from Meekel *et al.* [103] (feature, hardware platform, performances and attributes variability) or Bass *et al.* [71] who discusses about variability at the architectural level.

Central to the modeling of variability is the notion of *feature*, originally defined by Kang *et al.* as: *a prominent or distinctive user-visible aspect, quality or characteristic of a software system or systems* [99]. Based on this notion of *feature*, they proposed to use a *feature model* to model the variability in a SPL. A feature model consists of a *feature diagram* and other associated information: *constraints* and *dependency rules*. Feature diagrams provide a *graphical tree-like notation depicting the hierarchical organization of high level product functionalities* represented as features. The root of the tree refers to the complete system and is progressively decomposed into more refined features (tree nodes). Relations between nodes (features) are materialized by *decomposition edges* and *textual constraints*. Variability can be expressed in several ways. Presence or absence of a feature from a product is modeled using *mandatory* or *optional features*. Features are graphically represented as rectangles while some graphical elements (e.g., unfilled circle) are used to describe the variability (e.g., a feature may be optional).

Features can be organized into *feature groups*. Boolean operators *exclusive alternative (XOR)*, *inclusive alternative (OR)* or *inclusive (AND)* are used to select one, several or all the features from a feature group.



Dependencies between features can be modeled using *textual constraints*: *requires* (presence of a feature requires the presence of another), *mutex* (presence of a feature automatically excludes another). Feature attributes can be also used for modeling quantitative (e.g., numerical) information. Constraints over attributes and features can be specified as well.

Modeling variability allows an organization to capture and select which version of which variant of any particular aspect is wanted in the system [80]. To implement it cheaply, quickly and safely, redoing by hand the tedious weaving of every aspect is not an option: some form of automation is needed to leverage the modeling of variability [75], [87]. Model Driven Engineering (MDE) makes it possible to automate this weaving process [98]. This requires that models are no longer informal, and that the weaving process is itself described as a program (which is as a matter of facts an executable meta-model [106]) manipulating these models to produce for instance a detailed design that can ultimately be transformed to code, or to test suites [114], or other software artifacts.

### 3.1.3. Component-based software development

Component-based software development [124] aims at providing reliable software architectures with a low cost of design. Components are now used routinely in many domains of software system designs: distributed systems, user interaction, product lines, embedded systems, etc. With respect to more traditional software artifacts (e.g., object oriented architectures), modern component models have the following distinctive features [86]: description of requirements on services required from the other components; indirect connections between components thanks to ports and connectors constructs [101]; hierarchical definition of components (assemblies of components can define new component types); connectors supporting various communication semantics [83]; quantitative properties on the services [78].

In recent years component-based architectures have evolved from static designs to dynamic, adaptive designs (e.g., SOFA [83], Palladio [76], Frascati [107]). Processes for building a system using a statically designed architecture are made of the following sequential lifecycle stages: requirements, modeling, implementation, packaging, deployment, system launch, system execution, system shutdown and system removal. If for any reason after design time architectural changes are needed after system launch (e.g., because requirements changed, or the implementation platform has evolved, etc) then the design process must be reexecuted from scratch (unless the changes are limited to parameter adjustment in the components deployed).

Dynamic designs allow for *on the fly* redesign of a component based system. A process for dynamic adaptation is able to reapply the design phases while the system is up and running, without stopping it (this is different from a stop/redeploy/start process). Dynamic adaptation process supports *chosen adaptation*, when changes are planned and realized to maintain a good fit between the needs that the system must support and the way it supports them [100]. Dynamic component-based designs rely on a component meta-model that supports complex life cycles for components, connectors, service specification, etc. Advanced dynamic designs can also take platform changes into account at runtime, without human intervention, by adapting themselves [84], [127]. Platform changes and more generally environmental changes trigger *imposed adaptation*, when the system can no longer use its design to provide the services it must support. In order to support an eternal system [77], dynamic component based systems must separate architectural design and platform compatibility. This requires support for heterogeneity, since platform evolution can be partial.

The Models@runtime paradigm denotes a model-driven approach aiming at taming the complexity of dynamic software systems. It basically pushes the idea of reflection one step further by considering the reflection layer as a real model “something simpler, safer or cheaper than reality to avoid the complexity, danger and irreversibility of reality [118]”. In practice, component-based (and/or service-based) platforms offer reflection APIs that make it possible to introspect the system (to determine which components and bindings are currently in place in the system) and dynamic adaptation (by applying CRUD operations on these components and bindings). While some of these platforms offer rollback mechanisms to recover after an erroneous adaptation, the idea of Models@runtime is to prevent the system from actually enacting an erroneous adaptation. In other words, the “model at run-time” is a reflection model that can be uncoupled (for reasoning, validation, simulation purposes) and automatically resynchronized.

Heterogeneity is a key challenge for modern component based system. Until recently, component based techniques were designed to address a specific domain, such as embedded software for command and control, or distributed Web based service oriented architectures. The emergence of the Internet of Things paradigm calls for a unified approach in component based design techniques. By implementing an efficient separation of concern between platform independent architecture management and platform dependent implementations, *Models@runtime* is now established as a key technique to support dynamic component based designs. It provides DIVERSE with an essential foundation to explore an adaptation envelop at run-time.

Search Based Software Engineering [95] has been applied to various software engineering problems in order to support software developers in their daily work. The goal is to automatically explore a set of alternatives and assess their relevance with respect to the considered problem. These techniques have been applied to craft software architecture exhibiting high quality of services properties [92]. Multi Objectives Search based techniques [89] deal with optimization problem containing several (possibly conflicting) dimensions to optimize. These techniques provide DIVERSE with the scientific foundations for reasoning and efficiently exploring an envelope of software configurations at run-time.

### 3.1.4. Validation and verification

Validation and verification (V&V) theories and techniques provide the means to assess the validity of a software system with respect to a specific correctness envelop. As such, they form an essential element of DIVERSE's scientific background. In particular, we focus on model-based V&V in order to leverage the different models that specify the envelop at different moments of the software development lifecycle.

Model-based testing consists in analyzing a formal model of a system (*e.g.*, activity diagrams, which capture high-level requirements about the system, statecharts, which capture the expected behavior of a software module, or a feature model, which describes all possible variants of the system) in order to generate test cases that will be executed against the system. Model-based testing [126] mainly relies on model analysis, constraint solving [88] and search-based reasoning [102]. DIVERSE leverages in particular the applications of model-based testing in the context of highly-configurable systems and [128] interactive systems [104] as well as recent advances based on diversity for test cases selection [96].

Nowadays, it is possible to simulate various kinds of models. Existing tools range from industrial tools such as Simulink, Rhapsody or Telelogic to academic approaches like Omega [112], or Xholon<sup>0</sup>. All these simulation environments operate on homogeneous environment models. However, to handle diversity in software systems, we also leverage recent advances in heterogeneous simulation. Ptolemy [82] proposes a common abstract syntax, which represents the description of the model structure. These elements can be decorated using different directors that reflect the application of a specific model of computation on the model element. Metropolis [72] provides modeling elements amenable to semantically equivalent mathematical models. Metropolis offers a precise semantics flexible enough to support different models of computation. ModHel'X [94] studies the composition of multi-paradigm models relying on different models of computation.

Model-based testing and simulation are complemented by runtime fault-tolerance through the automatic generation of software variants that can run in parallel, to tackle the open nature of software-intensive systems. The foundations in this case are the seminal work about N-version programming [70], recovery blocks [116] and code randomization [74], which demonstrated the central role of diversity in software to ensure runtime resilience of complex systems. Such techniques rely on truly diverse software solutions in order to provide systems with the ability to react to events, which could not be predicted at design time and checked through testing or simulation.

### 3.1.5. Empirical software engineering

The rigorous, scientific evaluation of DIVERSE's contributions is an essential aspect of our research methodology. In addition to theoretical validation through formal analysis or complexity estimation, we also aim at applying state-of-the-art methodologies and principles of empirical software engineering. This approach encompasses a set of techniques for the sound validation contributions in the field of software engineering,

---

<sup>0</sup><http://www.primordion.com/Xholon/>

ranging from statistically sound comparisons of techniques and large-scale data analysis to interviews and systematic literature reviews [121], [119]. Such methods have been used for example to understand the impact of new software development paradigms [81]. Experimental design and statistical tests represent another major aspect of empirical software engineering. Addressing large-scale software engineering problems often requires the application of heuristics, and it is important to understand their effects through sound statistical analyses [69].

## 3.2. Research axis

Figure 1 illustrates the four dimensions of software diversity, which form the core research axis of DIVERSE: the **diversity of languages** used by the stakeholders involved in the construction of these systems; the **diversity of features** required by the different customers; the **diversity of runtime environments** in which software has to run and adapt; the **diversity of implementations** that are necessary for resilience through redundancy. These four axes share and leverage the scientific and technological results developed in the area of model-driven engineering in the last decade. This means that all our research activities are founded on sound abstractions to reason about specific aspects of software systems, compose different perspectives and automatically generate parts of the system.

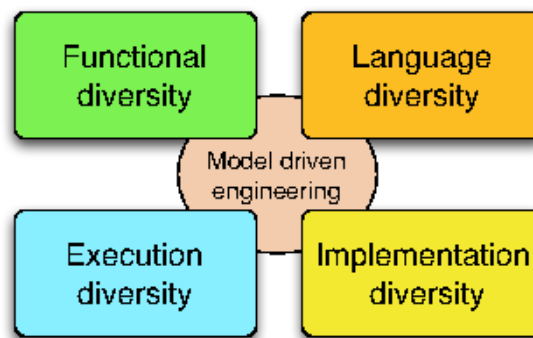


Figure 1. The four research axes of DIVERSE, which rely on a MDE scientific background

### 3.2.1. Software Language Engineering

The engineering of systems involves many different stakeholders, each with their own domain of expertise. Hence more and more organizations are adopting Domain Specific Modeling Languages (DSMLs) to allow domain experts to express solutions directly in terms of relevant domain concepts [120], [91]. This new trend raises new challenges about designing DSMLs, evolving a set of DSMLs and coordinating the use of multiple DSLs for both DSL designers and DSL users.

#### 3.2.1.1. Challenges

**Reusability** of software artifacts is a central notion that has been thoroughly studied and used by both academics and industrials since the early days of software construction. Essentially, designing reusable artifacts allows the construction of large systems from smaller parts that have been separately developed and validated, thus reducing the development costs by capitalizing on previous engineering efforts. However, it is still hardly possible for language designers to design typical language artifacts (e.g. language constructs, grammars, editors or compilers) in a reusable way. The current state of the practice usually prevents the reusability of language artifacts from one language to another, consequently hindering the emergence of real engineering techniques around software languages. Conversely, concepts and mechanisms that enable artifacts reusability abound in the software engineering community.

**Variability** in modeling languages occur in the definition of the abstract and concrete syntax as well as in the specification of the language's semantics. The major challenges met when addressing the need for variability are: (i) to set principles for modeling language units that support the modular specification of a modeling language; and (ii) to design mechanisms to assemble these units into a complete language, according to the set of authorized variation points for the modeling language family.

A new generation of complex software-intensive systems (for example smart health support, smart grid, building energy management, and intelligent transportation systems) gives new opportunities for leveraging modeling languages. The development of these systems requires expertise in diverse domains. Consequently, different types of stakeholders (e.g., scientists, engineers and end-users) must work in a coordinated manner on various aspects of the system across multiple development phases. DSMLs can be used to support the work of domain experts who focus on a specific system aspect, but they can also provide the means for coordinating work across teams specializing in different aspects and across development phases. The support and integration of DSMLs leads to what we call **the globalization of modeling languages**, *i.e.* the use of multiple languages for the coordinated development of diverse aspects of a system. One can make an analogy with world globalization in which relationships are established between sovereign countries to regulate interactions (e.g., travel and commerce related interactions) while preserving each country's independent existence.

### 3.2.1.2. Scientific objectives

We address reuse and variability challenges through the investigation of the time-honored concepts of substitutability, inheritance and components, evaluate their relevance for language designers and provide tools and methods for their inclusion in software language engineering. We will develop novel techniques for the modular construction of language extensions with support to model syntactical variability. From the semantics perspective, we investigate extension mechanisms for the specification of variability in operational semantics, focusing on static introduction and heterogeneous models of computation. The definition of variation points for the three aspects of the language definition provides the foundations for the novel concept Language Unit (LU) as well as suitable mechanisms to compose such units.

We explore the necessary breakthrough in software languages to support modeling and simulation of heterogeneous and open systems. This work relies on the specification of executable domain specific modeling languages (DSMLs) to formalize the various concerns of a software-intensive system, and of models of computation (MoCs) to explicitly model the concurrency, time and communication of such DSMLs. We develop a framework that integrates the necessary foundations and facilities for designing and implementing executable and concurrent domain-specific modeling languages. This framework also provides unique features to specify composition operators between (possibly heterogeneous) DSMLs. Such specifications are amenable to support the edition, execution, graphical animation and analysis of heterogeneous models. The objective is to provide both a significant improvement to MoCs and DSMLs design and implementation and to the simulation based validation and verification of complex systems.

We see an opportunity for the automatic diversification of programs' computation semantics, for example through the diversification of compilers or virtual machines. The main impact of this artificial diversity is to provide flexible computation and thus ease adaptation to different execution conditions. A combination of static and dynamic analysis could support the identification of what we call *plastic computation zones* in the code. We identify different categories of such zones: (i) areas in the code in which the order of computation can vary (e.g., the order in which a block of sequential statements is executed); (ii) areas that can be removed, keeping the essential functionality [122] (e.g., skip some loop iterations); (iii) areas that can be replaced by alternative code (e.g., replace a try-catch by a return statement). Once we know which zones in the code can be randomized, it is necessary to modify the model of computation to leverage the computation plasticity. This consists in introducing variation points in the interpreter to reflect the diversity of models of computation. Then, the choice of a given variation is performed randomly at run time.

### 3.2.2. Variability Modeling and Engineering

The systematic modeling of variability in software systems has emerged as an effective approach to document and reason about software evolution and heterogeneity (*cf.* Section 3.1.2). Variability modeling characterizes

an “envelope” of possible software variations. The industrial use of variability models and their relation to software artifact models require a complete engineering framework, including composition, decomposition, analysis, configuration and artifact derivation, refactoring, re-engineering, extraction, and testing. This framework can be used both to tame imposed diversity and to manage chosen diversity.

### 3.2.2.1. Challenges

A fundamental problem is that the **number of variants** can be exponential in the number of options (features). Already with 300 boolean configuration options, approximately  $10^{90}$  configurations exist – more than the estimated count of atoms in the universe. Domains like automotive or operating systems have to manage more than 10000 options (e.g., Linux). Practitioners face the challenge of developing billions of variants. It is easy to forget a necessary constraint, leading to the synthesis of unsafe variants, or to under-approximate the capabilities of the software platform. Scalable modelling techniques are therefore crucial to specify and reason about a very large set of variants.

Model-driven development supports two approaches to deal with the increasing number of concerns in complex systems: multi-view modeling, *i.e.* when modeling each concern separately, and variability modeling. However, there is little support to combine both approaches consistently. Techniques to integrate both approaches will enable the construction of a consistent set of views and variation points in each view.

The design, construction and maintenance of software families have a major impact on **software testing**. Among the existing challenges, we can cite: the selection of test cases for a specific variant; the evolution of test suites with integration of new variants; the combinatorial explosion of the number of software configurations to be tested. Novel model-based techniques for test generation and test management in a software product line context are needed to overcome state-of-the-art limits we already observed in some projects.

### 3.2.2.2. Scientific objectives

We aim at developing scalable reasoning techniques to **automatically analyze** variability models and their interactions with other views on the software intensive system (requirements, architecture, design, code). These techniques provide two major advancements in the state of the art: (1) an extension of the semantics of variability models in order to enable the definition of attributes (*e.g.*, cost, quality of service, effort) on features and to include these attributes in the reasoning; (2) an assessment of the consistent specification of variability models with respect to system views (since variability is orthogonal to system modeling, it is currently possible to specify the different models in ways that are semantically meaningless). The former aspect of analysis is tackled through constraint solving and finite-domain constraint programming, while the latter aspect is investigated through automatic search-based and learning-based techniques for the exploration of the space of interaction between variability and view models.

We aim at developing procedures to **reverse engineer** dependencies and features’ sets from existing software artefacts – be it source code, configuration files, spreadsheets (*e.g.*, product comparison matrices) or requirements. We expect to scale up (*e.g.*, for extracting a very large number of variation points) and guarantee some properties (*e.g.*, soundness of configuration semantics, understandability of ontological semantics). For instance, when building complex software-intensive systems, textual requirements are captured in very large quantities of documents. In this context, adequate models to formalize the organization of requirements documents and automated techniques to support impact analysis (in case of changes in the requirements) have to be developed.

### 3.2.3. Heterogeneous and dynamic software architectures

Flexible yet dependable systems have to cope with heterogeneous hardware execution platforms ranging from smart sensors to huge computation infrastructures and data centers. Evolution possibilities range from a mere change in the system configuration to a major architectural redesign, for instance to support addition of new features or a change in the platform architecture (*e.g.*, new hardware is made available, a running system switches to low bandwidth wireless communication, a computation node battery is running low, etc). In this context, we need to devise formalisms to reason about the impact of an evolution and about the transition from one configuration to another. It must be noted that this axis focuses on the use of models to drive the evolution

from design time to runtime. Models will be used to (i) systematically define predictable configurations and variation points through which the system will evolve; (ii) develop behaviors necessary to handle unforeseen evolution cases.

### 3.2.3.1. Challenges

The main challenge is to provide new homogeneous architectural modelling languages and efficient techniques that enable continuous software reconfiguration to react to changes. This work handles the challenges of handling the diversity of runtime infrastructures and managing the cooperation between different stakeholders. More specifically, the research developed in this axis targets the following dimensions of software diversity.

Platform architectural heterogeneity induces a first dimension of imposed diversity (type diversity). Platform reconfiguration driven by changing resources define another dimension of diversity (deployment diversity). To deal with these imposed diversity problems, we will rely on model based runtime support for adaptation, in the spirit of the dynamic distributed component framework developed by the Triskell team. Since the runtime environment composed of distributed, resource constrained hardware nodes cannot afford the overhead of traditional runtime adaptation techniques, we investigate the design of novel solutions relying on Models@runtime and on specialized tiny virtual machines to offer resource provisioning and dynamic reconfiguration.

Diversity can also be an asset to optimize software architecture. Architecture models must integrate multiple concerns in order to properly manage the deployment of software components over a physical platform. However, these concerns can contradict each other (*e.g.*, accuracy and energy). In this context, we investigate automatic solutions to explore the set of possible architecture models and to establish valid trade-offs between all concerns in case of changes.

### 3.2.3.2. Scientific objectives

**Automatic synthesis of optimal software architectures.** Implementing a service over a distributed platform (*e.g.*, a pervasive system or a cloud platform) consists in deploying multiple software components over distributed computation nodes. We aim at designing search-based solutions to (i) assist the software architect in establishing a good initial architecture (that balances between different factors such as cost of the nodes, latency, fault tolerance) and to automatically update the architecture when the environment or the system itself change. The choice of search-based techniques is motivated by the very large number of possible software deployment architectures that can be investigated and that all provide different trade-offs between qualitative factors. Another essential aspect that is supported by multi-objective search is to explore different architectural solutions that are not necessarily comparable. This is important when the qualitative factors are orthogonal to each other, such as security and usability for example.

**Flexible software architecture for testing and data management.** As the number of platforms on which software runs increases and different software versions coexist, the demand for testing environments also increases. For example, the number of testing environments to test a software patch or upgrade is the product of the number of execution environments the software supports and the number of coexisting versions of the software. Based on our first experiment on the synthesis of cloud environment using architectural models, our objective is to define a set of domain specific languages to catch the requirement and to design cloud environments for testing and data management of future internet systems from data centers to things. These languages will be interpreted to support dynamic synthesis and reconfiguration of a testing environment.

**Runtime support for heterogeneous environments.** Execution environments must provide a way to account or reserve resources for applications. However, current execution environments such as the Java Virtual Machine do not clearly define a notion of application: each framework has its own definition. For example, in OSGi, an application is a component, in JEE, an application is most of the time associated to a class loader, in the Multi-Tasking Virtual machine, an application is a process. The challenge consists in defining an execution environment that provides direct control over resources (CPU, Memory, Network I/O) independently from the definition of an application. We propose to define abstract resource containers to account and reserve resources on a distributed network of heterogeneous devices.



### 3.2.4. Diverse implementations for resilience

Open software-intensive systems have to evolve over their lifetime in response to changes in their environment. Yet, most verification techniques assume a closed environment or the ability to predict all changes. Dynamic changes and evolution cases thus represent a major challenge for these techniques that aim at assessing the correctness and robustness of the system. On the one hand, DIVERSE will adapt V&V techniques to handle diversity imposed by the requirements and the execution environment, on the other hand we leverage diversity to increase the robustness of software in face of unforeseen situations. More specifically, we address the following V&V challenges.

#### 3.2.4.1. Challenges

One major challenge to build flexible and open yet dependable systems is that current software engineering techniques require architects to foresee all possible situations the system will have to face. However, openness and flexibility also mean unpredictability: unpredictable bugs, attacks, environmental evolution, etc. Current fault-tolerance [116] and security [90] techniques provide software systems with the capacity of detecting accidental and deliberate faults. However, existing solutions assume that the set of bugs or vulnerabilities in a system does not evolve. This assumption does not hold for open systems, thus it is essential to revisit fault-tolerance and security solutions to account for diverse and unpredictable faults.

Diversity is known to be a major asset for the robustness of large, open, and complex systems (*e.g.*, economical or ecological systems). Following this observation, the software engineering literature provides a rich set of work that rely on implementation diversity in software systems in order to improve robustness to attacks or to changes in quality of service. These works range from N-version programming to obfuscation of data structures or control flow, to randomization of instruction sets. An essential and active challenge is to support the automatic synthesis and evolution of software diversity in open software-intensive systems. There is an opportunity to further enhance these techniques in order to cope with a wider diversity of faults, by multiplying the levels of diversity in the different software layers that are found in software-intensive systems (system, libraries, frameworks, application). This increased diversity must be based on artificial program transformations and code synthesis, which increase the chances of exploring novel solutions, better fitted at one point in time. The biological analogy also indicates that diversity should emerge as a side-effect of evolution, to prevent over-specialization towards one kind of diversity.

#### 3.2.4.2. Scientific objectives

The main objective is to address one of the main limitations of N-version programming for fault-tolerant systems: the manual production and management of software diversity. Through automated injection of artificial diversity we aim at systematically increasing failure diversity and thus increasing the chances of early error detection at run-time. A fundamental assumption for this work is that software-intensive systems can be “good enough” [117], [129].

**Proactive program diversification.** We aim at establishing novel principles and techniques that favor the emergence of multiple forms of software diversity in software-intensive systems, in conjunction with the software adaptation mechanisms that leverage this diversity. The main expected outcome is a set of meta-design principles that maintain diversity in systems and the experimental demonstration of the effects of software diversity. Higher levels of diversity in the system provide a pool of software solutions that can eventually be used to adapt to situations unforeseen at design time (bugs, crash, attacks, etc.). Principles of automated software diversification rely on the automated synthesis of variants in a software product line, as well as finer-grained program synthesis combining unsound transformations and genetic programming to explore the space of mutational robustness.

**Multi-tier software diversification.** We name multi-tier diversification the fact of diversifying several application software components simultaneously. The novelty of our proposal, with respect to the software diversity state of the art, is to diversify the application-level code (for example, diversify the business logic of the application), focusing on the technical layers found in web applications. The diversification of application software code is expected to provide a diversity of failures and vulnerabilities in web server deployment. Web server deployment usually adopts a form of the Reactor architecture pattern, for scalability purposes:

multiple copies of the server software stack, called request handlers, are deployed behind a load balancer. This architecture is very favorable for diversification, since by using the multiplicity of request handlers running in a web server we can simultaneously deploy multiple combinations of diverse software components. Then, if one handler is hacked or crashes the others should still be able to process client requests.



## EASE Project-Team

### 3. Research Program

#### 3.1. Collecting pertinent information

In our model, applications adapt their behavior (for instance, the level of automation) to the quality of their perception of the environment. This is important to alleviate the development constraint we usually have on automated systems. We "just" have to be sure a given process will always operate at the right automation level given the precision, the completeness or the confidence it has on its own perception. For instance, a car passing through crossing would choose its speed depending on the confidence it has gained during perception data gathering. When it has not enough information or when it could not trust it, it should reduce the automation level, therefore the speed, to only rely on its own sensors. Such adaptation capability shift requirements from the design and deployment (availability, robustness, accuracy, etc.) to the **assessment of the environment perception** we aim to facilitate in this first research axis.

*Data characterization.* The quality (freshness, accuracy, confidence, reliability, confidentiality, etc.) of the data are of crucial importance to assess the quality of the perception and therefore to ensure proper behavior. The way data is produced, consolidated, and aggregated while flowing to the consumer has an impact on its quality. Moreover part of these quality attributes requires to gather information at several communication layers from various entities. For this purpose, we want to design **lightweight cross-layer interactions** to collect relevant data. As a "frugality" principle should guide our approach, it is not appropriate to build all attributes we can imagine. It is therefore necessary to identify attributes relevant to the application and to have mechanisms to activate/deactivate at run-time the process to collect them.

*Data fusion.* Raw data should be directly used only to determine low-level abstraction. Further help in abstracting from low-level details can be provided by **data fusion** mechanisms. A good (re)construction of a meaningful information for the application reduces the complexity of the pervasive applications and helps the developers to concentrate on the application logic rather on the management of raw data. Moreover, the reactivity required in pervasive systems and the aggregation of large amounts of data (and its processing) are antagonists. We study **software services that can be deployed closer to the edge of the network**. The exploration of data fusion technics will be guided by different criteria: relevance of abstractions produced for pervasive applications, anonymization of exploited raw data, processing time, etc.

*Assessing the correctness of the behavior.* To ease the design of new applications and to align the development of new products with the ever faster standard developments, continuous integration could be used in parallel with continuous conformance and interoperability testing. We already participate in the design of new shared platforms that aims at facilitating this providing remote testing tools. Unfortunately, it is not possible to be sure that all potential peers in the surrounding have a conform behavior. Moreover, upon failure or security breach, a piece of equipment could stop to operate properly and lead to global mis-behavior. We want to propose conceptual tools for **testing at runtime devices in the environment**. The result of such conformance or interoperability tests could be stored safely in the environment by authoritative testing entity. Then application could interact with the device with a higher confidence. The confidence level of a device could be part of the quality attribute of the information it contributed to generate. The same set of tools could be used to identify misbehaving device for maintenance purpose or to trigger further testing.

#### 3.2. Building relevant abstraction for new interactions

The pervasive applications are often designed in an ad hoc manner depending on the targeted application area. Ressources (sensors / actuators, connected objets etc.) are often used in silos which complexify the implementation of rich pervasive computing scenarios. In the second research axis, we want to get away from technical aspects identifying **common and reusable system mechanisms** that could be used in various applications.

*Tagging the environment.* Information relative to environment could be stored by the application itself, but it could be complex to manage for mobile application since it could cross a large number of places with various features. Moreover the developer has to build its own representation of information especially when he wants to share information with other instances of the same application or with other applications. A promising approach is to store and to maintain this information associated to an object or to a place, in the environment itself. The infrastructure should provide services to application developers: add/retrieve information in the environment, share information and control who can access it, add computed properties to object for further usage. We want to study an **extensible model to describe and augment the environment**. Beyond a simple distributed storage, we have in mind a new kind of interaction between pervasive applications and changing environment and between applications themselves.

*Taking advantages of the spatial relationships.* To understand the world they have to interact with, pervasive applications often have to (re)built a model of it from the exchange they have with others or from their own observations. A part of the programmer's task consists in building a model of the spatial layout of the objects in the surrounding. The term *layout* can be understood in several ways: the co-location of multiple objects in the same vicinity, the physical arrangement of two objects relative to each other, or even the crossing of an object of a physical area to another, etc. Determining remotely these spatial properties (see figure 1 -a) is difficult without exchanging a lot of information. Properties related to the spatial layout are far easier to characterize locally. They could be abstracted from interaction pattern without any complex virtual representation of the environment (see figure 1 -b). We want to be able to rely on this type of spatial layout in a pervasive environment. In the prior years, the members of EASE already worked on **models for processing object interactions** in the physical world to automatically trigger processing. This was the case in particular of the spatial programming principle: physical space is treated as a tuple-space in which objects are automatically synchronized according to their spatial arrangement. We want to follow this approach by considering **richer and more expressive programming models**.

### 3.3. Acting on the environment

The conceptual tools we aim to study must be *frugal*: they use as less as possible resources, while having the possibility to use much more when it is required. Data needed by an application are not made available for "free"; for example, it costs energy to measure a characteristic of the environment, or to transmit it. So this "design frugality" requires a **fine-grained control** on how data is actually collected from the environment. The third research axis aims at designing solutions that give this control to application developers by **acting on the environment**.

*Acting on the data collection.* We want to be able to identify which information are really needed during the perception elaboration process. If a piece of data is missing to build a given information with the appropriate quality level, the data collection mechanism should find relevant information in the environment or modify the way it aggregates it. These could lead to a modification of the behavior of the network layer and the path the piece of data uses in the aggregation process.

*Acting on object interactions.* Objects in the environment could adapt their behavior in a way that strongly depends on the object itself and that is difficult to generalize. Beyond the specific behaviors of actuators triggered through specialized or standard interfaces, the production of information required by an application could necessitate an adaptation at the object level (eg. calibration, sampling). The environment should then be able to initiate such adaption transparently to the application, which may not know all objects it passes by.

*Adapting object behaviors.* The radio communication layers become more flexible and able to adapt the way they use energy to what is really required for a given transmission. We already study how beamforming technics could be used to adapt multicast strategy for video services. We want to show how playing with these new parameters of transmissions (eg. beamforming, power, ...) allows to control spatial relationships objects could have. There is a tradeoff to find between the capacity of the medium, the electromagnetic pollution and the reactivity of the environment. We plan to extend our previous work on interface selection and more generally on what we call **opportunistic networking**.

## FOCUS Project-Team

### 3. Research Program

#### 3.1. Foundations 1: Models

The objective of Focus is to develop concepts, techniques, and possibly also tools, that may contribute to the analysis and synthesis of CBUS. Fundamental to these activities is *modeling*. Therefore designing, developing and studying computational models appropriate for CBUS is a central activity of the project. The models are used to formalise and verify important computational properties of the systems, as well as to propose new linguistic constructs.

The models we study are in the process calculi (e.g., the  $\pi$ -calculus) and  $\lambda$ -calculus tradition. Such models, with their emphasis on algebra, well address compositionality—a central property in our approach to problems. Accordingly, the techniques we employ are mainly operational techniques based on notions of behavioural equivalence, and techniques based on algebra, mathematical logics, and type theory.

#### 3.2. Foundations 2: Foundational calculi and interaction

Modern distributed systems have witnessed a clear shift towards interaction and conversations as basic building blocks for software architects and programmers. The systems are made by components, that are supposed to interact and carry out dialogues in order to achieve some predefined goal; Web services are a good example of this. Process calculi are models that have been designed precisely with the goal of understanding interaction and composition. The theory and tools that have been developed on top of process calculi can set a basis with which CBUS challenges can be tackled. Indeed industrial proposals of languages for Web services such as BPEL are strongly inspired by process calculi, notably the  $\pi$ -calculus.

#### 3.3. Foundations 3: Type systems and logics

Type systems and logics for reasoning on computations are among the most successful outcomes in the history of the research in  $\lambda$ -calculus and (more recently) in process calculi. Type systems can also represent a powerful means of specifying dialogues among components of CBUS. For instance—again referring to Web services—current languages for specifying interactions only express basic connectivity, ignoring causality and timing aspects (e.g., an intended order on the messages), and the alternative is to use Turing Complete languages that are however undecidable. Types can come at hand here: they can express causality and order information on messages [53], [49], [54], while remaining decidable systems.

#### 3.4. Foundations 4: Implicit computational complexity

A number of elegant and powerful results have been recently obtained in implicit computational complexity in the  $\lambda$ -calculus in which ideas from Linear Logics enable a fine-grained control over computations. This experience can be profitable when tackling issues of CBUS related to resource consumption, such as resources allocation, access to resources, certification of bounds on resource consumption (e.g., ensuring that a service will answer to a request in time polynomial with respect to the size of the input data).

## INDES Project-Team

### 3. Research Program

#### 3.1. Parallelism, concurrency, and distribution

Concurrency management is at the heart of diffuse programming. Since the execution platforms are highly heterogeneous, many different concurrency principles and models may be involved. Asynchronous concurrency is the basis of shared-memory process handling within multiprocessor or multicore computers, of direct or fifo-based message passing in distributed networks, and of fifo- or interrupt-based event handling in web-based human-machine interaction or sensor handling. Synchronous or quasi-synchronous concurrency is the basis of signal processing, of real-time control, and of safety-critical information acquisition and display. Interfacing existing devices based on these different concurrency principles within *Hop* or other diffuse programming languages will require better understanding of the underlying concurrency models and of the way they can nicely cooperate, a currently ill-resolved problem.

#### 3.2. Web, functional, and reactive programming

We are studying new paradigms for programming Web applications that rely on multi-tier functional programming. We have created a Web programming environment named *Hop*. It relies on a single formalism for programming the server-side and the client-side of the applications as well as for configuring the execution engine.

*Hop* is a functional language based on the SCHEME programming language. That is, it is a strict functional language, fully polymorphic, supporting side effects, and dynamically type-checked. *Hop* is implemented as an extension of the BIGLOO compiler that we develop. In the past, we have extensively studied static analyses (type systems and inference, abstract interpretations, as well as classical compiler optimizations) to improve the efficiency of compilation in both space and time.

As a *Hop* DSL, we have created *HipHop*, a synchronous orchestration language for web and IoT applications. *HipHop* facilitates the design and programming of complex web/IoT applications by smoothly integrating three computation models and programming styles that have been historically developed in different communities and for different purposes: *i*) *Transformational programs* that simply compute output values from input values, with comparatively simple interaction with their environment; *ii*) asynchronous concurrent programs that perform interactions between their components or with their environment with uncontrollable timing, using typically network-based communication; and *iii*) synchronous reactive programs that react to external events in a conceptually instantaneous and deterministic way.

#### 3.3. Security of diffuse programs

The main goal of our security research is to provide scalable and rigorous language-based techniques that can be integrated into multi-tier compilers to enforce the security of diffuse programs. Research on language-based security has been carried on before in former Inria teams. In particular previous research has focused on controlling information flow to ensure confidentiality.

Typical language-based solutions to these problems are founded on static analysis, logics, provable cryptography, and compilers that generate correct code by construction. Relying on the multi-tier programming language *Hop* that tames the complexity of writing and analysing secure diffuse applications, we are studying language-based solutions to prominent web security problems such as code injection and cross-site scripting, to name a few.

## RMOD Project-Team

### 3. Research Program

#### 3.1. Software Reengineering

Strong coupling among the parts of an application severely hampers its evolution. Therefore, it is crucial to answer the following questions: How to support the substitution of certain parts while limiting the impact on others? How to identify reusable parts? How to modularize an object-oriented application?

Having good classes does not imply a good application layering, absence of cycles between packages and reuse of well-identified parts. Which notion of cohesion makes sense in presence of late-binding and programming frameworks? Indeed, frameworks define a context that can be extended by subclassing or composition: in this case, packages can have a low cohesion without being a problem for evolution. How to obtain algorithms that can be used on real cases? Which criteria should be selected for a given remodularization?

To help us answer these questions, we work on enriching Moose, our reengineering environment, with a new set of analyses [31], [30]. We decompose our approach in three main and potentially overlapping steps:

1. Tools for understanding applications,
2. Remodularization analyses,
3. Software Quality.

##### 3.1.1. Tools for understanding applications

**Context and Problems.** We are studying the problems raised by the understanding of applications at a larger level of granularity such as packages or modules. We want to develop a set of conceptual tools to support this understanding.

Some approaches based on Formal Concept Analysis (FCA) [59] show that such an analysis can be used to identify modules. However the presented examples are too small and not representative of real code.

##### **Research Agenda.**

FCA provides an important approach in software reengineering for software understanding, design anomalies detection and correction, but it suffers from two problems: (i) it produces lattices that must be interpreted by the user according to his/her understanding of the technique and different elements of the graph; and, (ii) the lattice can rapidly become so big that one is overwhelmed by the mass of information and possibilities [20]. We look for solutions to help people putting FCA to real use.

##### 3.1.2. Remodularization analyses

**Context and Problems.** It is a well-known practice to layer applications with bottom layers being more stable than top layers [47]. Until now, few works have attempted to identify layers in practice: Mudpie [61] is a first cut at identifying cycles between packages as well as package groups potentially representing layers. DSM (dependency structure matrix) [60], [55] seems to be adapted for such a task but there is no serious empirical experience that validates this claim. From the side of remodularization algorithms, many were defined for procedural languages [43]. However, object-oriented programming languages bring some specific problems linked with late-binding and the fact that a package does not have to be systematically cohesive since it can be an extension of another one [62], [34].

As we are designing and evaluating algorithms and analyses to remodularize applications, we also need a way to understand and assess the results we are obtaining.

**Research Agenda.** We work on the following items:

Layer identification. We propose an approach to identify layers based on a semi-automatic classification of package and class interrelationships that they contain. However, taking into account the wish or knowledge of the designer or maintainer should be supported.

Cohesion Metric Assessment. We are building a validation framework for cohesion/coupling metrics to determine whether they actually measure what they promise to. We are also compiling a number of traditional metrics for cohesion and coupling quality metrics to evaluate their relevance in a software quality setting.

### **3.1.3. Software Quality**

**Research Agenda.** Since software quality is fuzzy by definition and a lot of parameters should be taken into account we consider that defining precisely a unique notion of software quality is definitively a Grail in the realm of software engineering. The question is still relevant and important. We work on the two following items:

Quality models. We studied existing quality models and the different options to combine indicators — often, software quality models happily combine metrics, but at the price of losing the explicit relationships between the indicator contributions. There is a need to combine the results of one metric over all the software components of a system, and there is also the need to combine different metric results for any software component. Different combination methods are possible that can give very different results. It is therefore important to understand the characteristics of each method.

Bug prevention. Another aspect of software quality is validating or monitoring the source code to avoid the emergence of well known sources of errors and bugs. We work on how to best identify such common errors, by trying to identify earlier markers of possible errors, or by helping identifying common errors that programmers did in the past.

## **3.2. Language Constructs for Modular Design**

While the previous axis focuses on how to help remodularizing existing software, this second research axis aims at providing new language constructs to build more flexible and recomposable software. We will build on our work on traits [57], [32] and classboxes [21] but also start to work on new areas such as isolation in dynamic languages. We will work on the following points: (1) Traits and (2) Modularization as a support for isolation.

### **3.2.1. Traits-based program reuse**

**Context and Problems.** Inheritance is well-known and accepted as a mechanism for reuse in object-oriented languages. Unfortunately, due to the coarse granularity of inheritance, it may be difficult to decompose an application into an optimal class hierarchy that maximizes software reuse. Existing schemes based on single inheritance, multiple inheritance, or mixins, all pose numerous problems for reuse.

To overcome these problems, we designed a new composition mechanism called Traits [57], [32]. Traits are pure units of behavior that can be composed to form classes or other traits. The trait composition mechanism is an alternative to multiple or mixin inheritance in which the composer has full control over the trait composition. The result enables more reuse than single inheritance without introducing the drawbacks of multiple or mixin inheritance. Several extensions of the model have been proposed [29], [51], [22], [33] and several type systems were defined [35], [58], [52], [45].

Traits are reusable building blocks that can be explicitly composed to share methods across unrelated class hierarchies. In their original form, traits do not contain state and cannot express visibility control for methods. Two extensions, stateful traits and freezable traits, have been proposed to overcome these limitations. However, these extensions are complex both to use for software developers and to implement for language designers.



**Research Agenda: Towards a pure trait language.** We plan distinct actions: (1) a large application of traits, (2) assessment of the existing trait models and (3) bootstrapping a pure trait language.

- To evaluate the expressiveness of traits, some hierarchies were refactored, showing code reuse [24]. However, such large refactorings, while valuable, may not exhibit all possible composition problems, since the hierarchies were previously expressed using single inheritance and following certain patterns. We want to redesign from scratch the collection library of Smalltalk (or part of it). Such a redesign should on the one hand demonstrate the added value of traits on a real large and redesigned library and on the other hand foster new ideas for the bootstrapping of a pure trait-based language.

In particular we want to reconsider the different models proposed (stateless [32], stateful [23], and freezable [33]) and their operators. We will compare these models by (1) implementing a trait-based collection hierarchy, (2) analyzing several existing applications that exhibit the need for traits. Traits may be flattened [50]. This is a fundamental property that confers to traits their simplicity and expressiveness over Eiffel’s multiple inheritance. Keeping these aspects is one of our priority in forthcoming enhancements of traits.

- Alternative trait models. This work revisits the problem of adding state and visibility control to traits. Rather than extending the original trait model with additional operations, we use a fundamentally different approach by allowing traits to be lexically nested within other modules. This enables traits to express (shared) state and visibility control by hiding variables or methods in their lexical scope. Although the traits’ “flattening property” no longer holds when they can be lexically nested, the combination of traits with lexical nesting results in a simple and more expressive trait model. We formally specify the operational semantics of this combination. Lexically nested traits are fully implemented in AmbientTalk, where they are used among others in the development of a Morphic-like UI framework.
- We want to evaluate how inheritance can be replaced by traits to form a new object model. For this purpose we will design a minimal reflective kernel, inspired first from ObjVlisp [28] then from Smalltalk [38].

### 3.2.2. Reconciling Dynamic Languages and Isolation

**Context and Problems.** More and more applications require dynamic behavior such as modification of their own execution (often implemented using reflective features [42]). For example, F-script allows one to script Cocoa Mac-OS X applications and Lua is used in Adobe Photoshop. Now in addition more and more applications are updated on the fly, potentially loading untrusted or broken code, which may be problematic for the system if the application is not properly isolated. Bytecode checking and static code analysis are used to enable isolation, but such approaches do not really work in presence of dynamic languages and reflective features. Therefore there is a tension between the need for flexibility and isolation.

**Research Agenda: Isolation in dynamic and reflective languages.** To solve this tension, we will work on *Sure*, a language where isolation is provided by construction: as an example, if the language does not offer field access and its reflective facilities are controlled, then the possibility to access and modify private data is controlled. In this context, layering and modularizing the meta-level [25], as well as controlling the access to reflective features [26], [27] are important challenges. We plan to:

- Study the isolation abstractions available in erights (<http://www.erights.org>) [49], [48], and Java’s class loader strategies [44], [39].
- Categorize the different reflective features of languages such as CLOS [41], Python and Smalltalk [53] and identify suitable isolation mechanisms and infrastructure [36].
- Assess different isolation models (access rights, capabilities [54],...) and identify the ones adapted to our context as well as different access and right propagation.
- Define a language based on
  - the decomposition and restructuring of the reflective features [25],



- the use of encapsulation policies as a basis to restrict the interfaces of the controlled objects [56],
- the definition of method modifiers to support controlling encapsulation in the context of dynamic languages.

An open question is whether, instead of providing restricted interfaces, we could use traits to grant additional behavior to specific instances: without trait application, the instances would only exhibit default public behavior, but with additional traits applied, the instances would get extra behavior. We will develop *Sure*, a modular extension of the reflective kernel of Smalltalk (since it is one of the languages offering the largest set of reflective features such as pointer swapping, class changing, class definition,...) [53].

## AGORA Project-Team

### 3. Research Program

#### 3.1. Wireless network deployment

The deployment of networks has fostered a constant research effort for decades, continuously renewed by the evolution of networking technologies. Fundamentally, the deployment problem addresses the trade-off between the cost of the network to be minimized or fitted into a budget and the features and services provided by the system, that should reach a target level or be maximized. The variety of cost models and type of features gives rise to a wide scientific field. There are several cost factors of network infrastructure: components (number and capacity), energy, man power (installation and maintenance), etc. The features of the network matter as much as the metric to evaluate them. Coverage and capacity are basic features for wireless networks on which we will focus in the following. One recurrent question is therefore: What are the optimal number and position of network components to deploy so that a given territory is covered and enough networking capacity is provided?

Traditional telecommunication infrastructures were made of dedicated components, each of them providing a given set of functions. However, recently introduced paradigms yield issues on the deployment of network functions. Indeed, the last decade saw a trend towards adding more intelligence within the network. In the case of the access network, the concept of Cloud Radio Access Network (C-RAN) emerged. In the backhaul, the Evolved Packet Core (EPC) network can also benefit from virtualization techniques, as the convergence point for multiple access technologies, as imagined in the case of future 5G networks. The performance limits of a virtualized EPC remain unknown today: Is the delay introduced by this new architecture compatible with the requirements of the mobile applications? How to deploy the different network functions on generic hardware in order to maximize the quality of service?

**Network component deployment.** In this research direction, we address new issues of the optimal network deployment. In particular, we focus on the deployment of wireless sensor networks for environmental monitoring (e.g. atmospheric pollution). Most current air quality monitoring systems are using conventional measuring stations, equipped with multiple lab quality sensors. These systems are however massive, inflexible and expensive. An alternative – or complementary – solution is to use low-cost flexible wireless sensor networks. One of the main challenges is to introduce adequate models for the coverage of the phenomenon. Most of the state of the art considers a generic coverage formulation based on detection ranges which are not adapted to environmental sensing. For example, pollution propagation models should take into account the inherently stochastic weather conditions. An issue is to develop an adequate formulation and efficient integer linear programming (ILP) models and heuristics able to compute deployments at a relevant scale. In particular, it seems promising to adapt stochastic or robust optimization results of the operational research community in order to deal with uncertainty. Defining the quality of a coverage is also a modeling issue, which depends on the application considered. The detection of anomaly is close to a combinatorial problem. A more difficult objective is to deploy sensors in order to map the phenomenon by interpolation (or other reconstruction mechanisms). This challenge requires interdisciplinary research with fluid mechanics teams who develop numerical models of pollution propagation and practitioners like Atmo Auvergne-Rhône-Alpes.

Regarding the network connectivity, another challenge is to integrate suitable wireless link models accounting for the deployment environment. For example, modeling the integration of sensors in urban areas is challenging due to the presence of neighboring walls and obstacles, as well as moving vehicles and pedestrians that may induce field scattering. Also, the urban constraints and characteristics need to be carefully modeled and considered. Indeed, the urban environment yields constraints or facilities on the deployment of sensor nodes and gateways, such as their embedding within street furniture. Understanding the structure of these spatial constraints is necessary to develop efficient optimization methods able to compute on large scale scenarios.

**Network function deployment.** In this research direction, we do not address network virtualization per se, but the algorithmic and architectural challenges that virtualization brings in both radio access and core networks. As a first challenge, we focus on the evaluation of Cloud Radio Access Network solutions. The capacity of a C-RAN architecture and the way this compares to classical RAN is still an open question. The fact that C-RAN enables cooperation between the remote radio heads (RRH) served by the same base-band units (BBU) indicates an improved performance, but at the same time the resulting cells are much larger, which goes against the current trend of increasing capacity through the deployment of small cells. We propose to study the problem both from a user and a network perspective. On the user side, we use standard information theory tools, such as multiple-access channels to model C-RAN scenarios and understand their performance. On the network side, this translates in a resource allocation problem with cooperative base stations. We will extend our previous models for non-cooperative scenarios. Regarding the core network function deployment, we are interested in the specific case of Professional Mobile Radio (PMR) networks. These networks, used for public safety services and in scenarios like post-disaster relief, present the particularity of an EPC formed by a mobile wireless network. Due to its nature, the network can not be pre-planned, and the different EPC functions need to be autonomously deployed on the available network elements. We study the EPC function deployment problem as an optimization problem, constrained by the user capacity requests. User attachment mechanisms will also be proposed, adapted to the network function distribution, the global user demand, and the source/destination of the flows. These challenges are tackled as centralized optimization problems, then extended to the context of real-time decisions. Finally, in order to complete these theoretical works based on ILP models and heuristics, experiments using OpenAir Interface are used to evaluate our proposals.

### 3.2. Wireless data collection

With an anticipated 11-fold growth between 2014 and 2018, facing the growth of the mobile demand is the foremost challenge for mobile operators. In particular, a 100-fold increase in the number of supported connected devices, mostly newly connected objects with M2M traffic, is expected. A question therefore arises: how to cope with a dense set of M2M low bit rate traffics from energy and computing power constrained devices while classic cellular infrastructures are designed for the sparse high bit rate traffics from powerful devices?

A technological answer to the densification challenge is also embodied by long-range low-power networks such as SigFox, LoRa, NB-IoT, etc. In this context, the idea of offloading cellular traffic to different wireless access technologies is emerging as a very promising solution to relieve the traditional mobile network from its overwhelming load. In fact, offloading is already employed today, and, globally, 45% of total mobile data traffic was offloaded onto the fixed network through Wi-Fi or femtocells in 2013. Device-to-device (D2D) communications in hybrid networks, combining long-range cellular links and short-range technologies, opens even more possibilities. We aim at providing solutions that are missing for efficiently and practically mix multi-hop and cellular networks technologies.

**Cellular M2M.** Enabling a communication in a cellular network follows two major procedures: a resource allocation demand is first transmitted by the UE which, if successful, is followed by the actual data transmission phase, using dedicated resources allocated by the eNodeB (eNB) to the UE. This procedure was designed specifically for H2H traffic, which is bursty by nature, and it is based on the notions of session and call, activities that keep the user involved for a relatively long time and necessitate the exchange of a series of messages with the network. On the contrary, M2M traffic generates low amounts of data periodically or sporadically. Going through a signaling-heavy random access (RA) procedure to transmit one short message is strongly inefficient for both the M2M devices and the infrastructure.

In the perspective of 5G solutions, we are investigating mechanisms that regulate the M2M traffics in order to obtain good performances while keeping a reasonable quality of service (QoS) for human-to-human (H2H) terminals. The idea of piggybacking the M2M data transmission within one of the RA procedure messages is tempting and it is now considered as the best solution for this type of traffic. This means that the M2M data is transmitted on the shared resources of the RACH, and raises questions regarding the capacity of the RACH, which was not designed for these purposes. In this regard, our analysis of the access capacity of LTE-A

RACH procedure has to be adapted to multi-class scenarios, in order to understand the competition between M2M and H2H devices. Modeling based on Markov chains provides trends on system scale performances, while event-based simulations enable the analysis of the distribution of the performances over the different kinds of users. Combining both should give enough insights so as to design relevant regulation techniques and strategies. In particular two open questions that have to be tackled can be stated as: When should access resources be opened to M2M traffics without penalizing H2H performances? Does an eNodeB have a detailed enough knowledge of the system and transmit enough information to UE to regulate the traffics? The objective is to go to the analysis of achievable performances to actual protocols that take into account realistic M2M traffic patterns.

**Hybrid networks.** The first objective in this research axis is a realistic large-scale performance evaluation of Wi-Fi offloading solutions. While the mechanisms behind Wi-Fi offloading are now clear in the research community, their performance has only been tested in small-scale field tests, covering either small geographical areas (i.e. a few cellular base stations) and/or a small number of specific users (e.g. vehicular users). Instead, we evaluate the offloading performance at a city scale, building on real mobile network traces available in the team. First of all, through our collaboration with Orange Labs, we have access to an accurate characterization of the mobile traffic load at each base station in all major French cities. Second, a data collection application for Android devices has been developed in the team and used by hundreds of users in the Lyon metropolitan area. This application monitors and logs all the Wi-Fi access points in the coverage range of the smartphone, allowing us to build a map of Wi-Fi accessibility in some parts of the city. Combining these two data sources and completing them with simulation studies will allow an accurate evaluation of Wi-Fi offloading solutions over a large area.

On the D2D side, our focus is on the connected objects scenario, where we study the integration of short-range links and long-range technologies such as LTE, SigFox or LoRa. This requires the design of network protocols to discover and group the devices in a certain region. For this, we build on our expertise on clustering sensor and vehicular nodes. The important difference in this case is that the cellular network can assist the clustering formation process. The next step is represented by the selection of the devices that will be using the long-range links on behalf of the entire cluster. With respect to classical cluster head selection problems in ad-hoc networks, our problem distinguishes itself through device heterogeneity in terms of available communication technologies (not all devices have a long-range connection, or their quality is poor), energy resources (some devices might have energy harvesting capabilities) and expected lifetime. We will evaluate the proposed mechanisms both analytically (clustering problems are generally modeled by dominating set problems in graph theory) and through discrete-event simulation. Prototyping and experimental evaluation in cooperation with our industrial partners is also foreseen in this case.

### 3.3. Network data exploitation

Mobile devices are continuously interacting with the network infrastructure, and the associated geo-referenced events can be easily logged by the operators, for different purposes, including billing and resource management. This leads to the implicit possibility of monitoring a large percentage of the whole population with minimal cost: no other technology provides today an equivalent coverage. On the networking side, the exploitation of data collected within the cellular network can be the enabler of flexible and reconfigurable cellular systems. In order to enable this vision, algorithmic solutions are needed that drive, in concert with the variations in the mobile demand, the establishment, modification, release and relocation of any type of resources in the network. This raises, in turn, the fundamental problem of understanding the mobile demand, and linking it to the resource management processes. More precisely, we contribute to answer questions about the correlation between urban areas and mobile traffic usage, in particular the spatial and temporal causalities in the usage of the mobile network.

In a different type of architecture, the one of wireless sensor networks, the spatio-temporal characteristics of the data that are transported can also be leveraged to improve on the networking performances, e.g. capacity and energy consumption. In several applications (e.g. temperature monitoring, intrusion detection), wireless sensor nodes are prone to transmit redundant or correlated information. This wastes the bandwidth

and accelerates the battery depletion. Energy and network capacity savings can be obtained by leveraging spatial and temporal correlation in packet aggregation. Packet transmissions can be reduced with an overhead induced by distributed aggregation algorithms. We aim at designing data aggregation functions that preserve data accuracy and maximize the network lifetime with low assumptions on the network topology and the application.

**Mobile data analysis.** In this research axis, we delve deeper in the analysis of mobile traffic. In this sense, temporal and spatial usage profiles can be built, by including in our analysis datasets providing service-level usage information. Indeed, previous studies have been generally using call detail records (CDR) or, at best, aggregated packet traffic information. This data is already very useful in many research fields, but fine-grained usage data would allow an even better understanding of the spatiotemporal characteristics of mobile traffic. To achieve this, we exploit datasets made available by Orange Labs, providing information about the network usage for several different mobile services (web, streaming, download, mail, etc.).

To obtain even richer information, we combine this operator-side data with user-side data, collected by a crowdsensing application we developed within the PrivaMov research project. While covering hundreds of thousands of users, operator data only allows to localize the user at the cell level, and only when the user is connected to the network. The crowdsensing application we are using gathers precise GPS user localization data at a high frequency. Combining these two sources of data will allow us to gain insight in possible biases introduced by operator-side data and to infer microscopic properties which, correctly modeled, can be extended to the entire user population, even those for which we do not possess crowdsensed data.

Privacy preservation is an important topic in the field of mobile data analysis. Mobile traffic data anonymization techniques are currently proposed, mainly by adding noise or removing information from the original dataset. While we do not plan to develop anonymization algorithms, we collaborate with teams working on this topic (e.g. Inria Privatics) in order to assess the impact of anonymization techniques on the spatio-temporal properties of mobile traffic data. Through a statistical analysis of both anonymized and non-anonymized data, we hope to better understand the usability of anonymized data for different applications based on the exploration of mobile traffic data.

**Data aggregation.** Data-aggregation takes benefit from spatial and/or temporal correlation, while preserving the data accuracy. Such correlation comes from the physical phenomenon which is observed. Temporal aggregation is mainly addressed using temporal series (e.g. ARMA) whereas spatial aggregation is now led by compressive sensing solutions. Our objective is to get rid of the assumption of knowing of the network topology properties and the data traffic generated by the application, in particular for dense and massive wireless networks. Note that we focus on data-aggregation with a networking perspective, not with the background of information theory.

The rational design of an aggregation scheme implies understanding data dynamics (statistical characteristics, information representation), algorithmic optimization (aggregator location, minimizing the number of aggregators toward energy efficiency), and network dynamics (routing, medium sharing policies, node activity). We look for designing a complete aggregation chain including both intra-sensor aggregation and inter-sensor aggregation. For this, we characterize the raw data that are collected in order to understand the dynamics behind several key applications. The goal is to provide a taxonomy of the applications according to the data properties in terms of stationarity, dynamics, etc. Then, we aim to design temporal aggregation functions without knowledge of the network topology and without assumptions about the application data. Such functions should be able to self-adapt to the environment evolution. A related issue is the deployment of aggregators into the wireless network to allow spatial aggregation with respect to the energy consumption minimization, capacity saving maximization and distributed algorithm complexity. We therefore look to define dedicated protocols for each aggregation function family.

## COATI Project-Team

### 3. Research Program

#### 3.1. Research Program

Members of COATI have a strong expertise in the design and management of wired and wireless backbone, backhaul, broadband, software defined and complex networks. On the one hand, we cope with specific problems such as energy efficiency in backhaul and backbone networks, routing reconfiguration in connection oriented networks (MPLS, WDM), traffic aggregation in SONET networks, compact routing in large-scale networks, survivability to single and multiple failures, etc. These specific problems often come from questions of our industrial partners. On the other hand, we study fundamental problems mainly related to routing and reliability that appear in many networks (not restricted to our main fields of applications) and that have been widely studied in the past. However, previous solutions do not take into account the constraints of current networks/traffic such as their huge size and their dynamics. COATI thus puts a significant research effort in the following directions:

- **Service Function Chains (SFC):** we study the placement of Service Function Chains within the network considering the ordering constraints. Then, we focus firstly on energy efficiency and secondly on reliability and protection mechanisms. In a last step, we study reconfiguration of the SFCs in case of dynamic traffic with a make-before-break approach.
- **Larger networks:** Another challenge one has to face is the increase in size of practical instances. It is already difficult, if not impossible, to solve practical instances optimally using existing tools. Therefore, we have to find new ways to solve problems using reduction and decomposition methods, characterization of polynomial instances (which are surprisingly often the practical ones), or algorithms with acceptable practical performances.
- **Stochastic behaviors:** Larger topologies mean frequent changes due to traffic and radio fluctuations, failures, maintenance operations, growth, routing policy changes, etc. We aim at including these stochastic behaviors in our combinatorial optimization process to handle the dynamics of the system and to obtain robust designs of networks.

The methods and tools used in our studies come from discrete mathematics and combinatorial optimization, and COATI contributes to their improvements. Also, COATI works on graph-decomposition methods and various games on graphs which are essential for a better understanding of the structural and combinatorial properties of the problems, but also for the design of efficient exact or approximate algorithms. We contribute to the modelling of optimization problems in terms of graphs, study the complexity of the problems, and then we investigate the structural or metric properties of graphs that make these problems hard or easy. We exploit these properties in the design of algorithms in order to find the most efficient ways for solving the problems.

COATI also focuses on the theory of *directed graphs*. Indeed, graph theory can be roughly partitioned into two branches: the areas of undirected graphs and directed graphs. Even though both areas have numerous important applications, for various reasons, undirected graphs have been studied much more extensively than directed graphs. It is worth noticing that many telecommunication problems are modelled with directed graphs. Therefore, a deeper understanding of the theory of directed graphs will benefit to the resolution of telecommunication networks problems. For instance, the problem of finding disjoint paths becomes much more difficult in directed graphs and understanding the underlying structures of actual directed networks would help us to propose solutions.

Last, we have recently started investigating how tools from multi-agents based systems and machine learning theory could help solving some optimization problems in networks. The arrival of Emanuele Natale as a Junior Researcher (CNRS) in the team and of two new PhD students (Francesco D'Amore and Hicham Lesfari) will foster these investigations.

## DANTE Project-Team

### 3. Research Program

#### 3.1. Graph-based signal processing

**Participants:** Paulo Gonçalves, Rémi Gribonval, Marion Foare, Márton Karsai.

**Evolving networks can be regarded as "out of equilibrium" systems.** Indeed, their dynamics are typically characterized by non standard and intricate statistical properties, such as non-stationarity, long range memory effects, intricate space and time correlations.

Analyzing, modeling, and even defining adapted concepts for dynamic graphs is at the heart of DANTE. This is a largely open question that has to be answered by keeping a balance between specificity (solutions triggered by specific data sets) and generality (universal approaches disconnected from social realities). We will tackle this challenge from a graph-based signal processing perspective involving signal analysts and computer scientists, together with experts of the data domain application. One can distinguish two different issues in this challenge, one related to the graph-based organization of the data and the other to the time dependency that naturally exists in the dynamic graph object. In both cases, a number of contributions can be found in the literature, albeit in different contexts. In our application domain, high-dimensional data "naturally reside" on the vertices of weighted graphs. The emerging field of signal processing on graphs merges algebraic and spectral graph theoretic concepts with computational harmonic analysis to process such signals on graphs [74].

As for the first point, adapting well-founded signal processing techniques to data represented as graphs is an emerging, yet quickly developing field which has already received key contributions. Some of them are very general and delineate ambitious programs aimed at defining universal, generally unsupervised methods for exploring high-dimensional data sets and processing them. This is the case for instance of the "diffusion wavelets" and "diffusion maps" pushed forward at Yale and Duke [57]. Others are more traditionally connected with standard signal processing concepts, in the spirit of elaborating new methodologies via some bridging between networks and time series, see for instance [69] and references therein. Other viewpoints can be found as well, including multi-resolution Markov models [77], Bayesian networks or distributed processing over sensor networks [68]. Such approaches can be particularly successful for handling static graphs and unveiling aspects of their organization in terms of dependencies between nodes, grouping, etc. Incorporating possible time dependencies within the whole picture calls however for the addition of an extra dimension to the problem "as it would be the case when switching from one image to a video sequence", a situation for which one can imagine to take advantage of the whole body of knowledge attached to non-stationary signal processing [58].

The arrival of Rémi Gribonval in August 2019 brought a new dimension to the research program of this theme. Specialist of parsimonious representations of large data sets, R. Gribonval will develop at Dante a specific activity related to the sparsification of resources (computing and storage but also regarding the data volumes) in the context of machine and deep learning. This new orientation of Dante will be elaborated and fully integrated to the objectives of the future Inria project that will be proposed after Dante.

#### 3.2. Theory and Structure of dynamic Networks

**Participants:** Christophe Crespelle, Anthony Busson, Márton Karsai, Éric Guichard.

**Characterization of the dynamics of complex networks.** We need to focus on intrinsic properties of evolving/dynamic complex networks. New notions (as opposed to classical static graph properties) have to be introduced: rate of vertices or links appearances or disappearances, the duration of link presences or absences. Moreover, more specific properties related to the dynamics have to be defined and are somehow related to the way to model a dynamic graph.



Through the systematic analysis and characterization of static network representations of many different systems, researchers of several disciplines have unveiled complex topologies and heterogeneous structures, with connectivity patterns statistically characterized by heavy-tails and large fluctuations, scale-free properties and non trivial correlations such as high clustering and hierarchical ordering [71]. A large amount of work has been devoted to the development of new tools for statistical characterisation and modelling of networks, in order to identify their most relevant properties, and to understand which growth mechanisms could lead to these properties. Most of those contributions have focused on static graphs or on dynamic process (*e.g.* diffusion) occurring on static graphs. This has called forth a major effort in developing the methodology to characterize the topology and temporal behaviour of complex networks [71], [62], [78], [67], to describe the observed structural and temporal heterogeneities [55], [62], [56], to detect and measure emerging community structures [59], [75], [76], to see how the functionality of networks determines their evolving structure [66], and to determine what kinds of correlations play a role in their dynamics [63], [65], [70].

The challenge is now to extend this kind of statistical characterization to dynamical graphs. In other words, links in dynamic networks are temporal events, called contacts, which can be either punctual or last for some period of time. Because of the complexity of this analysis, the temporal dimension of the network is often ignored or only roughly considered. Therefore, fully taking into account the dynamics of the links into a network is a crucial and highly challenging issue.

Another powerful approach to model time-varying graphs is via activity driven network models. In this case, the only assumption relates to the distribution of activity rates of interacting entities. The activity rate is realistically broadly distributed and refers to the probability that an entity becomes active and creates a connection with another entity within a unit time step [73]. Even the generic model is already capable to recover some realistic features of the emerging graph, its main advantage is to provide a general framework to study various types of correlations present in real temporal networks. By synthesising such correlations (*e.g.* memory effects, preferential attachment, triangular closing mechanisms, ...) from the real data, we are able to extend the general mechanism and build a temporal network model, which shows certain realistic feature in a controlled way. This can be used to study the effect of selected correlations on the evolution of the emerging structure [64] and its co-evolution with ongoing processes like spreading phenomena, synchronisation, evolution of consensus, random walk etc. [64], [72]. This approach allows also to develop control and immunisation strategies by fully considering the temporal nature of the backgrounding network.

### 3.3. Distributed Algorithms for dynamic networks: regulation, adaptation and interaction

**Participants:** Thomas Begin, Anthony Busson, Isabelle Guérin Lassous, Philippe Nain.

**Dedicated algorithms for dynamic networks.** First, the dynamic network object itself trigger original algorithmic questions. It mainly concerns distributed algorithms that should be designed and deployed to efficiently measure the object itself and get an accurate view of its dynamic behavior. Such distributed measure should be “transparent”, that is, it should introduce no bias or at least a bias that is controllable and corrigible. Such problem is encountered in all distributed metrology measures / distributed probes: P2P, sensor network, wireless network, QoS routing... This question raises naturally the intrinsic notion of adaptation and control of the dynamic network itself since it appears that autonomous networks and traffic aware routing are becoming crucial.

Communication networks are dynamic networks that potentially undergo high dynamicity. The dynamicity exhibited by these networks results from several factors including, for instance, changes in the topology and varying workload conditions. Although most implemented protocols and existing solutions in the literature can cope with a dynamic behavior, the evolution of their behavior operates identically whatever the actual properties of the dynamicity. For instance, parameters of the routing protocols (*e.g.* hello packets transmission frequency) or routing methods (*e.g.* reactive / proactive) are commonly hold constant regardless of the nodes mobility. Similarly, the algorithms ruling CSMA/CA (*e.g.* size of the contention window) are tuned identically and they do not change according to the actual workload and observed topology.

Dynamicity in computer networks tends to affect a large number of performance parameters (if not all) coming from various layers (viz. physical, link, routing and transport). To find out which ones matter the most for our intended purpose, we expect to rely on the tools developed by the two former axes. These quantities should capture and characterize the actual network dynamicity. Our goal is to take advantage of this latter information in order to refine existing protocols, or even to propose new solutions. More precisely, we will attempt to associate “fundamental” changes occurring in the underlying graph of a network (reported through graph-based signal tools) to quantitative performance that are matter of interests for networking applications and the end-users. We expect to rely on available testbeds such as SensLab and FIT to experiment our solutions and ultimately validate our approach.

## DIANA Project-Team

### 3. Research Program

#### 3.1. Service Transparency

Transparency is to provide network users and application developers with reliable information about the current or predicted quality of their communication services, and about potential leakages of personal information, or of other information related to societal interests of the user as a “connected citizen” (e.g. possible violation of network neutrality, opinion manipulation). Service transparency therefore means to provide information meaningful to users and application developers, such as quality of experience, privacy leakages, or opinion manipulation, etc. rather than network-level metrics such as available bandwidth, loss rate, delay or jitter.

The Internet is built around a best effort routing service that does not provide any guarantee to end users in terms of quality of service (QoS). The simplicity of the Internet routing service is at the root of its huge success. Unfortunately, a simple service means unpredicted quality at the access. Even though a considerable effort is done by operators and content providers to optimise the Internet content delivery chain, mainly by over-provisioning and sophisticated engineering techniques, service degradation is still part of the Internet. The proliferation of wireless and mobile access technologies, and the versatile nature of Internet traffic, make end users quality of experience (QoE) forecast even harder. As a matter of fact, the Internet is missing a dedicated measurement plane that informs the end users on the quality they obtain and in case of substantial service degradation, on the origin of this degradation. Current state of the art activities are devoted to building a distributed measurement infrastructure to perform active, passive and hybrid measurements in the wired Internet. However, the problem is exacerbated with modern terminals such as smartphones or tablets that do not facilitate the task for end users (they even make it harder) as they focus on simplifying the interface and limiting the control on the network, whereas the Internet behind is still the same in terms of the quality it provides. Interestingly, this same observation explains the existing difficulty to detect and prevent privacy leaks. We argue that the lack of transparency for diagnosing QoE and for detecting privacy leaks have the same root causes and can be solved using common primitives. For instance, in both cases, it is important to be able to link data packets to an application. Indeed, as the network can only access data packets, there must be a way to bind these packets to an application (to understand users QoE for this application or to associate a privacy leak to an application). This is however a complex task as the traffic might be obfuscated or encrypted. Our objectives in the research direction are the following:

- Design and develop measurement tools providing transparency, in spite of current complexity
- Deploy those measurement tools at the Internet’s edge and make them useful for end users
- Propose measurements plane as an overlay or by exploiting in-network functionalities
- Adapt measurements techniques to network architectural change
- Provide measurements as native functionality in future network architecture

#### 3.2. Open network architecture

We are surrounded by personal content of all types: photos, videos, documents, etc. The volume of such content is increasing at a fast rate, and at the same time, the spread of such content among all our connected devices (mobiles, storage devices, set-top boxes, etc) is also increasing. All this complicates the control of personal content by the user both in terms of access and sharing with other users. The access of the personal content in a seamless way independently of its location is a key challenge for the future of networks. Proprietary solutions exist, but apart from fully depending on one of them, there is no standard plane in the Internet for a seamless access to personal content. Therefore, providing network architectural support to design and develop content access and sharing mechanisms is crucial to allow users control their own data over heterogeneous underlying network or cloud services.

On the other hand, privacy is a growing concern for states, administrations, and companies. Indeed, for instance the French CNIL (entity in charge of citizens privacy in computer systems) puts privacy at the core of its activities by defining rules on any stored and collected private data. Also, companies start to use privacy preserving solutions as a competitive advantage. Therefore, understanding privacy leaks and preventing them is a problem that can already find support. However, all end-users do not *currently* put privacy as their first concern. Indeed, in face of two services with one of higher quality, they usually prefer the highest quality one whatever the privacy implication. This was, for instance, the case concerning the Web search service of Google that is more accurate but less privacy preserving than Bing or Qwant. This is also the case for cloud services such as iCloud or Dropbox that are much more convenient than open source solutions, but very bad in terms of privacy. Therefore, to reach end-users, any privacy preserving solutions must offer a service equivalent to the best existing services.

We consider that it will be highly desirable for Internet users to be able to *easily* move their content from a provider to another and therefore not to depend on a content provider or a social network monopoly. This requires that the network provides built-in architectural support for content networking.

In this research direction, we will define a new *service abstraction layer* (SAL) that could become the new waist of the network architecture with network functionalities below (IP, SDN, cloud) and applications on top. SAL will define different services that are of use to all Internet users for accessing and sharing data (seamless content localisation and retrieval, privacy leakage protection, transparent vertical and horizontal handover, etc.). The biggest challenge here is to cope in the same time with large number of content applications requirements and high underlying networks heterogeneity while still providing efficient applications performance. This requires careful definition of the services primitives and the parameters to be exchanged through the service abstraction layer.

Two concurring factors make the concept behind SAL feasible and relevant today. First, the notion of scalable network virtualization that is a required feature to deploy SAL in real networks today has been discussed recently only. Second, the need for new services abstraction is recent. Indeed, more than fifteen years ago the Internet for the end-users was mostly the Web. Only ten years ago smartphones came into the picture of the Internet boosting the number of applications with new functionalities and risks. Since a few years, many discussions in the network communities took place around the actual complexity of the Internet and the difficulty to develop applications. Many different approaches have been discussed (such as CCN, SDN) that intend to solve only part of the complexity. SAL takes a broader architectural look at the problem and considers solutions such as CCN as mere use cases. Our objectives in this research direction include the following:

- Identify common key networking services required for content access and sharing
- Detect and prevent privacy leaks for content communication
- Enhance software defined networks for large scale heterogeneous environments
- Design and develop open Content Networking architecture
- Define a service abstraction layer as the thin waist for the future content network architecture
- Test and deploy different applications using SAL primitives on heterogeneous network technologies

### 3.3. Methodology

We follow an experimental approach that can be described in the following techniques:

- **Measurements:** the aim is to get a better view of a problem in quantifiable terms. Depending on the field of interest, this may involve large scale distributed systems crawling tools; active probing techniques to infer the status and properties of a complex and non controllable system as the Internet; or even crowdsourcing-based deployments for gathering data on real-users environments or behaviours.
- **Experimental evaluation:** once a new idea has been designed and implemented, it is of course very desirable to assess and quantify how effective it can be, before being able to deploy it on any realistic scale. This is why a wide range of techniques can be considered for getting early, yet as significant as possible, feedback on a given paradigm or implementation. The spectrum for such techniques span from simulations to real deployments in protected and/or controlled environments.

## DIONYSOS Project-Team

### 3. Research Program

#### 3.1. Introduction

The scientific foundations of our work are those of network design and network analysis. Specifically, this concerns the principles of packet switching and in particular of IP networks (protocol design, protocol testing, routing, scheduling techniques), and the mathematical and algorithmic aspects of the associated problems, on which our methods and tools are based.

These foundations are described in the following paragraphs. We begin by a subsection dedicated to Quality of Service (QoS) and Quality of Experience (QoE), since they can be seen as unifying concepts in our activities. Then we briefly describe the specific sub-area of model evaluation and about the particular multidisciplinary domain of network economics.

#### 3.2. Quality of Service and Quality of Experience

Since it is difficult to develop as many communication solutions as possible applications, the scientific and technological communities aim towards providing general *services* allowing to give to each application or user a set of properties nowadays called “Quality of Service” (QoS), a terminology lacking a precise definition. This QoS concept takes different forms according to the type of communication service and the aspects which matter for a given application: for performance it comes through specific metrics (delays, jitter, throughput, etc.), for dependability it also comes through appropriate metrics: reliability, availability, or vulnerability, in the case for instance of WAN (Wide Area Network) topologies, etc.

QoS is at the heart of our research activities: We look for methods to obtain specific “levels” of QoS and for techniques to evaluate the associated metrics. Our ultimate goal is to provide tools (mathematical tools and/or algorithms, under appropriate software “containers” or not) allowing users and/or applications to attain specific levels of QoS, or to improve the provided QoS, if we think of a particular system, with an optimal use of the resources available. Obtaining a good QoS level is a very general objective. It leads to many different areas, depending on the systems, applications and specific goals being considered. Our team works on several of these areas. We also investigate the impact of network QoS on multimedia payloads to reduce the impact of congestion.

Some important aspects of the behavior of modern communication systems have subjective components: the quality of a video stream or an audio signal, *as perceived by the user*, is related to some of the previous mentioned parameters (packet loss, delays, ...) but in an extremely complex way. We are interested in analyzing these types of flows from this user-oriented point of view. We focus on the *user perceived quality*, in short, PQ, the main component of what is nowadays called Quality of Experience (in short, QoE), to underline the fact that, in this case, we want to center the analysis on the user. In this context, we have a global project called PSQA, which stands for Pseudo-Subjective Quality Assessment, and which refers to a technology we have developed allowing to automatically measure this PQ.

Another special case to which we devote research efforts in the team is the analysis of qualitative properties related to interoperability assessment. This refers to the act of determining if end-to-end functionality between at least two communicating systems is as required by the base standards for those systems. Conformance is the act of determining to what extent a single component conforms to the individual requirements of the standard it is based on. Our purpose is to provide such a formal framework (methods, algorithms and tools) for interoperability assessment, in order to help in obtaining efficient interoperability test suites for new generation networks, mainly around IPv6-related protocols. The interoperability test suites generation is based on specifications (standards and/or RFCs) of network components and protocols to be tested.

### 3.3. Stochastic modeling

The scientific foundations of our modeling activities are composed of stochastic processes theory and, in particular, Markov processes, queuing theory, stochastic graphs theory, etc. The objectives are either to develop numerical solutions, or analytical ones, or possibly discrete event simulation or Monte Carlo (and Quasi-Monte Carlo) techniques. We are always interested in model evaluation techniques for dependability and performability analysis, both in static (network reliability) and dynamic contexts (depending on the fact that time plays an explicit role in the analysis or not). We look at systems from the classical so-called *call level*, leading to standard models (for instance, queues or networks of queues) and also at the *burst level*, leading to *fluid models*.

In recent years, our work on the design of the topologies of WANs led us to explore optimization techniques, in particular in the case of very large optimization problems, usually formulated in terms of graphs. The associated methods we are interested in are composed of simulated annealing, genetic algorithms, TABU search, etc. For the time being, we have obtained our best results with GRASP techniques.

Network pricing is a good example of a multi-disciplinary research activity half-way between applied mathematics, economy and networking, centered on stochastic modeling issues. Indeed, the Internet is facing a tremendous increase of its traffic volume. As a consequence, real users complain that large data transfers take too long, without any possibility to improve this by themselves (by paying more, for instance). A possible solution to cope with congestion is to increase the link capacities; however, many authors consider that this is not a viable solution as the network must respond to an increasing demand (and experience has shown that demand of bandwidth has always been ahead of supply), especially now that the Internet is becoming a commercial network. Furthermore, incentives for a fair utilization between customers are not included in the current Internet. For these reasons, it has been suggested that the current flat-rate fees, where customers pay a subscription and obtain an unlimited usage, should be replaced by usage-based fees. Besides, the future Internet will carry heterogeneous flows such as video, voice, email, web, file transfers and remote login among others. Each of these applications requires a different level of QoS: for example, video needs very small delays and packet losses, voice requires small delays but can afford some packet losses, email can afford delay (within a given bound) while file transfer needs a good average throughput and remote login requires small round-trip times. Some pricing incentives should exist so that each user does not always choose the best QoS for her application and so that the final result is a fair utilization of the bandwidth. On the other hand, we need to be aware of the trade-off between engineering efficiency and economic efficiency; for example, traffic measurements can help in improving the management of the network but is a costly option. These are some of the various aspects often present in the pricing problems we address in our work. More recently, we have switched to the more general field of network economics, dealing with the economic behavior of users, service providers and content providers, as well as their relations.

## **DYOGENE Project-Team**

### **3. Research Program**

#### **3.1. Initial research axes**

The following research axes have been defined in 2013 when the project-team was created.

- Algorithms for network performance analysis, led by A. Bouillard and A. Busic.
- Stochastic geometry and information theory for wireless network, led by F. Baccelli and B. Błaszczyszyn.
- The cavity method for network algorithms, led by M. Lelarge.

Our scientific interests keep evolving. Research areas which received the most of our attention in 2019 are summarized in the following sections.

#### **3.2. Distributed network control and smart-grids**

Theory and algorithms for distributed control of networks with applications to the stabilization of power grids subject to high volatility of renewable energy production are being developed by A. Busic in collaboration with Sean Meyn [Prof. at University of Florida and Inria International Chair].

#### **3.3. Mathematics of wireless cellular networks**

A comprehensive approach involving information theory, queueing and stochastic geometry to model and analyze the performance of large cellular networks, validated and implemented by Orange is being led by B. Błaszczyszyn in collaboration with F. Baccelli and M. K. Karray [Orange Labs]. A new collaboration between the Standardization and Research Lab at Nokia Bell Labs and ERC NEMO led by F. Baccelli has been started in 2019.

#### **3.4. High-dimensional statistical inference and distributed learning**

We computed information theoretic bounds for unsupervised and semi-supervised learning and proved complexity bounds for distributed optimization of convex functions using a network of computing units.

#### **3.5. Stochastic Geometry**

In collaboration with Mir-Omid Haji-Mirsadeghi [Sharif University, Tehran, Iran] and Ali Khezeli [School of Mathematical Sciences, Tehran, Iran] F. Baccelli develops a theory of unimodular random metric spaces.

The distortion properties of unconstrained one-bit compression were analyzed by F. Baccelli in collaboration with E. O'Reilly [Caltech] using high dimensional hyperplane tessellations.

In collaboration with D. Yogeshwaran [Indian Statistical Institute, Bangalore] and J. E. Yukich [Lehigh University] B. Błaszczyszyn develops a limit theory (Laws of Large Numbers and Central Limit Theorems) for functionals of spatially correlated point processes.



## EVA Project-Team

### 3. Research Program

#### 3.1. Pitch

##### Designing Tomorrow's Internet of (Important) Things

Inria-EVA is a leading research team in low-power wireless communications. The team pushes the limits of low-power wireless mesh networking by applying them to critical applications such as industrial control loops, with harsh reliability, scalability, security and energy constraints. Grounded in real-world use cases and experimentation, EVA co-chairs the IETF 6TiSCH and LAKE standardization working groups, co-leads Berkeley's OpenWSN project and works extensively with Analog Devices' SmartMesh IP networks. Inria-EVA is the birthplace of the Wattson Elements startup and the Falco solution. The team is associated with Prof. Glaser's (UC Berkeley) and Prof. Kerkez (U. Michigan) through the REALMS associate research team, and with OpenMote through a long-standing Memorandum of Understanding.

#### 3.2. Physical Layer

We study how advanced physical layers can be used in low-power wireless networks. For instance, collaborative techniques such as multiple antennas (e.g. Massive MIMO technology) can improve communication efficiency. The core idea is to use massive network densification by drastically increasing the number of sensors in a given area in a Time Division Duplex (TDD) mode with time reversal. The first period allows the sensors to estimate the channel state and, after time reversal, the second period is to transmit the data sensed. Other techniques, such as interference cancellation, are also possible.

#### 3.3. Wireless Access

Medium sharing in wireless systems has received substantial attention throughout the last decade. HiPERCOM2 has provided models to compare TDMA and CSMA. HiPERCOM2 has also studied how network nodes must be positioned to optimize the global throughput.

EVA pursues modeling tasks to compare access protocols, including multi-carrier access, adaptive CSMA (particularly in VANETs), as well as directional and multiple antennas. There is a strong need for determinism in industrial networks. The EVA team focuses particularly on scheduled medium access in the context of deterministic industrial networks; this involves optimizing the joint time slot and channel assignment. Distributed approaches are considered, and the EVA team determines their limits in terms of reliability, latency and throughput. Furthermore, adaptivity to application or environment changes are taken into account.

#### 3.4. Coexistence of Wireless Technologies

Wireless technologies such as cellular, low-power mesh networks, (Low-Power) WiFi, and Bluetooth (low-energy) can reasonably claim to fit the requirements of the IoT. Each, however, uses different trade-offs between reliability, energy consumption and throughput. The EVA team studies the limits of each technology, and will develop clear criteria to evaluate which technology is best suited to a particular set of constraints.

Coexistence between these different technologies (or different deployments of the same technology in a common radio space) is a valid point of concern.

The EVA team aims at studying such coexistence, and, where necessary, propose techniques to improve it. Where applicable, the techniques will be put forward for standardization. Multiple technologies can also function in a symbiotic way.

For example, to improve the quality of experience provided to end users, a wireless mesh network can transport sensor and actuator data in place of a cellular network, when and where cellular connectivity is poor.

The EVA team studies how and when different technologies can complement one another. A specific example of a collaborative approach is Cognitive Radio Sensor Networks (CRSN).

### 3.5. Energy-Efficiency and Determinism

Reducing the energy consumption of low-power wireless devices remains a challenging task. The overall energy budget of a system can be reduced by using less power-hungry chips, and significant research is being done in that direction. That being said, power consumption is mostly influenced by the algorithms and protocols used in low-power wireless devices, since they influence the duty-cycle of the radio.

EVA will search for energy-efficient mechanisms in low-power wireless networks. One new requirement concerns the ability to predict energy consumption with a high degree of accuracy. Scheduled communication, such as the one used in the IEEE 802.15.4 TSCH (Time Slotted CHannel Hopping) standard, and by IETF 6TiSCH, allows for a very accurate prediction of the energy consumption of a chip. Power conservation will be a key issue in EVA.

To tackle this issue and match link-layer resources to application needs, EVA's 5-year research program dealing with Energy-Efficiency and Determinism centers around 3 studies:

- Performance Bounds of a TSCH network. We propose to study a low-power wireless TSCH network as a Networked Control System (NCS), and use results from the NCS literature. A large number of publications on NCS, although dealing with wireless systems, consider wireless links to have perfect reliability, and do not consider packet loss. Results from these papers can not therefore be applied directly to TSCH networks. Instead of following a purely mathematical approach to model the network, we propose to use a non-conventional approach and build an empirical model of a TSCH network.
- Distributed Scheduling in TSCH networks. Distributed scheduling is attractive due to its scalability and reactivity, but might result in a sub-optimal schedule. We continue this research by designing a distributed solution based on control theory, and verify how this solution can satisfy service level agreements in a dynamic environment.

### 3.6. Network Deployment

Since sensor networks are very often built to monitor geographical areas, sensor deployment is a key issue. The deployment of the network must ensure full/partial, permanent/intermittent coverage and connectivity. This technical issue leads to geometrical problems which are unusual in the networking domain.

We can identify two scenarios. In the first one, sensors are deployed over a given area to guarantee full coverage and connectivity, while minimizing the number of sensor nodes. In the second one, a network is re-deployed to improve its performance, possibly by increasing the number of points of interest covered, and by ensuring connectivity. EVA will investigate these two scenarios, as well as centralized and distributed approaches. The work starts with simple 2D models and will be enriched to take into account more realistic environment: obstacles, walls, 3D, fading.

### 3.7. Data Gathering and Dissemination

A large number of WSN applications mostly do data gathering (a.k.a "convergecast"). These applications usually require small delays for the data to reach the gateway node, requiring time consistency across gathered data. This time consistency is usually achieved by a short gathering period.

In many real WSN deployments, the channel used by the WSN usually encounters perturbations such as jamming, external interferences or noise caused by external sources (e.g. a polluting source such as a radar) or other coexisting wireless networks (e.g. WiFi, Bluetooth). Commercial sensor nodes can communicate on multiple frequencies as specified in the IEEE 802.15.4 standard. This reality has given birth to the multichannel communication paradigm in WSNs.

Multichannel WSNs significantly expand the capability of single-channel WSNs by allowing parallel transmissions, and avoiding congestion on channels or performance degradation caused by interfering devices.

In EVA, we will focus on raw data convergecast in multichannel low-power wireless networks. In this context, we are interested in centralized/distributed algorithms that jointly optimize the channel and time slot assignment used in a data gathering frame. The limits in terms of reliability, latency and bandwidth will be evaluated. Adaptivity to additional traffic demands will be improved.

### **3.8. Self-Learning Networks**

To adapt to varying conditions in the environment and application requirements, the EVA team investigate self-learning networks. Machine learning approaches, based on experts and forecasters, are investigated to predict the quality of the wireless links in a WSN. This allows the routing protocol to avoid using links exhibiting poor quality and to change the route before a link failure. Additional applications include where to place the aggregation function in data gathering. In a content delivery network (CDN), it is very useful to predict popularity, expressed by the number of requests per day, for a multimedia content. The most popular contents are cached near the end-users to maximize the hit ratio of end-users' requests. Thus the satisfaction degree of end-users is maximized and the network overhead is minimized.

### **3.9. Internet of Things Security**

Existing Internet threats might steal our digital information. Tomorrow's threats could disrupt power plants, home security systems, hospitals. The Internet of Things is bridging our digital security with personal safety. Popular magazines are full of stories of hacked devices (e.g. drone attack on Philips Hue), IoT botnets (e.g. Mirai), and inherent insecurity.

*Why has the IoT industry failed to adopt the available computer security techniques and best practices?* Our experience from research, industry collaborations, and the standards bodies has shown that the main challenges are:

1. The circumvention of the available technical solutions due to their inefficiency.
2. The lack of a user interface for configuring the product in the field resulting in default parameters being (re)used.
3. Poorly tested software, often lacking secure software upgrade mechanisms.

Our research goal is to contribute to a more secure IoT, by proposing technical solutions to these challenges for low-end IoT devices with immediate industrial applicability and transfer potential. We complement the existing techniques with the missing pieces to move towards truly usable and secure IoT systems.

## FUN Project-Team

### 3. Research Program

#### 3.1. Introduction

We will focus on wireless ubiquitous networks that rely on constrained devices, i.e. with limited resources in terms of storage and computing capacities. They can be sensors, small robots, RFID readers or tags. A wireless sensor retrieves a physical measure such as light. A wireless robot is a wireless sensor that in addition has the ability to move by itself in a controlled way. A drone is a robot with the ability to manoeuvre in 3D (in the air or in the water). RFID tags are passive items that embed a unique identifier for a place or an object allowing accurate traceability. They can communicate only in the vicinity of an RFID reader. An RFID reader can be seen as a special kind of sensor in the network which data is the one read on tags. These devices may run on batteries that are not envisaged to be changed or recharged. These networks may be composed of ten to thousands of such heterogeneous devices for which energy is a key issue.

Today, most of these networks are homogeneous, i.e. composed of only one kind of devices. They have mainly been studied in application and technology silos. Because of this, they are approaching fundamental limitations especially in terms of topology deployment, management and communications, while exploiting the complementarity of heterogeneous devices and communication technologies would enlarge their capacities and the set of applications. Finally, these networks must work efficiently even in dynamic and realistic situations, i.e. they must consider by design the different dynamic parameters and automatically self-adapt to their variations.

Our overall goal is represented by Figure 1 . We will investigate wireless ubiquitous IoT services for constrained devices by smartly combining **different frequency bands** and **different medium access and routing techniques** over **heterogeneous devices** in a **distributed** and **opportunistic** fashion. Our approach will always deal with **hardware constraints** and take care of **security** and **energy** issues to provide protocols that ride on **synergy** and **self-organization** between devices.

*The goal of the FUN project team is to provide these next generation networks with a set of innovative and distributed self-organizing cooperative protocols to raise them to a new level of scalability, autonomy, adaptability, manageability and performance. We aim to break these silos to exploit the full synergy between devices, making them cooperate in a single holistic network. We will consider them as networks of heterogeneous devices rather than a collection of heterogeneous networks.*

To realize the full potential of these ubiquitous networks, there is a need to provide them with a set of tools that allow them to (i) (self-)deploy, (ii) self-organize, (iii) discover and locate each other, resources and services and (iv) communicate. These tools will be the basics for enabling cooperation, co-existence and witnessing a global efficient behavior. The deployment of these mechanisms is challenging since it should be achieved in spite of several limitations. The main difficulties are to provide such protocols in a **secured** and **energy-efficient** fashion in spite of:

- dynamic topology changes due to various factors such as the unreliability of the wireless medium, the wireless interferences between devices, node mobility and energy saving mechanisms;
- hardware constraints in terms of CPU and memory capacities that limit the operations and data each node can perform/collect;
- lacks of interoperability between applicative, hardware and technological silos that may prevent from data exchange between different devices.

##### 3.1.1. Objectives and methodology

To reach our overall goal, we will pursue the two following objectives. These two objectives are orthogonal and can be carried on jointly:

1. Providing realistic complete self-organizing tools *e.g. vertical perspective*.
2. Going to heterogeneous energy-efficient performing wireless networks *e.g. horizontal perspective*.

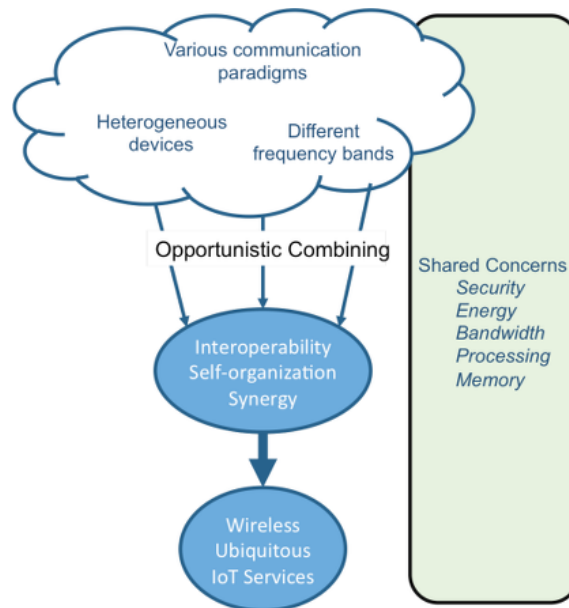


Figure 1. FUN's overall goal.

We give more details on these two objectives below. To achieve our main objectives, we will mainly apply the methodology depicted in Figure 2 combining both theoretical analysis and experimental validation. Mathematical tools will allow us to properly dimension a problem, formally define its limitations and needs to provide suitable protocols in response. Then, they will allow us to qualify the outcome solutions before we validate and stress them in real scenarios with regards to applications requirements. For this, we will realize proofs-of-concept with real scenarios and real devices. Differences between results and expectations will be analyzed in return in order to well understand them and integrate them by design for a better protocol self-adaptation capability.

### 3.2. Vertical Perspective

As mentioned, future ubiquitous networks evolve in dynamic and unpredictable environments. Also, they can be used in a large scope of applications that have several expectations in terms of performance and different contextual limitations. In this heterogeneous context, IoT devices must support multiple applications and relay traffic with non-deterministic pattern.

To make our solutions practical and efficient in real conditions, we will adopt the dual approach both *top-down* and *bottom-up*. The *top-down* approach will ensure that we consider the application (such as throughput, delay, energy consumption, etc.) and environmental limitations (such as deployment constraints, etc.). The *bottom-up* approach will ensure that we take account of the physical and hardware characteristics such as memory, CPU, energy capacities but also physical interferences and obstacles. With this integrated perspective, we will be in capacity to design well adapted **cross-layer** integrated protocols [59]. We will design jointly routing and MAC layers by taking dynamics occurring at the physical layer into account with a constant concern for energy and security. We will investigate new adaptive frequency hopping techniques combined with routing protocols [59], [45].

This vision will also allow us to integrate external factors by design in our protocols, in an opportunistic way. Yet, we will leverage on the occurrence of any of these phenomena rather than perceiving them as obstacles

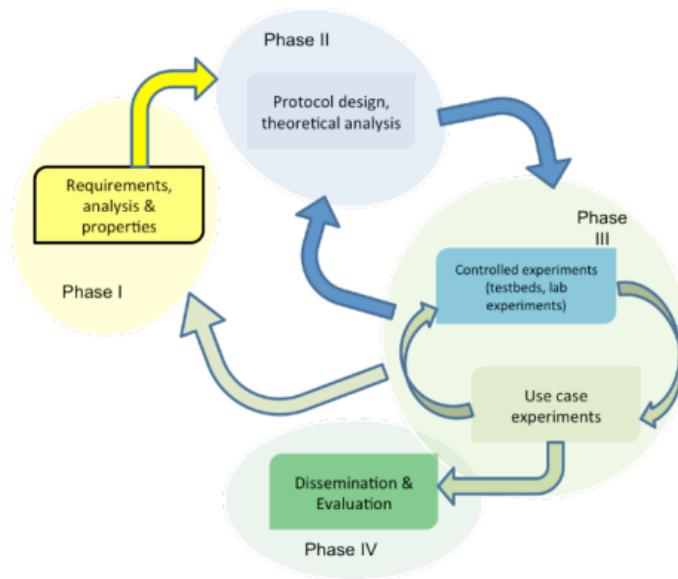


Figure 2. Methodology to be applied in FUN.

or limitations. As an example, we will rely on node undergone mobility to enhance routing performance as we have started to investigate in [55], [37]. On the same idea, when specific features are available like controlled mobility, we will exploit it to improve connectivity or coverage quality like in [49], [57], [31], [25].

### 3.3. Horizontal perspective

We aim at designing efficient tools for a plethora of wireless devices supporting highly heterogeneous technologies. We will thus investigate these networks from a horizontal perspective, e.g. by considering heterogeneity in low level communications layers.

Given the spectrum scarcity, they will probably need to coexist in the same frequency bands and sometimes for different purposes (RFID tag reading may use the same frequency bands as the wireless sensors). One important aspect to consider in this setting is how these different access technologies will interact with each other, and what are the mechanisms needed to be put in place to guarantee that all services obtain the required share of resources when needed. This problem appears in different application domains, ranging from traffic offloading to unlicensed bands by cellular networks and the need to coexist with WiFi and radars, from a scenario in which multiple-purpose IoT clouds coexist in a city [56]. We will thus explore the dynamics of these interactions and devise ways to ensure smooth coexistence while considering the heterogeneity of the devices involved, the access mechanisms used as well as the requirements of the services provided.

To face the spectrum scarcity, we will also investigate new alternative communication paradigms such as phonon-based or light-based communications as we have initiated in [42] and we will work on the coexistence of these technologies with traditional communication techniques, specifically by investigating efficient switching techniques from one communication technology to the other (they were most focused on the security aspects, to prevent jamming attacks). Resilience and reliability of the whole system will be the key factors to be taken into account [43], [39], [18].

As a more prospective activity, we consider exploring software and communication security for IoT. This is challenging given that existing solutions do not address systems that are both constrained and networked [44].

Finally, in order to contribute to a better interoperability between all these technologies, we will continue to contribute to standardization bodies such as IETF and EPC Global.



## **GANG Project-Team**

### **3. Research Program**

#### **3.1. Graph and Combinatorial Algorithms**

We focus on two approaches for designing algorithms for large graphs: decomposing the graph and relying on simple graph traversals.

##### **3.1.1. Graph Search**

We more deeply study multi-sweep graph searches. In this domain a graph search only yields a total ordering of the vertices which can be used by the subsequent graph searches. This technique can be used on huge graphs and does not need extra memory. We have already obtained preliminary results in this direction and many well-known graph algorithms can be put in this framework. The idea behind this approach is that each sweep discovers some structure of the graph. At the end of the process either we have found the underlying structure (for example an interval representation for an interval graph) or an approximation of it (for example in hard discrete optimization problems). We envision applications to exact computations of centers in huge graphs, to underlying combinatorial optimization problems, but also to networks arising in biology.

##### **3.1.2. Graph Decomposition**

In order to summarize a graph into a more compact and more human-readable form, we introduced the hub-laminar decomposition. It is suitable for graphs that are dominated by long isometric cycles or shortest paths, called laminar, which meet only at their extremities, called hubs. Computing this decomposition is NP-hard but a canonical approximation may be computed under some hypotheses on the distances between hubs. It provides a distance labelling for the decomposable graphs. We also investigated the case where the decomposition is reduced to a single cycle, yielding the problem of finding the longest isometric cycle, which is NP-complete and for which a first approximation algorithm was proposed in ENTCS.

##### **3.1.3. Graph Exploration**

In the course of graph exploration, a mobile agent is expected to regularly visit all the nodes of an unknown network, trying to discover all its nodes as quickly as possible. Our research focuses on the design and analysis of agent-based algorithms for exploration-type problems, which operate efficiently in a dynamic network environment, and satisfy imposed constraints on local computational resources, performance, and resilience. Our recent contributions in this area concern the design of fast deterministic algorithms for teams of agents operating in parallel in a graph, with limited or no persistent state information available at nodes. We plan further studies to better understand the impact of memory constraints and of the availability of true randomness on efficiency of the graph exploration process.

#### **3.2. Distributed Computing**

The distributed computing community can be viewed as a union of two sub-communities. This is also true in our team. Although they have interactions, they are disjoint enough not to leverage each other's results. At a high level, one is mostly interested in timing issues (clock drifts, link delays, crashes, etc.) while the other one is mostly interested in spatial issues (network structure, memory requirements, etc.). Indeed, one sub-community is mostly focusing on the combined impact of asynchronism and faults on distributed computation, while the other addresses the impact of network structural properties on distributed computation. Both communities address various forms of computational complexity, through the analysis of different concepts. This includes, e.g., failure detectors and wait-free hierarchy for the former community and compact labeling schemes, and computing with advice for the latter community. We have an ambitious project to achieve the reconciliation between the two communities by focusing on the same class of problems, the yes/no-problems, and establishing the scientific foundations for building up a consistent theory of computability and complexity

for distributed computing. The main question addressed is therefore: is the absence of globally coherent computational complexity theories covering more than fragments of distributed computing, inherent to the field? One issue is obviously the types of problems located at the core of distributed computing. Tasks like consensus, leader election, and broadcasting are of very different nature. They are not *yes-no* problems, neither are they minimization problems. Coloring and Minimal Spanning Tree are optimization problems but we are often more interested in constructing an optimal solution than in verifying the correctness of a given solution. Still, it makes full sense to analyze the *yes-no* problems corresponding to checking the validity of the output of tasks. Another issue is the power of individual computation. The FLP impossibility result as well as Linial's lower bound hold independently of the individual computational power of the involved computing entities. For instance, the individual power of solving NP-hard problems in constant time would not help overcoming these limits, which are inherent to the fact that computation is distributed. A third issue is the abundance of models for distributed computing frameworks, from shared memory to message passing, spanning all kinds of specific network structures (complete graphs, unit-disk graphs, etc.) and/or timing constraints (from complete synchronism to full asynchronism). There are however models, typically the wait-free model and the LOCAL model, which, though they do not claim to reflect accurately real distributed computing systems, enable focusing on some core issues. Our ongoing research program is to carry many important notions of Distributed Computing into a *standard* computational complexity.

### **3.3. Network Algorithms and Analysis**

Based on our scientific expertise in both graph algorithms and distributed algorithms, we plan to analyze the behavior of various networks such as future Internet, social networks, overlay networks resulting from distributed applications or online social networks.

#### **3.3.1. Information Dissemination**

One of the key aspects of networks resides in the dissemination of information among the nodes. We aim at analyzing various procedures of information propagation from dedicated algorithms to simple distributed schemes such as flooding. We also consider various models, e.g. where noise can alter information as it propagates or where memory of nodes is limited.

#### **3.3.2. Routing Paradigms**

We try to explore new routing paradigms such as greedy routing in social networks for example. We are also interested in content centric networking where routing is based on content name rather than content address. One of our targets is multiple path routing: how to design forwarding tables providing multiple disjoint paths to the destination?

#### **3.3.3. Beyond Peer-to-Peer**

Based on our past experience of peer-to-peer application design, we would like to broaden the spectrum of distributed applications where new efficient algorithms can be designed and their analysis can be performed. We especially target online social networks as we see them as collaborative tools for exchanging information. A basic question resides in making the right connections for gathering filtered and accurate information with sufficient coverage.

#### **3.3.4. SAT and Forwarding Information Verification**

As forwarding tables of networks grow and are sometimes manually modified, the problem of verifying them becomes critical and has recently gained interest. Some problems that arise in network verification such as loop detection for example, may be naturally encoded as Boolean Satisfiability problems. Beside theoretical interest in complexity proofs, this encoding allows one to solve these problems by taking advantage of efficient Satisfiability testing solvers. Indeed, SAT solvers have proved to be very efficient in solving problems coming from various areas (Circuit Verification, Dependency and Conflicts in Software distributions...) and encoded in Conjunctive Normal Form. To test an approach using SAT solvers in network verification, one needs to collect data sets from a real network and to develop good models for generating realistic networks. The technique of encoding and the solvers themselves need to be adapted to this kind of problems. All this represents a rich experimental field of future research.

### **3.3.5. Network Analysis**

Finally, we are interested in analyzing the structural properties of practical networks. This can include diameter computation or ranking of nodes. As we mostly consider large networks, we are often interested in efficient heuristics. Ideally, we target heuristics that give exact answers and are reasonably fast in practice although short computation time is not guaranteed for all networks. We have already designed such heuristics for diameter computation; understanding the structural properties that enable short computation time in practice is still an open question.

### **3.3.6. Network Parameter**

Betweenness centrality is a graph parameter that has been successfully applied to network analysis. In the context of computer networks, it was considered for various objectives, ranging from routing to service placement. However, as observed by Maccari et al. [INFOCOM 2018], research on betweenness centrality for improving protocols was hampered by the lack of a usable, fully distributed algorithm for computing this parameter. In [21], we resolved this issue by designing an efficient algorithm for computing betweenness centrality, which can be implemented by minimal modifications to any distance-vector routing protocol based on Bellman-Ford. The convergence time of our implementation is shown to be proportional to the diameter of the network.

## **GANG Project-Team**

### **3. Research Program**

#### **3.1. Graph and Combinatorial Algorithms**

We focus on two approaches for designing algorithms for large graphs: decomposing the graph and relying on simple graph traversals.

##### **3.1.1. Graph Search**

We more deeply study multi-sweep graph searches. In this domain a graph search only yields a total ordering of the vertices which can be used by the subsequent graph searches. This technique can be used on huge graphs and does not need extra memory. We have already obtained preliminary results in this direction and many well-known graph algorithms can be put in this framework. The idea behind this approach is that each sweep discovers some structure of the graph. At the end of the process either we have found the underlying structure (for example an interval representation for an interval graph) or an approximation of it (for example in hard discrete optimization problems). We envision applications to exact computations of centers in huge graphs, to underlying combinatorial optimization problems, but also to networks arising in biology.

##### **3.1.2. Graph Decomposition**

In order to summarize a graph into a more compact and more human-readable form, we introduced the hub-laminar decomposition. It is suitable for graphs that are dominated by long isometric cycles or shortest paths, called laminar, which meet only at their extremities, called hubs. Computing this decomposition is NP-hard but a canonical approximation may be computed under some hypotheses on the distances between hubs. It provides a distance labelling for the decomposable graphs. We also investigated the case where the decomposition is reduced to a single cycle, yielding the problem of finding the longest isometric cycle, which is NP-complete and for which a first approximation algorithm was proposed in ENTCS.

##### **3.1.3. Graph Exploration**

In the course of graph exploration, a mobile agent is expected to regularly visit all the nodes of an unknown network, trying to discover all its nodes as quickly as possible. Our research focuses on the design and analysis of agent-based algorithms for exploration-type problems, which operate efficiently in a dynamic network environment, and satisfy imposed constraints on local computational resources, performance, and resilience. Our recent contributions in this area concern the design of fast deterministic algorithms for teams of agents operating in parallel in a graph, with limited or no persistent state information available at nodes. We plan further studies to better understand the impact of memory constraints and of the availability of true randomness on efficiency of the graph exploration process.

#### **3.2. Distributed Computing**

The distributed computing community can be viewed as a union of two sub-communities. This is also true in our team. Although they have interactions, they are disjoint enough not to leverage each other's results. At a high level, one is mostly interested in timing issues (clock drifts, link delays, crashes, etc.) while the other one is mostly interested in spatial issues (network structure, memory requirements, etc.). Indeed, one sub-community is mostly focusing on the combined impact of asynchronism and faults on distributed computation, while the other addresses the impact of network structural properties on distributed computation. Both communities address various forms of computational complexity, through the analysis of different concepts. This includes, e.g., failure detectors and wait-free hierarchy for the former community and compact labeling schemes, and computing with advice for the latter community. We have an ambitious project to achieve the reconciliation between the two communities by focusing on the same class of problems, the yes/no-problems, and establishing the scientific foundations for building up a consistent theory of computability and complexity

for distributed computing. The main question addressed is therefore: is the absence of globally coherent computational complexity theories covering more than fragments of distributed computing, inherent to the field? One issue is obviously the types of problems located at the core of distributed computing. Tasks like consensus, leader election, and broadcasting are of very different nature. They are not *yes-no* problems, neither are they minimization problems. Coloring and Minimal Spanning Tree are optimization problems but we are often more interested in constructing an optimal solution than in verifying the correctness of a given solution. Still, it makes full sense to analyze the *yes-no* problems corresponding to checking the validity of the output of tasks. Another issue is the power of individual computation. The FLP impossibility result as well as Linial's lower bound hold independently of the individual computational power of the involved computing entities. For instance, the individual power of solving NP-hard problems in constant time would not help overcoming these limits, which are inherent to the fact that computation is distributed. A third issue is the abundance of models for distributed computing frameworks, from shared memory to message passing, spanning all kinds of specific network structures (complete graphs, unit-disk graphs, etc.) and/or timing constraints (from complete synchronism to full asynchronism). There are however models, typically the wait-free model and the LOCAL model, which, though they do not claim to reflect accurately real distributed computing systems, enable focusing on some core issues. Our ongoing research program is to carry many important notions of Distributed Computing into a *standard* computational complexity.

### **3.3. Network Algorithms and Analysis**

Based on our scientific expertise in both graph algorithms and distributed algorithms, we plan to analyze the behavior of various networks such as future Internet, social networks, overlay networks resulting from distributed applications or online social networks.

#### **3.3.1. Information Dissemination**

One of the key aspects of networks resides in the dissemination of information among the nodes. We aim at analyzing various procedures of information propagation from dedicated algorithms to simple distributed schemes such as flooding. We also consider various models, e.g. where noise can alter information as it propagates or where memory of nodes is limited.

#### **3.3.2. Routing Paradigms**

We try to explore new routing paradigms such as greedy routing in social networks for example. We are also interested in content centric networking where routing is based on content name rather than content address. One of our targets is multiple path routing: how to design forwarding tables providing multiple disjoint paths to the destination?

#### **3.3.3. Beyond Peer-to-Peer**

Based on our past experience of peer-to-peer application design, we would like to broaden the spectrum of distributed applications where new efficient algorithms can be designed and their analysis can be performed. We especially target online social networks as we see them as collaborative tools for exchanging information. A basic question resides in making the right connections for gathering filtered and accurate information with sufficient coverage.

#### **3.3.4. SAT and Forwarding Information Verification**

As forwarding tables of networks grow and are sometimes manually modified, the problem of verifying them becomes critical and has recently gained interest. Some problems that arise in network verification such as loop detection for example, may be naturally encoded as Boolean Satisfiability problems. Beside theoretical interest in complexity proofs, this encoding allows one to solve these problems by taking advantage of efficient Satisfiability testing solvers. Indeed, SAT solvers have proved to be very efficient in solving problems coming from various areas (Circuit Verification, Dependency and Conflicts in Software distributions...) and encoded in Conjunctive Normal Form. To test an approach using SAT solvers in network verification, one needs to collect data sets from a real network and to develop good models for generating realistic networks. The technique of encoding and the solvers themselves need to be adapted to this kind of problems. All this represents a rich experimental field of future research.

### **3.3.5. Network Analysis**

Finally, we are interested in analyzing the structural properties of practical networks. This can include diameter computation or ranking of nodes. As we mostly consider large networks, we are often interested in efficient heuristics. Ideally, we target heuristics that give exact answers and are reasonably fast in practice although short computation time is not guaranteed for all networks. We have already designed such heuristics for diameter computation; understanding the structural properties that enable short computation time in practice is still an open question.

### **3.3.6. Network Parameter**

Betweenness centrality is a graph parameter that has been successfully applied to network analysis. In the context of computer networks, it was considered for various objectives, ranging from routing to service placement. However, as observed by Maccari et al. [INFOCOM 2018], research on betweenness centrality for improving protocols was hampered by the lack of a usable, fully distributed algorithm for computing this parameter. In [21], we resolved this issue by designing an efficient algorithm for computing betweenness centrality, which can be implemented by minimal modifications to any distance-vector routing protocol based on Bellman-Ford. The convergence time of our implementation is shown to be proportional to the diameter of the network.

## MARACAS Team

### 3. Research Program

#### 3.1. General description

As presented in the first section, *Computing Networks* is a concept generalizing the study of multi-user systems under the communication perspective. This problematic is partly addressed in the aforementioned references. Optimizing *Computing Networks* relies on exploiting simultaneously multi-user communication capabilities, in the one hand, and storage and computing resources in the other hand. Such optimization needs to cope with various constraints such as energy efficiency or energy harvesting, delays, reliability or network load.

The notion of reliability (used in MARACAS acronym) is central when considered in the most general sense : ultimately, the reliability of a *Computing Network* measures its capability to perform its intended role under some confidence interval. Figure 1 represents the most important performance criteria to be considered to achieve reliable communications. These metrics fit with those considered in 5G and beyond technologies [63].

On the theoretical side, multi-user information theory is a keystone element. It is worth noting that classical information theory focuses on the power-bandwidth tradeoff usually referred as Energy Efficiency-Spectral Efficiency (EE-SE) tradeoff (green arrow on 1 ). However, the other constraints can be efficiently introduced by using a non-asymptotic formulation of the fundamental limits [62], [64] and in association with other tools devoted to the analysis of random processes (queuing theory, ...).

**Maracas aims at studying *Computing Networks* from a communication point of view, using the foundations of information theory in association with other theoretical tools related to estimation theory and probability theory.**

In particular, Maracas combines techniques from communication and information theory with statistical signal processing, control theory, and game theory. Wireless networks is the emblematic application for Maracas, but other scenarios are appealing for us, such as molecular communications, smart grids or smart buildings.

Several teams at Inria are addressing computing networks, but working on this problem with an emphasis on communication aspects is unique within Inria.

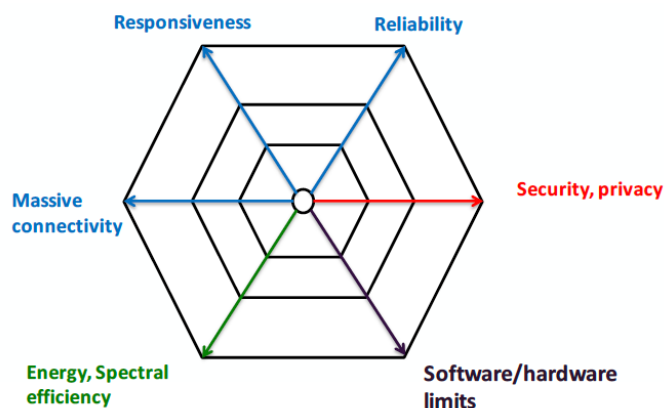


Figure 1. Main metrics for future networks (5G and beyond)



The complexity of *Computing Networks* comes first from the high dimensionality of the problem: i) thousands of nodes, each with up to tens setting parameters and ii) tens variable objective functions to be minimized/maximized.

In addition, the necessary decentralization of the decision process, the non stationary behavior of the network itself (mobility, ON/OFF Switching) and of the data flows, and the necessary reduction of costly feedback and signaling (channel estimation, topology discovering, medium access policies...) are additional features that increase the problem complexity.

**The original positioning of Maracas holds in his capability to address three complementary challenges :**

1. **to develop a sound mathematical framework inspired by information theory.**
2. **to design algorithms, achieving performance close to these limits.**
3. **to test and validate these algorithms on experimental testbeds.**

### 3.2. Research program

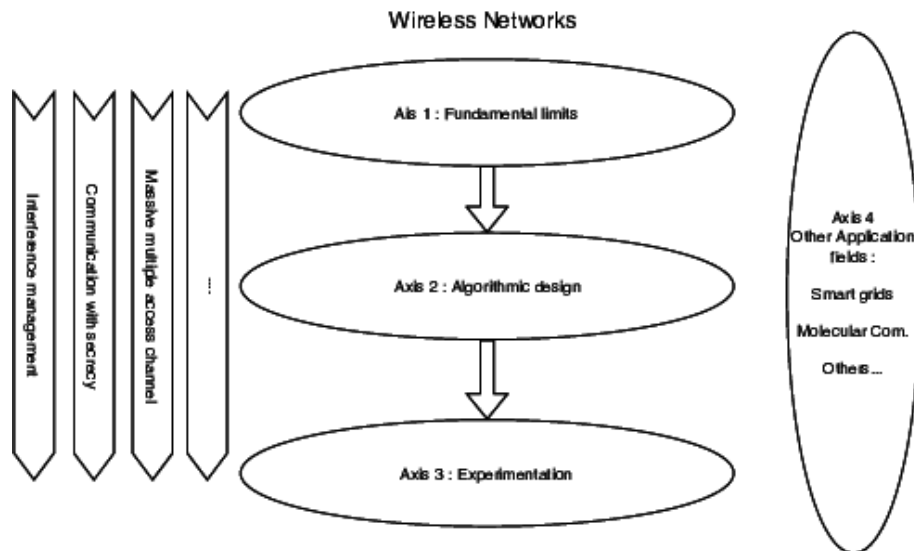


Figure 2. Maracas organization

Our research is organized in 4 research axes:

- **Axis 1 - Fundamental Limits of Reliable Communication Systems:** Information theory is revisited to integrate reliability in the wide sense. The non-asymptotic theory which made progress recently and attracted a lot of interest in the information theory community is a good starting point. But for addressing computing network in a wide sense, it is necessary to go back to the foundation of communication theory and to derive new results, e.g. for non Gaussian channels [8] or for multi-constrained systems [17].

This also means revisiting the fundamental estimation-detection problem [65] in a general multi-criteria, multi-user framework to derive tractable and meaningful bounds.

As mentioned in the introduction, *Computing Networks* also relies on a data-centric vision, where transmission, storage and processing are jointly optimized. The strategy of *caching at the edge* [57] proposed for cellular networks shows the high potential of considering simultaneously data and network properties. Maracas is willing to extend his skills on source coding aspects to tackle with a data-oriented modeling of *Computing Networks*.

- **Axis 2 - Algorithms and protocols:** Our second objective is to elaborate new algorithms and protocols able to achieve or at least to approach the aforementioned fundamental limits. While the exploration of fundamental limits is helpful to determine the most promising strategies (e.g. relaying, cooperation, interference alignment) to increase system performance, the transformation of these degrees of freedom into real protocols is a non trivial issue. One reason is the exponentially growing complexity of multi-user communication strategies, with the number of users, due to the necessity of some coordination, feedback and signaling. The general problem is a decentralized and dynamic multi-agents multi-criteria optimization problem and the general formulation is a non-linear and non-convex large scale problem.

The conventional research direction aims at reducing the complexity by relaxing some constraints or by reducing the number of degrees of freedom. For instance, topology interference management is a seducing model used to reduce feedback needs in decentralized wireless networks leading to original and efficient algorithms [67], [59].

Another emerging research direction relies on using machine learning techniques [54] as a natural evolution of cognitive radio based approaches. Machine learning in the wide sense is not new in radio networks, but the most important works in the past were devoted to reinforcement learning approaches. The use of deep learning (DL) is much more recent, with two important issues : i) identifying the right problems that really need DL algorithms and ii) providing extensive data sets from simulation and real experiments. Our group started to work on this topic in association with Nokia in the joint research lab. As we are not currently expert in deep learning, our primary objective is to identify the strategic problems and to collaborate in the future with Inria experts in DL, and in the long term to contribute not only to the application of these techniques, but also to improve their design according to the constraints of computing networks.

- **Axis 3 - Experimental validation :** With the rapid evolution of network technologies, and their increasing complexity, experimental validation is necessary for two reasons: to get data, and to validate new algorithms on real systems.

Maracas activity leverages on the FIT/CorteXlab platform (<http://www.cortexlab.fr/>), and our strong partnerships with leading industry including Nokia Bell Labs, Orange labs, Sigfox or Sequans. Beyond the platform itself which offers a worldwide unique and remotely accessible testbed , Maracas also develops original experimentations exploiting the reproducibility, the remote accessibility, and the deployment facilities to produce original results at the interface of academic and industrial research [1], [10]. FIT/CorteXlab uses the GNU Radio environment to evaluate new multi-user communication systems.

Our experimental work is developed in collaboration with other Inria teams especially in the Rhone-Alpes centre but also in the context of the future SILECS project <https://www.silecs.net/> which will implement the convergence between FIT and Grid'5000 infrastructures in France, in cooperation with European partners and infrastructures. SILECS is a unique framework which will allow us to test our algorithms, to generate data, as required to develop a data-centric approach for computing networks.

Last but not least, software radio technologies are leaving the confidentiality of research laboratories and are made available to a wide public market with cheap (few euros) programmable equipment, allowing to setup non standard radio systems. The existence of home-made and non official radio systems with legacy ones could prejudice the deployment of Internet of things. Developing efficient algorithms able to detect, analyse and control the spectrum usage is an important issue. Our research on FIT/CorteXlab will contribute to this know-how.

- **Axis 4 - Other application fields** : Even if the wireless network context is still challenging and provides interesting problems, Maracas targets to broaden its exploratory playground from an application perspective. We are looking for new communication systems, or simply other multi-user decentralized systems, for which the theory developed in the context of wireless networks can be useful. Basically, Maracas might address any problem where multi-agents are trying to optimize their common behavior and where the communication performance is critical (e.g. vehicular communications, multi-robots systems, cyberphysical systems). Following this objective, we already studied the problem of missing data recovery in smart grids [11] and the original paradigm of molecular communications [6].

Of course, the objective of this axis is not to address random topics but to exploit our scientific background on new problems, in collaboration with other academic teams or industry. This is a winning strategy to develop new partnerships, in collaboration with other Inria teams.

## NEO Project-Team

### 3. Research Program

#### 3.1. Stochastic Operations Research

Stochastic Operations Research is a collection of modeling, optimization and numerical computation techniques, aimed at assessing the behavior of man-made systems driven by random phenomena, and at helping to make decisions in such a context.

The discipline is based on applied probability and focuses on effective computations and algorithms. Its core theory is that of Markov chains over discrete state spaces. This family of stochastic processes has, at the same time, a very large modeling capability and the potential of efficient solutions. By “solution” is meant the calculation of some *performance metric*, usually the distribution of some random variable of interest, or its average, variance, etc. This solution is obtained either through exact “analytic” formulas, or numerically through linear algebra methods. Even when not analytically or numerically tractable, Markovian models are always amenable to “Monte-Carlo” simulations with which the metrics can be statistically measured.

An example of this is the success of classical Queueing Theory, with its numerous analytical formulas. Another important derived theory is that of the Markov Decision Processes, which allows to formalize *optimal* decision problems in a random environment. This theory allows to characterize the optimal decisions, and provides algorithms for calculating them.

Strong trends of Operations Research are: a) an increasing importance of multi-criteria multi-agent optimization, and the correlated introduction of Game Theory in the standard methodology; b) an increasing concern of (deterministic) Operations Research with randomness and risk, and the consequent introduction of topics like Chance Constrained Programming and Stochastic Optimization. Data analysis is also more and more present in Operations Research: techniques from statistics, like filtering and estimation, or Artificial Intelligence like clustering, are coupled with modeling in Machine Learning techniques like Q-Learning.

## RESIST Team

### 3. Research Program

#### 3.1. Overview

The Resist project aims at designing, implementing and validating novel models, algorithms and tools to **make networked systems elastic and resilient so as to enhance their scalability and security**, assuming users, applications and devices whose volume and heterogeneity will continue to increase.

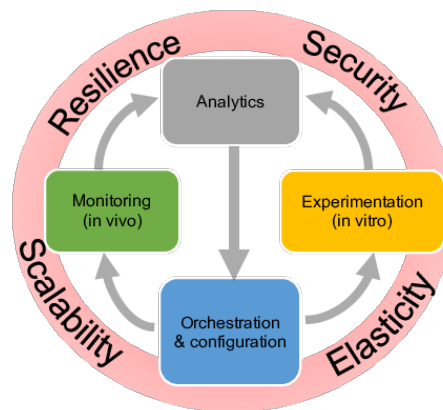


Figure 1. The Resist project

**Softwarization of networks** and **data analytics** are key enablers to design intelligent methods to orchestrate – *i.e.* configure in a synchronized and distributed manner – both network and system resources. Intelligent **orchestration** leverages relevant data for decision-making using **data analytics**. Input data reflecting the past, current and even future (predicted) states of the system are used to build relevant knowledge. Two approaches are pursued to generate knowledge and to validate orchestration decisions. First, a running system can be **monitored in vivo**. Second, **in vitro experimentation** in a controlled environment (simulators, emulators and experimental platforms) is helpful to reproduce a running system with a high reliability and under different hypotheses. Monitoring and experimentation are steered and configured through orchestration according to the two intertwined loops illustrated in Figure 1 .

Accordingly Resist is thus structured into four main research objectives (activities) namely Monitoring, Experimentation, Analytics and Orchestration.

#### 3.2. Monitoring

The evolving nature of the Internet ecosystem and its continuous growth in size and heterogeneity call for a better understanding of its characteristics, limitations, and dynamics, both locally and globally so as to improve application and protocol design, detect and correct anomalous behaviors, and guarantee performance.

To face these scalability issues, **appropriate monitoring models, methods and algorithms are required for data collection, analysis and sharing** from which knowledge about Internet traffic and usage can be extracted. Measuring and collecting traces necessitate user-centered and data-driven paradigms to cover the wide scope of heterogeneous user activities and perceptions. In this perspective, we propose monitoring algorithms and architectures for large scale environments involving mobile and Internet of Things (IoT) devices.

Resist also assesses **the impact of the Internet infrastructure evolution integrating network softwarization on monitoring**, for example the need for dedicated measurement methodologies. We take into account not only the technological specifics of such paradigms for their monitoring but also the ability to use them for collecting, storing and processing monitoring data in an accurate and cost-effective manner.

Crowd-sourcing and third-party involvement are gaining in popularity, paving the way for massively distributed and collaborative monitoring. We thus investigate opportunistic mobile crowdsensing in order to collect user activity logs along with contextual information (social, demographic, professional) to effectively measure end-users' **Quality of Experience**. However, collaborative monitoring raises serious concerns regarding trust and sensitive data sharing (open data). Data anonymization and sanitization need to be carefully addressed.

### 3.3. Experimentation

Of paramount importance in our target research context is experimental validation using testbeds, simulators and emulators. In addition to using various existing experimentation methodologies, Resist contributes in **advancing the state of the art in experimentation methods and experimental research practices**, particularly focusing on elasticity and resilience.

We develop and deploy testbeds and emulators for **experimentation with new networking paradigms** such as SDN and NFV, to enable large-scale in-vitro experiments combining all aspects of Software-Defined Infrastructures (server virtualization, SDN/NFV, storage). Such fully controlled environments are particularly suitable for our experiments on resilience, as they ease the management of fault injection features.

We are playing a central role in the development of the Grid'5000 testbed [44] and our objective is to reinforce our collaborations with other testbeds, towards a **testbed federation** in order to enable experiments to scale to multiple testbeds, providing a diverse environment reflecting the Internet itself.

Moreover, our research focuses on extending the infrastructure virtualization capabilities of our Distem [47] emulator, which provides a flexible software-based experimental environment.

Finally, methodological aspects are also important for ensuring **trustworthy and reproducible experiments**, and raises many challenges regarding testbed design, experiment description and orchestration, along with automated or assisted provenance data collection [45].

### 3.4. Analytics

A large volume of data is processed as part of the operations and management of networked systems. These include traditional monitoring data generated by network components and components' configuration data, but also data generated by dedicated network and system probes.

**Understanding and predicting security incidents or system ability to scale** requires the elaboration of novel **data analytics techniques** capable to cope with large volumes of data generated from various sources, in various formats, possibly incomplete, non-fully described or even encrypted.

We use machine learning techniques (*e.g.* Topological Data Analysis or multilayer perceptrons) and leverage our domain knowledge to fine-tune them. For instance, machine learning on network data requires the definition of new distance metrics capable to capture the properties of network configurations, packets and flows similarly to edge detection in image processing. Resist contributes to developing and making publicly available an **analytics framework dedicated to networked systems** to support Intelligence-Defined Networked Systems.

Specifically, the goal of the Resist analytics framework is to facilitate the extraction of knowledge useful for **detecting, classifying or predicting security or scalability issues**. The extracted knowledge is then leveraged for orchestration purposes to achieve system elasticity and guarantee its resilience. Indeed, predicting when, where and how issues will occur is very helpful in deciding the provisioning of resources at the right time and place. Resource provisioning can be done either reactively to solve the issues or proactively to prepare the networked system for absorbing the incident (resiliency) in a timely manner thanks to its elasticity.

While the current trend is towards centralization where the collected data is exported to the cloud for processing, we seek to extend this model by also developing and evaluating novel approaches in which **data analytics is seamlessly embedded within the monitored systems**. This combination of big data analytics with network softwarization enablers (SDN, NFV) can enhance the scalability of the monitoring and analytics infrastructure.

### 3.5. Orchestration

The ongoing transformations in the Internet ecosystem including network softwarization and cloudification bring new management challenges in terms of service and resource orchestration. Indeed, the growing sophistication of Internet applications and the complexity of services deployed to support them require novel models, architectures and algorithms for their automated **configuration** and **provisioning**. Network applications are more and more instantiated through the **composition of services, including virtualized hardware and software resources**, that are offered by **multiple providers** and are subject to changes and updates over time. In this dynamic context, efficient orchestration becomes fundamental for ensuring performance, resilience and security of such applications. We are investigating the chaining of different functions for supporting the security protection of smart devices, based on the networking behavior of their applications.

From a resilience viewpoint, this orchestration at the network level allows the dynamic **reconfiguration of resources** to absorb the effects of congestions, such as link-flooding behaviors. The goal is to drastically reduce the effects of these congestions by imposing dynamic policies on all traffic where the network will adapt itself until it reaches a stable state. We also explore mechanisms for **detecting and remediating potential dysfunctions** within a virtualized network. Corrective operations can be performed through dynamically composed VNFs (Virtualized Network Functions) based on available resources, their dependencies (horizontal and vertical), and target service constraints. We also conduct research on verification methods for automatically assessing and validating the composed chains.

From a security viewpoint, this orchestration provides **prevention mechanisms** that capture adversaries' intentions early and **enforces security policies** in advance through the available resources, to be able to proactively mitigate their attacks. We mainly rely on the results obtained in our research activity on security analytics to build such policies, and the orchestration part focuses on the required algorithms and methods for their automation.



## SOCRATE Project-Team

### 3. Research Program

#### 3.1. Flexible Radio Front-End

These are the research axis as they were proposed at the creation of the Socrate Team.

This axis mainly deals with the radio front-end of software radio terminals. In order to ensure a high flexibility in a global wireless network, each node is expected to offer as many degrees of freedom as possible. For instance, the choice of the most appropriate communication resource (frequency channel, spreading code, time slot,...), the interface standard or the type of antenna are possible degrees of freedom. The *multi-\** paradigm denotes a highly flexible terminal composed of several antennas providing MIMO features to enhance the radio link quality, which is able to deal with several radio standards to offer interoperability and efficient relaying, and can provide multi-channel capability to optimize spectral reuse. On the other hand, increasing degrees of freedom can also increase the global energy consumption, therefore for energy-limited terminals a different approach has to be defined.

In this research axis, we expect to demonstrate optimization of flexible radio front-end by fine grain simulations, and also by the design of home made prototypes. Of course, studying all the components deeply would not be possible given the size of the team, we are currently not working in new technologies for DAC/ADC and power amplifiers which are currently studied by hardware oriented teams. The purpose of this axis is to build system level simulation taking into account the state of the art of each key component.

#### 3.2. Multi-User Communications

While the first and the third research axes deal with the optimization of the cognitive radio nodes themselves from system and programming point of view, an important complementary objective is to consider the radio nodes in their environments. Indeed, cognitive radio does not target the simple optimization of point to point transmissions, but the optimization of simultaneous concurrent transmissions. The tremendous development of new wireless applications and standards currently observed calls for a better management of the radio spectrum with opportunistic radio access, cooperative transmissions and interference management. This challenge has been identified as one of the most important issue for 5G to guarantee a better exploitation of the spectrum. In addition, mobile internet is going to support a new revolution that is the *tactile internet*, with real time interactions between the virtual and the real worlds, requiring new communication objectives to be met such as low latency end to end communications, distributed learning techniques, in-the-network computation, and many more. The future network will be heterogeneous in terms of technologies, type of data flows and QoS requirements. To address this revolution two work directions have naturally formed within the axis. The first direction concerns the theoretical study of fundamental limits in wireless networks. Introduced by Claude Shannon in the 50s and heavily developed up to today, Information Theory has provided a theoretical foundation to study the performance of wireless communications, not from a practical design view point, but using the statistical properties of wireless channels to establish the fundamental trade-offs in wireless communications. Beyond the classical *energy efficiency - spectral efficiency* tradeoff, information theory and its many derivations, i.e., network information theory, may also help to address additional questions such as determining the optimal rates under decentralized policies, asymptotic behavior when the density of nodes increases, latency controlled communication with finite block-length theory, etc. In these cases, information theory is often associated to other theoretical tools such as game theory, stochastic geometry, control theory, graph theory and many others.

Our first research direction consists in evaluating specific multi-user scenarios from a network information theory perspective, inspired by practical scenarios from various applicative frameworks (e.g. 5G, Wifi, sensor networks, IoT, etc.), and to establish fundamental limits for these scenarios. The second research direction is related to algorithmic and protocol design (PHY/MAC), applied to practical scenarios. Exploiting signal processing, linear algebra inspired models and distributed algorithms, we develop and evaluate various distributed algorithms allowing to improve many QoS metrics such as communication rates, reliability, stability, energy efficiency or computational complexity.

It is clear that both research directions are symbiotic with respect to each other, with the former providing theoretical bounds that serves as a reference to the performance of the algorithms created in the later. In the other way around, the later offers target scenarios for the former, through identifying fundamental problems that are interesting to be studied from the fundamental side. Our contributions of the year in these two directions are summarized further in the document.

### 3.3. Software Radio Programming Model

Finally the third research axis is concerned with software aspect of the software radio terminal. We have currently two actions in this axis, the first one concerns the programming issues in software defined radio devices, the second one focusses on low power devices: how can they be adapted to integrate some reconfigurability.

The expected contributions of Socrate in this research axis are :

- The design and implementation of a “middleware for SDR”, probably based on a Virtual Machine.
- Prototype implementations of novel software radio systems, using chips from Leti and/or Lyrtech software radio boards.
- Development of a *smart node*: a low-power Software-Defined Radio node adapted to WSN applications.
- Methodology clues and programming tools to program all these prototypes.

### 3.4. Evolution of the Socrate team

In 2018 the Socrate team which was originally conceived to develop software defined radio has decided to split in two teams: the Maracas team will consist of the activities of Socrate Axis 2 and be directed by Jean-Marie Gorce, and the Socrate team which will consist in the Axis 1 and 3 of the current version of Socrate. This change is explicit since september 2018 as the Maracas team is created.

The advent of non-volatile memory technologies (NVRAM) is causing a major evolution in all software layers. On the one hand, the non-volatility of data in the event of a breakdown necessarily leads to fatal inconsistencies if the memory is not managed correctly. On the other hand, these memories have very different performances from the usual DRAM, which tends to the appearance of hybrid and complex memory hierarchies. Many technological and scientific challenges are to be faced in all software layers to deal with these two sets of issues. Above all, the answers to be provided depend on the calculation system considered and for what purpose it is constructed.

Within the framework of very low consumption sensors and devices, the Socrate team proposed, with Sytare [28], a software solution allowing to develop embedded applications on platforms supporting an intermittent power supply (TPS – *Transiently Powered System*) and integrating NVRAM as illustrated in Figure 3 . The IPL ZEP (<https://project.inria.fr/iplzep/>) was also launched by Socrate last year to respond to various scientific challenges related to this issue.

The recent advance in harvesting technologies provides new research direction to Socrate which have skills in radio propagation and low power radio (wake-up radio for instance [31]). Fig 3 , illustrates the *future ultra-low sensor* as envisioned by Socrate.

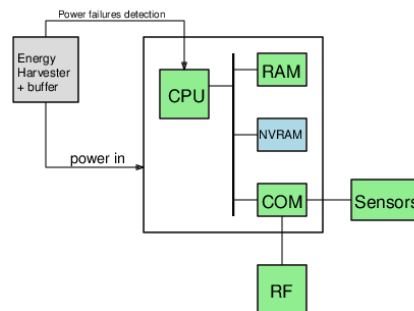


Figure 3. Architecture targeted by Socrate: low energy wireless sensor with peripherals and non volatile memory

## TRIBE Project-Team

### 3. Research Program

#### 3.1. Research program

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

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

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