

RESEARCH CENTRE

**Inria Lyon Centre**

IN PARTNERSHIP WITH:

Université Claude Bernard (Lyon 1), Ecole  
normale supérieure de Lyon, CNRS

2024

ACTIVITY REPORT

Project-Team

AVALON

**Algorithms and Software Architectures for  
Distributed and HPC Platforms**

IN COLLABORATION WITH: Laboratoire de l'Informatique du Parallélisme  
(LIP)

**DOMAIN**

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

**THEME**

**Distributed and High Performance  
Computing**

*Inria*

# Contents

<b>Project-Team AVALON</b>	<b>1</b>
<b>1 Team members, visitors, external collaborators</b>	<b>2</b>
<b>2 Overall objectives</b>	<b>3</b>
2.1 Presentation	3
2.2 Objectives	4
<b>3 Research program</b>	<b>4</b>
3.1 Energy Application Profiling and Modeling	4
3.2 Data-intensive Application Profiling, Modeling, and Management	5
3.3 Resource-Agnostic Application Description Model	5
3.4 Application Mapping and Scheduling	6
3.4.1 Application Mapping and Software Deployment	6
3.4.2 Non-Deterministic Workflow Scheduling	6
<b>4 Application domains</b>	<b>6</b>
4.1 Overview	6
4.2 Climatology	7
4.3 Astrophysics	7
4.4 Bioinformatics	7
<b>5 Social and environmental responsibility</b>	<b>7</b>
5.1 Footprint of research activities	7
<b>6 Highlights of the year</b>	<b>8</b>
<b>7 New software, platforms, open data</b>	<b>8</b>
7.1 New software	8
7.1.1 Halley	8
7.1.2 XKBLAS	8
7.1.3 execo	9
7.2 New platforms	9
7.2.1 Platform: Grid'5000	9
7.2.2 Platform: SLICES-FR	10
7.2.3 Platform: SLICES	10
<b>8 New results</b>	<b>10</b>
8.1 Energy Efficiency in Large Scale Distributed Systems	10
8.1.1 Estimating the power consumption of bare metal water-cooled servers	10
8.1.2 Deep Reinforcement Learning for Energy-efficient Selection of Embedded Services at the Edge	11
8.1.3 S-ORCA : A social-based consolidation approach to reduce Cloud infrastructures energy consumption	11
8.1.4 Estimating the environmental impact of Generative-AI services	11
8.1.5 A Methodology and a Toolbox to Explore Dataset related to the Environmental Impact of HTTP Requests	12
8.1.6 Exploring RAPL as a Power Capping Leverage for Power-Constrained Infrastructures	12
8.1.7 Many locks, few leverages! Digital sufficiency: the situation is serious but not hopeless	13
8.1.8 Sufficient generative AI: an oxymoron?	13
8.1.9 Fine-grained methodology to assess environmental impact of a set of digital services	13
8.2 Edge, Cloud and Distributed Resource Management	14
8.2.1 SkyData: Autonomous Data paradigm	14
8.2.2 Numerics in the Cloud	14

8.2.3	A specialized model and implementation of an actuarial chatbot based on Federated Learning	15
8.3	HPC Applications and Runtimes	15
8.3.1	Measuring and Interpreting Dependent Task-based Applications Performances.	15
8.3.2	Evaluation of mix MPI + Dependent OpenMP task programming models on Fugaku.	15
8.3.3	Handling dynamicity of HPC applications designed by a task-based component model	16
8.3.4	New high level programming framework: Experiments with Kokkos	16
8.3.5	Taskgrind: Heavyweight Dynamic Binary Instrumentation for Parallel Programs Analysis	16
8.3.6	ETH4HPC Strategic Research Agenda	17
<b>9</b>	<b>Bilateral contracts and grants with industry</b>	<b>17</b>
9.1	Bilateral grants with industry	17
<b>10</b>	<b>Partnerships and cooperations</b>	<b>17</b>
10.1	International initiatives	17
10.1.1	Participation in other International Programs	17
10.2	International research visitors	18
10.2.1	Visits of international scientists	18
10.3	European initiatives	18
10.3.1	Horizon Europe	18
10.3.2	Other european programs/initiatives	19
10.4	National initiatives	19
<b>11</b>	<b>Dissemination</b>	<b>23</b>
11.1	Promoting scientific activities	23
11.1.1	Scientific events: organisation	23
11.1.2	Scientific events: selection	23
11.1.3	Journal	24
11.1.4	Invited talks	24
11.1.5	Scientific expertise	24
11.1.6	Research administration	25
11.2	Teaching - Supervision - Juries	25
11.2.1	Teaching	25
11.2.2	Supervision	26
11.2.3	Juries	27
11.3	Popularization	27
11.3.1	Productions (articles, videos, podcasts, serious games, ...)	27
11.3.2	Participation in Live events	27
<b>12</b>	<b>Scientific production</b>	<b>27</b>
12.1	Publications of the year	27
12.2	Cited publications	30

## Project-Team AVALON

*Creation of the Project-Team: 2014 July 01*

### Keywords

#### Computer sciences and digital sciences

- A1.1.1. – Multicore, Manycore
- A1.1.2. – Hardware accelerators (GPGPU, FPGA, etc.)
- A1.1.4. – High performance computing
- A1.1.5. – Exascale
- A1.3.5. – Cloud
- A1.3.6. – Fog, Edge
- A1.6. – Green Computing
- A2.1.6. – Concurrent programming
- A2.1.7. – Distributed programming
- A2.1.10. – Domain-specific languages
- A2.2.8. – Code generation
- A2.5.2. – Component-based Design
- A2.6.2. – Middleware
- A2.6.4. – Ressource management
- A3.1.3. – Distributed data
- A4.4. – Security of equipment and software
- A7.1. – Algorithms
- A7.1.1. – Distributed algorithms
- A7.1.2. – Parallel algorithms
- A8.2.1. – Operations research
- A8.9. – Performance evaluation

#### Other research topics and application domains

- B1.1.7. – Bioinformatics
- B4.5. – Energy consumption
- B4.5.1. – Green computing
- B6.1.1. – Software engineering
- B9.5.1. – Computer science
- B9.7. – Knowledge dissemination
- B9.7.1. – Open access
- B9.7.2. – Open data
- B9.8. – Reproducibility

# 1 Team members, visitors, external collaborators

## Research Scientists

- Christian Perez [Team leader, INRIA, Senior Researcher, HDR]
- Simon Delamare [CNRS, Researcher]
- Thierry Gautier [INRIA, Researcher, from Oct 2024, HDR]
- Thierry Gautier [INRIA, Researcher, until Aug 2024, HDR]
- Laurent Lefevre [INRIA, Researcher, HDR]

## Faculty Members

- Yves Caniou [UNIV LYON I, Associate Professor]
- Eddy Caron [UNIV LYON I, Professor, HDR]
- Olivier Glück [UNIV LYON I, Associate Professor]
- Elise Jeanneau [UNIV LYON I, Associate Professor]

## Post-Doctoral Fellows

- Jerry Lacmou Zeutouo [INRIA, Post-Doctoral Fellow, until Aug 2024]
- Lucien Ndjie Ngale [ENS DE LYON, Post-Doctoral Fellow, until Nov 2024]
- Romain Pereira [INRIA, Post-Doctoral Fellow, from Feb 2024]

## PhD Students

- Maxime Agusti [OVH, CIFRE]
- Adrien Berthelot [OCTO TECHNOLOGY, CIFRE, Phd defense 12/11/2024]
- Emile Egreteau–Druet [INRIA, from Dec 2024]
- Simon Lambert [CIRIL GROUP, CIFRE]
- Vladimir Ostapenco [INRIA]
- Thomas Stavis [INRIA, from Oct 2024]
- Gabriel Suau [CEA, CIFRE]

## Technical Staff

- Annour Saad Allamine [INRIA, Engineer, from Dec 2024]
- Johanna Desprez [INRIA, Engineer, from Apr 2024]
- Pierre Jacquot [INRIA, Engineer]
- Jean Christophe Mignot [CNRS, Engineer]
- Emeline Pegon [CNRS, Engineer, from Aug 2024]
- Pierre-Etienne Polet [INRIA, Engineer, from Mar 2024]
- Pierre-Etienne Polet [INRIA, Engineer, until Feb 2024]

- Dominique Ponsard [CNRS, Engineer]
- Jean-Camille Seck [INRIA, Engineer, from Oct 2024]
- Anass Serhani [INRIA, Engineer]

### **Interns and Apprentices**

- Hamza Aabirrouche [UNIV LYON I, Intern, from Aug 2024]
- Annour Saad Allamine [UNIV LYON I, Intern, from Mar 2024 until Oct 2024]
- Yuliy Daniel [INRIA, Intern, from May 2024 until Jul 2024]
- Cyril Devaux [INRIA, Apprentice, from Oct 2024]
- Cyril Devaux [INRIA, Intern, from May 2024 until Jul 2024]
- Maxime Just [ENS DE LYON, Intern, from Oct 2024]
- Maxime Just [ENS DE LYON, Intern, from Feb 2024 until Jul 2024]
- Thomas Stavis [INRIA, Intern, from Feb 2024 until Jul 2024]

### **Administrative Assistant**

- Chrystelle Mouton [INRIA]

### **External Collaborator**

- Doreid Ammar [AIVANCITY]

## **2 Overall objectives**

### **2.1 Presentation**

The fast evolution of hardware capabilities in terms of wide area communication, computation and machine virtualization leads to the requirement of another step in the abstraction of resources with respect to parallel and distributed applications. These large scale platforms based on the aggregation of large clusters (Grids), datacenters (Clouds) with IoT (Edge/Fog), or high performance machines (Supercomputers) are now available to researchers of different fields of science as well as to private companies. This variety of platforms and the way they are accessed also have an important impact on how applications are designed (*i.e.*, the programming model used) as well as how applications are executed (*i.e.*, the runtime/middleware system used). The access to these platforms is driven through the use of multiple services providing mandatory features such as security, resource discovery, load-balancing, monitoring, *etc.*

The goal of the AVALON team is to execute parallel and/or distributed applications on parallel and/or distributed resources while ensuring user and system objectives with respect to performance, cost, energy, security, *etc.* Users are generally not interested in the resources used during the execution. Instead, they are interested in how their application is going to be executed: the duration, its cost, the environmental footprint involved, *etc.* This vision of utility computing has been strengthened by the cloud concepts and by the short lifespan of supercomputers (around three years) compared to application lifespan (tens of years). Therefore a major issue is to design models, systems, and algorithms to execute applications on resources while ensuring user constraints (price, performance, *etc.* ) as well as system administrator constraints (maximizing resource usage, minimizing energy consumption, *etc.* ).

## 2.2 Objectives

To achieve the vision proposed in the previous section, the AVALON project aims at making progress on four complementary research axes: energy, data, programming models and runtimes, application scheduling.

**Energy Application Profiling and Modeling** AVALON will improve the profiling and modeling of scientific applications with respect to energy consumption. In particular, it will require to improve the tools that measure the energy consumption of applications, virtualized or not, at large scale, so as to build energy consumption models of applications.

**Data-intensive Application Profiling, Modeling, and Management** AVALON will improve the profiling, modeling, and management of scientific applications with respect to CPU and data intensive applications. Challenges are to improve the performance prediction of parallel regular applications, to model and simulate (complex) intermediate storage components, and data-intensive applications, and last to deal with data management for hybrid computing infrastructures.

**Programming Models and Runtimes** AVALON will design component-based models to capture the different facets of parallel and distributed applications while being resource agnostic, so that they can be optimized for a particular execution. In particular, the proposed component models will integrate energy and data modeling results. AVALON in particular targets OpenMP runtime as a specific use case and contributes to improve it for multi-GPU nodes.

**Application Mapping and Scheduling** AVALON will propose multi-criteria mapping and scheduling algorithms to meet the challenge of automating the efficient utilization of resources taking into consideration criteria such as performance (CPU, network, and storage), energy consumption, and security. AVALON will in particular focus on application deployment, workflow applications, and security management in clouds.

All our theoretical results will be validated with software prototypes using applications from different fields of science such as bioinformatics, physics, cosmology, *etc.* The experimental testbeds GRID'5000 and SLICES will be our platforms of choice for experiments.

## 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. 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

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 of the 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.* ) can be seen as a first piece of a 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 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.

## 4 Application domains

### 4.1 Overview

The AVALON team targets applications with large computing and/or data storage needs, which are still difficult to program, deploy, and maintain. Those applications can be parallel and/or distributed applications, such as large scale simulation applications or code coupling applications. Applications can also be workflow-based as commonly found in distributed systems such as grids or clouds.

The team aims at not being restricted to a particular application field, thus avoiding any spotlight. The team targets different HPC and distributed application fields, which brings use cases with different issues. This will be eased with our participation to the Joint Laboratory for Extreme Scale Computing (JLESC), to BioSyL, a federative research structure about Systems Biology of the University of Lyon, or to the SKA project. Last but not least, the team has a privileged connection with CC-IN2P3 that opens up collaborations, in particular in the astrophysics field.

In the following, some examples of representative applications that we are targeting are presented. In addition to highlighting some application needs, they also constitute some of the use cases that will be used to validate our theoretical results.

## 4.2 Climatology

The world's climate is currently changing due to the increase of the greenhouse gases in the atmosphere. Climate fluctuations are forecasted for the years to come. For a proper study of the incoming changes, numerical simulations are needed, using general circulation models of a climate system. Simulations can be of different types: HPC applications (e.g., the NEMO framework [28] for ocean modelization), code-coupling applications (e.g., the OASIS coupler [33] for global climate modeling), or workflows (long term global climate modeling).

As for most applications the team is targeting, the challenge is to thoroughly analyze climate-forecasting applications to model their needs in terms of programming model, execution model, energy consumption, data access pattern, and computing needs. Once a proper model of an application has been set up, appropriate scheduling heuristics can be designed, tested, and compared. The team has a long tradition of working with CERFACS on this topic, since for example in the LEGO (2006-09) and SPADES (2009-12) French ANR projects.

## 4.3 Astrophysics

Astrophysics is a major field to produce large volumes of data. For instance, the **Vera C. Rubin Observatory** will produce 20 TB of data every night, with the goals of discovering thousands of exoplanets and of uncovering the nature of dark matter and dark energy in the universe. The **Square Kilometer Array** will produce 9 Tbits/s of raw data. One of the scientific projects related to this instrument called Evolutionary Map of the Universe is working on more than 100 TB of images. The **Euclid Imaging Consortium** will generate 1 PB data per year.

The SKA project () is an international effort to build and operate the world's largest radiotelescopes covering all together the wide frequency range between 50 MHz and 15.4 GHz. The scale of the SKA project represents a huge leap forward in both engineering and research & development towards building and delivering a unique Observatory, whose construction has officially started on July 2021. The SKA Observatory is the second intergovernmental organisation for ground-based astronomy in the world, after the European Southern Observatory. AVALON participates to the activities of the SCOOP team in SKAO's SAFE framework that deals with platforms related issues such as application benchmarking and profiling, hardware-software co-design.

## 4.4 Bioinformatics

Large-scale data management is certainly one of the most important applications of distributed systems in the future. Bioinformatics is a field producing such kinds of applications. For example, DNA sequencing applications make use of MapReduce skeletons.

The AVALON team is a member of **BioSyL**, a Federative Research Structure attached to University of Lyon. It gathers about 50 local research teams working on systems biology. AVALON is in particular collaborating with the Inria **Beagle** team on artificial evolution and computational biology as the challenges are around high performance computation and data management.

# 5 Social and environmental responsibility

## 5.1 Footprint of research activities

Through its research activities on energy efficiency and on energy and environmental impacts reductions, Avalon tries to reduce some impacts of distributed systems.

Avalon deals with frugality in clouds with the leadership of the FrugalCloud challenge (*Défi*) between Inria and OVHcloud. Laurent Lefevre is also involved in the steering committee of the **EcoInfo** GDS CRNS group which deals with eco-responsibility of ICT. Avalon is also involved in the sustainable management of large scale experimental infrastructures like Slices. Laurent Lefevre has proposed a Green Slices methodology which is under review in order to deal with the life cycle of such infrastructures. Laurent Lefevre and Emeline Pegon are strongly involved in the Alt-Impact programme on digital sufficiency, between Ademe, CNRS and Inria.

## 6 Highlights of the year

The Interim Supervisory Board of SLICES RI validated the Statutes (without the appendices) as well the Scientific and Technical description. This is an important step towards the creation of SLICES-ERIC.

## 7 New software, platforms, open data

### 7.1 New software

#### 7.1.1 Halley

**Name:** Halley

**Keywords:** Software Components, HPC

**Scientific Description:** Halley is an implementation of the COMET component model that enables to efficiently compose independent parallel codes using both classical use/provide ports but also dataflow oriented ports that are used to generate tasks for multi-core shared-memory machines.

**Functional Description:** Halley transforms a COMET assembly into a L2C assembly that contains some special components that deal with the data flow section. In particular, a dataflow section of COMET generates a "scheduler" L2C component that contains the code that is in charged of creating its tasks.

**Release Contributions:** In 2024, dynamic workflow support was added. So, for example, the input data size of a meta-task can now depend on the output of a previous meta-task.

**Publications:** [tel-01663718](#), [hal-01518730](#), [hal-01566288](#), [hal-01901806](#)

**Contact:** Christian Perez

**Participants:** Jérôme Richard, Christian Perez, Jerry Lacmou Zeutouo

#### 7.1.2 XKBLAS

**Name:** XKBLAS

**Keywords:** BLAS, Dense linear algebra, GPU

**Functional Description:** XKBLAS is yet another BLAS library (Basic Linear Algebra Subroutines) that targets multi-GPUs architecture thanks to the XKaapi runtime and with block algorithms from PLASMA library. XKBLAS is able to exploit large multi-GPUs node with sustained high level of performance. The library offers a wrapper library able to capture calls to BLAS (C or Fortran). The internal API is based on asynchronous invocations in order to enable overlapping between communication by computation and also to better composed sequences of calls to BLAS.

This current version of XKBlas is the first public version and contains only BLAS level 3 algorithms, including XGEMMT:

XGEMM XGEMMT: see MKL GEMMT interface XTRSM XTRMM XSYMM XSYRK XSYR2K XHEMM XHERK XHER2K

For classical precision Z, C, D, S.

**Release Contributions:** 0.1.x versions: calls to BLAS kernels must be initiate by the same thread that initializes the XKBlas library. 0.2.x versions: better support for libblas\_wrapper and improved scheduling heuristic to take into account memory hierarchy between GPUs 0.4.x versions: add support for AMD GPU 0.5.x : better support for AMD GPU (MI250x). Add capacity to clustering GPUs and CPU threads.

**News of the Year:** New development to clustering both GPUs and CPU threads invoking BLAS kernels in order to reduce communication between resources. Better support for AMD GPU (configuration and performance). Presentation during the MUMPS annual meeting at ANSYS, 2024/06/20. Seminar presentation at EDF, 2024/11/26.

**URL:** <https://gitlab.inria.fr/xkblas/versions>

**Contact:** Thierry Gautier

**Participants:** Thierry Gautier, João Vicente Ferreira Lima

### 7.1.3 execo

**Keywords:** Toolbox, Deployment, Orchestration, Python

**Functional Description:** Execo offers a Python API for asynchronous control of local or remote, standalone or parallel, unix processes. It is especially well suited for quickly and easily scripting workflows of parallel/distributed operations on local or remote hosts: automate a scientific workflow, conduct computer science experiments, perform automated tests, etc. The core python package is execo. The execo\_g5k package provides a set of tools and extensions for the Grid5000 testbed. The execo\_engine package provides tools to ease the development of computer sciences experiments.

**Release Contributions:** Release 2.8.1 on October 21, 2024 (list of changes since version 2.7: adapt to changes in oarstat output format (compatibility with old and new output formats), g5k api cache stored as json instead of pickle, support clusters names ending with numbers (eg. abacus-X), canonical\_host\_name handles interface / kavlan / ipv6 and support cluster names ending with numbers + fix for ifname != ethX (eg. fpgaX), add get\_host\_interface, extend planning API to allow requests at node level additionally to cluster and site level, spawn process lifecycle handlers in separate threads to avoid blocking + refactoring, handle encoding (py3+) when writing to Process, full redesign of the Processes expect implementation, add get\_cluster\_queues and get\_cluster\_jobtypes, add KaconsoleProcess, add substitutions to filenames in stdout/stderr handlers, scp commands in Get / Put as list and shell=False to (securely) handle spaces in path, fix eating 100% of one core iterating through high number of fd to close them in conductor, fix various regexes within invalid escape sequences, add option in execo\_engine for using pty to copy\_outputs(), fix corner case in process args handling in Remote)

**URL:** <https://gitlab.inria.fr/mimbert/execo>

**Contact:** Matthieu Imbert

**Participants:** Florent Chuffart, Laurent Pouilloux, Matthieu Imbert

## 7.2 New platforms

### 7.2.1 Platform: Grid'5000

**Participants:** Simon Delamare, Pierre Jacquot, Matthieu Imbert, Laurent Lefèvre, Christian Perez, Jean-Camille Seck, Quentin Assuncao, Cyril Devaux.

#### FUNCTIONAL DESCRIPTION

The Grid'5000 experimental platform is a scientific instrument to support computer science research related to distributed systems, including parallel processing, high performance computing, cloud computing, operating systems, peer-to-peer systems and networks. It is distributed on 10 sites in France and Luxembourg, including Lyon. Grid'5000 is a unique platform as it offers to researchers many and varied hardware resources and a complete software stack to conduct complex experiments, ensure reproducibility and ease understanding of results.

- Contact: Laurent Lefèvre
- URL: [www.grid5000.fr/](http://www.grid5000.fr/)

### 7.2.2 Platform: SLICES-FR

**Participants:** Simon Delamare, Pierre Jacquot, Matthieu Imbert, Laurent Lefèvre, Christian Perez, Jean-Camille Seck, Quentin Assuncao, Cyril Devaux.

**FUNCTIONAL DESCRIPTION** The SLICES-FR infrastructure aims at providing an experimental platform for experimental computer Science (Internet of things, clouds, HPC, big data, *etc.* ). This new infrastructure will supersede two existing infrastructures, Grid'5000 and FIT.

- Contact: Christian Perez
- URL: [www.slices-fr.eu/](http://www.slices-fr.eu/)

### 7.2.3 Platform: SLICES

**Participants:** Simon Delamare, Pierre Jacquot, Laurent Lefèvre, Christian Perez.

**FUNCTIONAL DESCRIPTION** SLICES is an European effort that aims at providing a flexible platform designed to support large-scale, experimental research focused on networking protocols, radio technologies, services, data collection, parallel and distributed computing and in particular cloud and edge-based computing architectures and services. SLICES-FR is the The French node of SLICES.

- Contact: Christian Perez
- URL: [www.slices-ri.eu](http://www.slices-ri.eu/)

## 8 New results

### 8.1 Energy Efficiency in Large Scale Distributed Systems

#### 8.1.1 Estimating the power consumption of bare metal water-cooled servers

**Participants:** Maxime Agusti, Eddy Caron, Laurent Lefèvre.

Numerous cloud providers offer physical servers for rental in bare metal paradigm. This mode gives customers total control over hardware resources, but limits cloud providers' visibility of their usage. Accurately measuring server energy consumption in this context represents a major challenge, as installing physical energy meters is both costly and complex. Existing energy models are generally based on system usage data, which is incompatible with the general privacy policies of bare-metal server contracts. To deal with these problems, it is imperative to develop new approaches for estimating the energy consumption of these servers. This paper presents an original non-intrusive method for estimating the energy consumption of a server cooled by direct-chip liquid-cooling, based on the coolant temperature and the processor temperature obtained via IPMI. Our approach is evaluated on an experiment carried out on 19 bare metal servers of a production infrastructure equipped with physical wattmeters.[3]

### 8.1.2 Deep Reinforcement Learning for Energy-efficient Selection of Embedded Services at the Edge

**Participants:** Laurent Lefèvre, Doreid Ammar, Hugo Hadjur.

Edge computing helps to release the tension at the center of IoT systems' networks, thus reducing the latency, optimizing the bandwidth, and providing new privacy and security solutions, among others. Despite its benefits, edge computing faces unique challenges, including latency, security, and resource constraints. Among these challenges, energy consumption has emerged in the research community, and the global objective is to "do more with less". Researchers explore diverse strategies to enhance sustainability, from hardware optimizations to intelligent algorithms. The quest for energy efficiency and to reduce several impacts aligns with broader efforts to create an environmentally conscious technology landscape. In this paper, we present a task selection model for the edge. We focus on energy consumption and aim to maximize the value given by tasks, all the while minimizing the energy consumed. To do so, we develop a computer environment to simulate outdoor energy-harvesting edge devices, contribute to research reproducibility by recreating a photovoltaic energy harvesting prediction model, and train deep reinforcement learning models to select the best set of tasks at the edge. Our best deep reinforcement learning model, which uses Trust Region Policy Optimization, outperforms our best heuristic and is a robust task selector under varying external conditions.

Some experimental parts of this work occur in the CPER LECO/GreenCube project and some parts are financially supported by aivancity School for Technology, Business & Society Paris-Cachan. This work, within the the Ph.D. of Hugo Hadjur is co-advised by Doreid Ammar (Academic Dean and Professor at aivancity School for Technology, Business & Society Paris-Cachan and external member of Avalon team) and Laurent Lefevre [7].

### 8.1.3 S-ORCA : A social-based consolidation approach to reduce Cloud infrastructures energy consumption

**Participants:** Eddy Caron, Laurent Lefèvre, Simon Lambert.

Information and Communication Technologies (ICT) and Data Centres (DC) have nowadays considerable environmental impacts. The number of applications and services hosted in the cloud is significant, and the associated number of infrastructure is steadily increasing. Cloud Service Providers (CSP) need to respond to this growing demand and size their infrastructures accordingly but users misbehaviour and expectations in terms of service quality lead to oversized infrastructures. Infrastructures are sized to meet peak-demand, resulting in poor resource utilization and additional power consumption. But as their behaviour can increase DC energy consumption and environmental footprint, users can also help reducing them. In this paper, we study how users provided with a simple tool can reduce the power consumption of a virtualization cluster in a cloud company. Using a Virtual Machine (VM) shutdown policy, users can directly contribute to the mitigation of the power consumption of the infrastructure. Part of the paper is dedicated to profile the users and understand their behaviour when it comes to powering off their VMs. To further reduce the energy consumption of the cluster, we combine the VM shutdown policy with a simple consolidation heuristic. Simulations show a 23.95% power consumption reduction, with an additional 8.72% reduction thanks to the users. A production implementation was conducted and results in a 12.58% power consumption reduction over one week [8, 23].

### 8.1.4 Estimating the environmental impact of Generative-AI services

**Participants:** Adrien Berthelot, Eddy Caron, Mathilde Jay, Laurent Lefèvre.

Generative AI (Gen-AI) represents a major growth potential for the digital industry, a new stage in digital transformation through its many applications. Unfortunately, by accelerating the growth of digital technology, Gen-AI is contributing to the significant and multiple environmental damage caused by its sector. The question of the sustainability of IT must include this new technology and its applications, by measuring its environmental impact. To best respond to this challenge, we propose various ways of improving the measurement of Gen-AI's environmental impact. Whether using life-cycle analysis methods or direct measurement experiments, we illustrate our methods by studying Stable Diffusion a Gen-AI image generation available as a service. By calculating the full environmental costs of this Gen-AI service from end to end, we broaden our view of the impact of these technologies. We show that Gen-AI, as a service, generates an impact through the use of numerous user terminals and networks. We also show that decarbonizing the sources of electricity for these services will not be enough to solve the problem of their sustainability, due to their consumption of energy and rare metals. This consumption will inevitably raise the question of feasibility in a world of finite resources. We therefore propose our methodology as a means of measuring the impact of Gen-AI in advance. Such estimates will provide valuable data for discussing the sustainability or otherwise of Gen-AI solutions in a more transparent and comprehensive way [4, 21] [1]. This result is a joint work explored during the PhD of Adrien Berthelot [16] co-advised by Laurent Lefevre and Eddy Caron and during the PhD of Mathilde Jay [17] co-advised by Laurent Lefevre and Denis Trystram (UGA).

#### 8.1.5 A Methodology and a Toolbox to Explore Dataset related to the Environmental Impact of HTTP Requests

**Participants:** Mathilde Jay, Laurent Lefèvre.

EcoIndex has been proposed to evaluate the absolute environmental performance of a given URL using a score ranging from 0 to 100 (the higher, the better). In this article, we make a critical analysis of the initial approach and propose alternatives that no longer calculate a plain score but allow the query to be situated among other queries. The generalized critiques come with statistics and rely on extensive experiments (first contribution). Then, we move on to low-cost Machine Learning (ML) approaches (second contribution) and a transition before obtaining our final results (third contribution). Our research aims to extend the initial idea of analytical computation, i.e., a relation between three variables, in the direction of algorithmic ML computations. The fourth contribution corresponds to a discussion on our implementation, available on a GitHub repository. Along with the paper, we invite the reader to examine the question: What attributes make sense for our problem?, or equivalently, what is a relevant data policy for studying digital environmental impacts? Beyond computational questions, it is important for the scientific community to focus on this question in particular. We currently promote using wellestablished ML techniques because of their potential, which we discuss in this research. However, we also question techniques for their frugality or otherwise. We also want to encourage synergy between technical expertise and business knowledge because this is fundamental for advancing the data project.[6]

#### 8.1.6 Exploring RAPL as a Power Capping Leverage for Power-Constrained Infrastructures

**Participants:** Laurent Lefèvre, Vladimir Ostapenco.

Data centers are very energy-intensive facilities whose power provision is challenging and constrained by power bounds. In modern data centers, servers account for a significant portion of the total power consumption. In this context, the ability to limit the instant power consumption of an individual computing node is an important requirement. There are several energy and power capping techniques that can be used to limit compute node power consumption, such as Intel RAPL. Although it is nowadays mainly utilized for energy measurement, Intel RAPL (Running Average Power Limit) was originally designed for power limitation purposes. Some works use Intel RAPL for power limitation in a limited context without



full knowledge of the inner workings of this technology and what is done behind the scenes to enforce the power constraint. Furthermore, Intel has not revealed any details about its internal implementation. It is unclear exactly how Intel RAPL technology operates and what effects it has on application performance and power consumption. In this work, we conduct a thorough analysis of Intel RAPL technology as a power capping leverage on a variety of heterogeneous nodes for a selection of CPU and memory intensive workloads. For this purpose, we first validate Intel RAPL power capping mechanism using a high-precision external power meter and investigate properties such as accuracy, power limit granularity, and settling time. Then, we attempt to determine which mechanisms are employed by RAPL to adjust power consumption [9]. This work was explored during the PhD of Vladimir Ostapenco [18] co-advised by Laurent Lefevre and Anne-Cécile Orgerie (Magellan team, IRISA).

#### 8.1.7 Many locks, few levers! Digital sufficiency: the situation is serious but not hopeless

**Participants:** Laurent Lefèvre.

We are attached to digital technology as a work tool, a research subject, a leisure companion, an interface with the administration, a social bonding tool... it's part of our lives. But on the one hand, it is far from fulfilling its promises in terms of ecology, and on the other, it contributes to locking in a certain choice of society, in which the user loses his autonomy in the face of tools. Yet a digital world that is useful, accessible, repairable, easy to dismantle, recyclable, durable and robust would be a desirable option for everyone. In this article, we review a number of socio-technical, cultural, economic, political and institutional obstacles. The idea here is not to be exhaustive, but to highlight a few structural impediments to the implementation of sobriety-enhancing actions, from the point of view of reducing the negative impacts of digital technology. On the other hand, we also observe some emerging signs that could bring about change. We conclude by highlighting the tensions arising from the locks on the one hand and the levers on the other, as revealing opposing societal visions and value systems.[5]

#### 8.1.8 Sufficient generative AI: an oxymoron?

**Participants:** Laurent Lefèvre.

The onslaught of Artificial Intelligence tools is shaking up the technological, financial and economic landscape, creating a bubble of more or less desirable possibilities and echoing the questions posed by the literature. Digital technology has entered a new phase of growth, driven mainly by the development, training and use of deep neural networks and, in particular, generative AI, which are at the root of today's spectacular results. The sheer number of parameters (several hundred billion) and the immense training corpus exploited mean that even the most seasoned of us can be bluffed. Of course, this unbridled growth is not without consequences. While it is difficult to pinpoint the precise impact of AI, mainly due to a lack of data on usage and infrastructure, we do know that it uses hardware, software, energy and human resources on a massive scale. AI tools are growing much faster than studies on their impact, and we are moving forward blindly at a time when the environmental crisis is worsening exponentially. AI is also making inroads into higher education in a variety of ways.[14]

#### 8.1.9 Fine-grained methodology to assess environmental impact of a set of digital services

**Participants:** Adrien Berthelot, Eddy Caron, Laurent Lefèvre.

The digital sector does not have a neglectable footprint, contributing to global climate change and overcommitment of planetary boundaries. To face environmental challenges, digital services could provide help and even mitigate other sector environmental impact. But, as the assessment of the



environmental impacts of digital services remains difficult, digital services are still uncertain solutions. Global assessments of the digital sector already exist but are still often uncompleted, limited to carbon emission, use phase of the equipment and are made at a too large scale to be applicable. As digital services rely on highly mutualized equipment and infrastructure, it is difficult to assess the cost of a single service. In this paper, we propose a methodology to calculate the complete environmental impact of a service. We take a service provider's perspective on how, and with what data from service hosting, to calculate and reduce the service footprint. From a host point of view, we distribute the footprint between the services, allowing to consider precisely the impacts. With such data on the environmental impacts of the hosted services, we propose strategies to prioritize optimization effort and decommissioning policies. This preliminary work is available in [20].

## 8.2 Edge, Cloud and Distributed Resource Management

### 8.2.1 SkyData: Autonomous Data paradigm

**Participants:** Eddy Caron, Elise Jeanneau, Laurent Lefèvre, Etienne Mauffret, Christian Perez.

With the rise of Data as a Service, companies understood that whoever controls the data has the power. The past few years have exposed some of the weaknesses of traditional data management systems. For example, application owner can collect and use data to their own advantage without the user's consent. We defined the SkyData concept, which revolves around autonomous data evolving in a distributed system. This new paradigm is a complete break from traditional data management systems. This paradigm is born from the many issues associated with traditional data management systems, such as resells or private information collected without consent, for example. Self managed data, or SKDs, are agents endowed with data, capabilities and goals to achieve. They are free to behave as they wish and try to accomplish their goals as efficiently as possible. They use learning algorithms to improve their decision making and learn new capabilities and services. We introduced how SKDs could be developed and provided some insight on useful capabilities.

In 2024, we explore algorithms that reach a trade-off between data autonomy and cohesion in a given subset of replicas. We propose a deterministic algorithm that ensures cohesion under a costly assumption on the number of failures. Alternatively, we investigate the use of an eventual leader election mechanism in a probabilistic algorithm where a designated leader manages coordination and acts as a communication relay. Compared to the naïve approach of broadcasting new positions to all replicas after each migration, we show experimentally that this approach reduces the loss of cohesion in most scenarios, even without assumption on the number of failures.

### 8.2.2 Numerics in the Cloud

We have defined some guidelines for critical applications to reduce the arithmetic numerical issue and we provide additional guidelines dedicated to Cloud platform.

Using a simple experiment we shown how the result of a floating-point computation can be affected when the program is compiled and executed in different environments (different processors, with different floating-point extensions and different compiler options), which is to be expected when running applications on a Cloud. Our example is simply based on floating-point summation, which is well known to be “not as easy to compute accurately as it seems” in the literature. However, this experiment is really meant to illustrate the difficulty to guarantee reproducible results, but not to exhibit real accuracy problems. With some care, porting numerical software from a micro-controller to the Cloud, or directly writing applications to the Cloud, can be achieved provided certain recommendations we have defined are followed.

We built an automated, flexible testing framework designed to evaluate the numerical stability of an application across diverse configurations and Cloud platforms. By leveraging advanced DevOps practices, this environment enables:

- The evaluation of numerical stability under varying hardware and software setups, and deployment methods (containerized or native).
- Among the many available Cloud providers, this report focuses on cross-Cloud consistency tests conducted on Grid'5000, AWS, and Azure as a proof of concept.
- The preliminary work for cost-performance analyses to identify optimal Cloud platforms.

We created a scalable pipeline was implemented using tools such as Jenkins, Terraform, Ansible, Docker, and GitLab. The framework supports detailed configuration tests and generates structured outputs for further numerical and cost-effectiveness evaluations. This testing environment forms the backbone of future analytical efforts, enabling accurate and reproducible results across multiple configurations and Cloud environments.

### 8.2.3 A specialized model and implementation of an actuarial chatbot based on Federated Learning

In this work we focused on developing a language model specialized in the actuarial domain using modern machine learning techniques, including federated learning (FL). The primary objective is to create a chatbot based on large language models (LLMs) capable of meeting actuaries' specific needs in risk modeling, financial forecasting, and other technical tasks.

Training language models is a crucial process where parameters like batch size, learning rate, and optimizers are adjusted to improve performance. Specific fine-tuning techniques, like instruction-based adaptation and alignment with human preferences via Reinforcement Learning from Human Feedback (RLHF), are employed to tailor LLMs to the actuarial field.

In the context of this work, actuarial schools and other financial institutions become clients of the federated system, contributing to the training of a decentralized actuarial model. The central server aggregates updates from local models trained on each client's private data, ensuring data remains secure throughout the process. To optimize performance and reduce communication costs between clients and the server, techniques for compressing updates are applied, including sketched and structured updates.

In this exploratory work, we have gained expertise in LLMs and federated learning.

## 8.3 HPC Applications and Runtimes

### 8.3.1 Measuring and Interpreting Dependent Task-based Applications Performances.

**Participants:** Thierry Gautier, Romain Pereira.

Breaking down the parallel time into work, idleness, and overheads is crucial for assessing the performance of HPC applications, but difficult to measure in asynchronous dependent tasking runtime systems. No existing tools allow its measurement portably and accurately. The paper [10] introduces POT: a tool-suite for dependent task-based applications performance measurement. We focus on its low-disturbance methodology consisting of task modeling, discrete-event tracing, and post-mortem simulation-based analysis. It supports the OMPT standard OpenMP specifications. The paper evaluates the precision of POT's parallel time breakdown analysis on LLVM and MPC implementations and shows that measurement bias may be neglected above 16 $\mu$ s workload per task, portably across two architectures and OpenMP runtime systems.

### 8.3.2 Evaluation of mix MPI + Dependent OpenMP task programming models on Fugaku.

**Participants:** Thierry Gautier, Romain Pereira.

The adoption of ARM processor architectures is on the rise in the HPC ecosystem. Fugaku supercomputer is a homogeneous ARMbased machine, and is one among the most powerful machine in the world.

In the programming world, dependent task-based programming models are gaining tractions due to their many advantages: dynamic load balancing, implicit expression of communication/computation overlap, early-bird communication posting. MPI and OpenMP are two widespread programming standards that make possible task-based programming at a distributed memory level. Despite its many advantages, mixed-use of the standard programming models using dependent tasks is still under-evaluated on large-scale machines. In the paper [11], we provide an overview on mixing OpenMP dependent tasking model with MPI with the state-of-the-art software stack (GCC-13, Clang17, MPC-OMP). We provide the level of performances to expect by porting applications to such mixed-use of the standard on the Fugaku supercomputers, using two benchmarks (Cholesky, HPCCG) and a proxy-application (LULESH). We show that software stack, resource binding and communication progression mechanisms are factors that have a significant impact on performance. On distributed applications, performances reaches up to 80% of efficiency for task-based applications like HPCCG. We also point-out a few areas of improvements in OpenMP runtimes.

### 8.3.3 Handling dynamicity of HPC applications designed by a task-based component model

**Participants:** Jerry Lacmou Zeutouo, Christian Perez, Thierry Gautier, Romain Pereira.

We extended the COMET component model with the support of dynamic dependencies in its data-flow model. From a meta-task based data-flow, the COMET compiler generates the OpenMP code that will submit the tasks as well as the associated dependencies. The limitation of COMET was that these dependencies were to be known when submitting the tasks of all the data-flow. Hence, a task could not depend on a value compute by another task. We have extended the COMET model with the support of dynamic dependencies and have modified accordingly its runtime. A major difficulty was to generate the code that handle those dynamically-known dependencies under HPC constraints. We evaluated the relevance and performance of three models of dependencies (flat, nested, and weak dependencies) provided by OpenMP related runtimes (LLVM, MPC, and OmpSs-2).

### 8.3.4 New high level programming framework: Experiments with Kokkos

**Participants:** Thierry Gautier, Gabriel Suau.

The paper [2] describes the implementation of DONUT, a small multi-group S N -DG transport solver that aims at providing efficient and portable sweep kernels on shared-memory architectures for Cartesian and hexagonal geometries. DONUT heavily relies on the Kokkos C++ library for portability and genericity. First encouraging performance results are presented for multicore CPU architectures.

### 8.3.5 Taskgrind: Heavyweight Dynamic Binary Instrumentation for Parallel Programs Analysis

**Participants:** Romain Pereira.

Determinacy races are concurrent programming hazards occurring when two accesses on the same memory address are not ordered, and at least one is writing. Their presence hints at a correctness error, particularly under asynchronous task-based parallel programming models. The paper [12] introduces Taskgrind: a Valgrind tool for memory access analysis of parallel programming models such as Cilk or OpenMP. We illustrate the tool's capabilities with a determinacy-race analysis and confront it with state-of-the-art tools. Results show fewer false negatives and memory overheads on a set of microbenchmarks and LULESH, with meaningful error reports toward assisting programmers when parallelizing programs.

### 8.3.6 ETH4HPC Strategic Research Agenda

**Participants:** Christian perez.

This white paper [22] is released as part of the ETP4HPC's Strategic Research Agenda 6. High-performance computing (HPC) applications achieve extreme levels of performance on large-scale systems by utilizing a wide range of tools, including compilers, runtime/middleware, APIs for memory access, debuggers and performance profilers, as well as high-level frameworks and domain-specific languages (DSLs).

## 9 Bilateral contracts and grants with industry

### 9.1 Bilateral grants with industry

**Participants:** Eddy Caron, Thierry Gautier, Laurent Lefevre.

**Bosch** We have a collaboration with Bosch and AriC (a research team of the LIP laboratory, jointly supported by CNRS, ENS de Lyon, Inria and Université Claude Bernard (Lyon 1)). We conducted a study to provide guidelines for writing portable floating-point software in Cloud environments. With some care, porting numerical software from a micro-controller to the Cloud, or directly writing applications to the Cloud, can be eased with the help of some recommendations. It is in fact not more difficult than porting software from a micro-controller to any general-purpose processor.

**CEA** We have a collaboration with CEA INSTN/SFRES / Saclay. This collaboration is based on the co-advising of a CEA PhD. The research of the PhD student (Gabriel Suau) focuses on high performance codes for neutron transport. One of the goal of the PhD is to work on better integration of Kokkos with a task based model.

**Octo technology** We have a collaboration with Octo Technology (Part of Accenture). This collaboration is sealed through a CIFRE PhD grant. The research of the PhD student (Adrien Berthelot) focuses on accelerated and driven evaluation of the environmental impacts of an Information System with the full set of digital services

**SynAApps** We have a collaboration with SynAApps (part of Cyril Group). This collaboration is sealed through a CIFRE PhD grant. The research of the PhD student (Simon Lambert) focuses on forecast and dynamic resource provisioning on a virtualization infrastructure.

## 10 Partnerships and cooperations

### 10.1 International initiatives

#### 10.1.1 Participation in other International Programs

##### JLESC

**Participants:** Thierry Gautier, Jerry Lacmou Zeutouo, Romain Pereira, Christian Perez.

**Title:** Joint Laboratory for Extreme Scale Computing

**Partner Institutions:** NCSA (US), ANL (US), Inria (FR), Jülich Supercomputing Centre (DE), BSC (SP), Riken (JP).

**Date/Duration:** 2014-

**Summary:** The purpose of the Joint Laboratory for Extreme Scale Computing (JLESC) is to be an international, virtual organization whose goal is to enhance the ability of member organizations and investigators to make the bridge between Petascale and Extreme computing. The founding partners of the JLESC are INRIA and UIUC. Further members are ANL, BSC, JSC and R-CCS. UTK is a research member. JLESC involves computer scientists, engineers and scientists from other disciplines as well as from industry, to ensure that the research facilitated by the Laboratory addresses science and engineering's most critical needs and takes advantage of the continuing evolution of computing technologies.

## SKA

**Participants:** Anass Serhani, Laurent Lefevre, Christian Perez.

**Title:** Square Kilometer Array Organization(SKA)

**Summary:** The Avalon team collaborates with SKA Organization (an IGO) whose mission is to build and operate cutting-edge radio telescopes to transform our understanding of the Universe, and deliver benefits to society through global collaboration and innovation.

## 10.2 International research visitors

### 10.2.1 Visits of international scientists

#### UiT, Tromso, Norway

**Participants:** Laurent Lefevre, Eddy Caron.

In the context of the PHC Aurora project with University of Tromso (Norway) on the topic "Exploring Energy Monitoring and Leveraging Energy Efficiency on End-to-end Worst Edge-Fog- Cloud Continuum for Extreme Climate Environments Observatories", we have received a norwegian delegation from UiT (Tromso) (one week in december 2024).

## 10.3 European initiatives

### 10.3.1 Horizon Europe

#### SLICES-PP

**Participants:** Christian Perez, Laurent Lefevre, Pierre Jacquot.

**Title:** Scientific Large-scale Infrastructure for Computing/Communication Experimental Studies - Preparatory Phase

**Partners:** Institut National de Recherche en Informatique et Automatique (INRIA), France Sorbonne Université (SU), France Universiteit van Amsterdam (UvA), Netherlands University of Thessaly (UTH), Greece Consiglio Nazionale delle Ricerche (CNR), Italy Instytut Chemii Bioorganicznej Polskiej Nauk (PSNC), Poland Mandat International (MI), Switzerland IoT Lab (IoTLAB), Switzerland Universidad Carlos III de Madrid (UC3M), Spain Interuniversitair Micro-Electronica Centrum (IMEC),

Belgium UCLan Cyprus (UCLAN), Cyprus EURECOM, France Számítástechnikai és Automatizálási Kutatóintézet (SZTAKI), Hungary Consorzio Interuniversitario Nazionale per l'Informatica (CINI), Italy Consorzio Nazionale Interuniversitario per le Telecomunicazioni (CNIT), Italy Université du Luxembourg (Uni.Lu), Luxembourg Technical Universitaet Muenchen (TUM), Germany Euskal Herriko Unibertsitatea (EHU), Spain Kungliga Tekniska Hogskolan (KTH), Sweden Oulun Yliopisto (UOULU), Finland EBOS Technologies Ltd (EBOS), Cyprus Simula Research Laboratory AS (SIMULA), Norway Centre National de la Recherche Scientifique (CNRS), France Institut Mines-Télécom (IMT), France Université de Geneve (UniGe), Switzerland

**Duration:** From September 1, 2022 to Decembre 31, 2025

**Summary:** The digital infrastructures research community continues to face numerous new challenges towards the design of the Next Generation Internet. This is an extremely complex ecosystem encompassing communication, networking, data-management and data-intelligence issues, supported by established and emerging technologies such as IoT, 5/6G, cloud-to-edge computing. Coupled with the enormous amount of data generated and exchanged over the network, this calls for incremental as well as radically new design paradigms. Experimentally-driven research is becoming worldwide a de-facto standard, which has to be supported by large-scale research infrastructures to make results trusted, repeatable and accessible to the research communities. SLICES-RI (Research Infrastructure), which was recently included in the 2021 ESFRI roadmap, aims to answer these problems by building a large infrastructure needed for the experimental research on various aspects of distributed computing, networking, IoT and 5/6G networks. It will provide the resources needed to continuously design, experiment, operate and automate the full lifecycle management of digital infrastructures, data, applications, and services. Based on the two preceding projects within SLICES-RI, SLICES-DS (Design Study) and SLICES-SC (Starting Community), the SLICES-PP (Preparatory Phase) project will validate the requirements to engage into the implementation phase of the RI lifecycle. It will set the policies and decision processes for the governance of SLICES-RI: i.e. the legal and financial frameworks, the business model, the required human resource capacities and training programme. It will also settle the final technical architecture design for implementation. It will engage member states and stakeholders to secure commitment and funding needed for the platform to operate. It will position SLICE

### 10.3.2 Other european programs/initiatives

#### Extrem Edge-Fog-Cloud continuum

**Participants:** Eddy Caron, Laurent Lefevre

**Title:** Extrem Edge-Fog-Cloud continuum

**Date/Duration:** 2023-2024

**Partner Institution:** University of Tromso, Norway

**Summary:** Laurent Lefevre co-leads a PHC Aurora project with University of Tromso (Norway) on the topic " Exploring Energy Monitoring and Leveraging Energy Efficiency on End-to-end Worst Edge-Fog- Cloud Continuum for Extreme Climate Environnements Observatories". The aim of this collaboration is to explore a best effort end-to-end energy monitoring system for a worst-case continuum, to discuss a long-term infrastructure continuum for the edge and to design an experimental validation of usable energy leverages at the edge.

## 10.4 National initiatives

### Priority Research Programmes and Equipments (PEPR)

### PEPR Cloud – Taranis

**Participants:** Christian Perez, Yves Caniou, Eddy Caron, Elise Jeanneau, Laurent Lefevre, Johanna Desprez.

**Title:** Taranis : Model, Deploy, Orchestrate, and Optimize Cloud Applications and Infrastructure

**Partners:** Inria, CNRS, IMT, U. Grenoble-Alpes, CEA, U. Rennes, ENSL, U. Lyon I, U. Lille, INSA Lyon, INSA Rennes, Grenoble INP

**Date:** Sep 2023 – Aug 2030.

**Summary:** New infrastructures, such as Edge Computing or the Cloud-Edge-IoT computing continuum, make cloud issues more complex as they add new challenges related to resource diversity and heterogeneity (from small sensor to data center/HPC, from low power network to core networks), geographical distribution, as well as increased dynamicity and security needs, all under energy consumption and regulatory constraints.

In order to efficiently exploit new infrastructures, we propose a strategy based on a significant abstraction of the application structure description to further automate application and infrastructure management. Thus, it will be possible to globally optimize the resources used with respect to multi-criteria objectives (price, deadline, performance, energy, etc.) on both the user side (applications) and the provider side (infrastructures). This abstraction also includes the challenges related to the abstraction of application reconfiguration and to automatically adapt the use of resources.

The Taranis project addresses these issues through four scientific work packages, each focusing on a phase of the application lifecycle: application and infrastructure description models, deployment and reconfiguration, orchestration, and optimization.

### PEPR Cloud – CareCloud

**Participants:** Laurent Lefevre, Eddy Caron, Olivier Glück, Thomas Stavis.

**Title:** Understanding, improving, reducing the environmental impacts of Cloud Computing

**Partners:** CNRS, Inria, Univ. Toulouse, IMT

**Date:** Sept 2023 - Aug 2030

**Summary:** The CARECloud project (understanding, improving, reducing the environmental impacts of Cloud Computing) aims to drastically reduce the environmental impacts of cloud infrastructures. Cloud infrastructures are becoming more and more complex: both in width, with more and more distributed infrastructures, whose resources are scattered as close as possible to the user (edge, fog, continuum computing) and in depth, with an increasing software stacking between the hardware and the user's application (operating system, virtual machines, containers, orchestrators, micro-services, etc.) The first objective of the project is to understand how these infrastructures consume energy in order to identify sources of waste and to design new models and metrics to qualify energy efficiency. The second objective focuses on the energy efficiency of cloud infrastructures, i.e., optimizing their consumption during the usage phase. In particular, this involves designing resource allocation and energy lever orchestration strategies: mechanisms that optimize energy consumption (sleep modes, dynamic adjustment of the size of virtual resources, optimization of processor frequency, etc.). Finally, the third objective targets digital sobriety in order to sustainably reduce the environmental impact of clouds. Indeed, current clouds offer high availability and very high fault tolerance, at the cost of significant energy expenditure, particularly due to redundancy and oversizing. This third objective aims to design infrastructures that are more energy and IT resource efficient, resilient to electrical intermittency, adaptable to the production of electricity from renewable energy sources and tolerant of the disconnection of a highly decentralized part of the infrastructure



**PEPR Cloud – Silecs**

**Participants:** Simon Delamare, Pierre Jacquot, Laurent Lefevre, Christian Perez.

**Title:** Super Infrastructure for Large-Scale Experimental Computer Science for Cloud/Edge/IoT

**Partners:** Inria, CNRS, IMT, U. Lille, INSA Lyon, U. Strasbourg, U. Grenoble-Alpes, Sorbonne U., U. Toulouse, Nantes U., Renater.

**Date:** Sept 2023 - Aug 2030

**Summary:** The infrastructure component of the PEPR Cloud (SILECS) will structure the Cloud/Fog/Edge/IoT aspects of the SLICES-FR (Super Infrastructure for Large-Scale Experimental Computer Science) platform, the French node of the ESFRI SLICES-RI action. SILECS will enable the prototyping and conduct of reproducible experiments of any hardware and software element of current and future digital environments at all levels of the Cloud IoT continuum, addressing the experimental needs of the other PEPR components. SILECS will be complemented within SLICES-FR by funding from the PEPR Networks of the Future, which focuses on specific aspects of 5G and beyond technologies. There will therefore be continuous and coordinated strong interactions between the two PEPRs

**PEPR 5G Network of the Future – JEN**

**Participants:** Laurent Lefevre, Doreid Ammar, Emile Egreteau- -Druet.

**Title:** JEN: Network of the Future – Just Enough Networks

**Partners:** CEA, CNRS, ENSL, ESIEE, IMT, INPB, Inria, INSAL

**Date:** 2023-2028

**Summary:** In the NF-JEN project, partners propose to develop just enough networks: network whose dimension, performance, resource usage and energy consumption are just enough to satisfy users' needs. Along with designing energy-efficient and sober networks, we will provide multi-indicators models that could help policy-makers and inform the public debate.

**PEPR NumPEX – Exa-Soft**

**Participants:** Thierry Gautier, Christian Perez, Jerry Lacmou Zeutouo, Pierre-Etienne Polet.

**Title:** Exa-Soft: HPC software and tools

**Partners:** Inria, CEA, CNRS, U. Paris-Saclay, Telcom SudParis, Bordeaux INP, ENSIIE, U. Bordeaux, U. Grenoble-Alpes, U. Rennes I, U. Strasbourg, U. Toulouse

**Date:** 2023-2029

**Summary:** Though significant efforts have been devoted to the implementation and optimization of several crucial parts of a typical HPC software stack, most HPC experts agree that exascale supercomputers will raise new challenges, mostly because the trend in exascale compute-node hardware is toward heterogeneity and scalability: Compute nodes of future systems will have a combination of regular CPUs and accelerators (typically GPUs), along with a diversity of GPU architectures.

Meeting the needs of complex parallel applications and the requirements of exascale architectures raises numerous challenges which are still left unaddressed.



As a result, several parts of the software stack must evolve to better support these architectures. More importantly, the links between these parts must be strengthened to form a coherent, tightly integrated software suite.

Our project aims at consolidating the exascale software ecosystem by providing a coherent, exascale-ready software stack featuring breakthrough research advances enabled by multidisciplinary collaborations between researchers.

## ANR

### SkyData

**Participants:** Eddy Caron, Elise Jeanneau, Etienne Mauffret, Lucien Ndjie Ngale, Laurent Lefèvre, Christian Perez, Maxime Just.

**Title:** SkyData: A new data paradigm: Intelligent and Autonomous Data

**Partners:** LIP, VERIMAHG, LIP6

**Date:** 01.2023-01.2027.

**Summary:** Nowadays, who controls the data, controls the world or at least the IT world. Usually Data are managed through a middleware, but in this project, we propose a new data paradigm without any data manager. We want to endow the data with autonomous behaviors and thus create a new entity, so-called Self-managed data. We plan to develop a distributed and autonomous environment, that we call SKYDATA, where the data are regulated by themselves. This change of paradigm represents a huge and truly innovative challenge! This goal must be built on the foundation of a strong theoretical study and knowledge on autonomic computing, since Self-managed data will now have to obtain and compute the services they need in autonomy. We also plan to actually develop a SKYDATA framework prototype and a green-IT use case that focuses data energy consumption. SKYDATA will be compliant with GDPR through the targeted datas and some internal process.

## French Joint Laboratory

### ECLAT

**Participants:** Anass Serhani, Laurent Lefèvre, Christian Perez.

**Partner Institution(s):** CNRS, Inria, Eviden, Observatoire de la Côte d'Azur, Observatoire de Paris-PSL

**Date/Duration:** 2023-

**Summary** ECLAT is a joint laboratory gathering 14 laboratories to support the French contribution to the SKAO observatory.

## Inria Large Scale Initiative

### FrugalCloud: Défi Inria OVHCloud

**Participants:** Eddy Caron, Laurent Lefèvre, Christian Perez.

**Summary** A joint collaboration between Inria and OVH Cloud company on the topic challenge of frugal cloud has been launched in October 2021. It addresses several scientific challenge on the eco-design of cloud frameworks and services for large scale energy and environmental impact reduction. Laurent Lefèvre is the scientific animator of this project. Some Avalon PhD students are involved in this Inria Large Scale Initiative (Défi) : Maxime Agusti and Vladimir Ostanpenko.

## Alt-Impact program

**Participants:** Laurent Lefèvre, Emeline Pegon.

**Summary** Alt Impact is a program supported by ADEME, CNRS and INRIA, designed to raise public awareness of the environmental impact of digital technology. Our mission is to provide information in a clear and accessible way, with verified, up-to-date and entertaining content. In addition to providing information, we offer practical solutions that are easy to implement on a day-to-day basis, so that everyone, whether an individual or an organization, and whatever their level of knowledge, can take concrete action to reduce their digital ecological footprint.

At the same time, we are pursuing our objective of accelerating and supporting the transition to digital sufficiency, with a focus on measuring and managing it, through the identification and sharing of reliable data and tools, as well as supporting actions aimed at integrating digital sufficiency into the strategies of local authorities and businesses.

## 11 Dissemination

### 11.1 Promoting scientific activities

#### 11.1.1 Scientific events: organisation

##### General chair, scientific chair

- Laurent Lefevre was co General Chair of the [PAISE 2024: 6th Workshop on Parallel AI and Systems for the Edge](#), during the IPDPS2024 conference (San Francisco, USA, May 27-31 2024)

##### Member of the organizing committees

- Christian Perez was member of the Organizing Committee of [JCAD, the French Journées Calcul Données](#) (Bordeaux, 4-6 Nov 2024).
- Laurent Lefevre was co-organizer of the [FACTO Camp: discovering IoT"](#) (St Paul en Jarez, France, October 20-25, 2024)
- Laurent Lefevre was co organizer of the [Eco-ICT 2024 school : GDR RSD / FrugalCloud challenge school on Energy and Environmental Impacts reduction in Digital Services, Networks and Distributed Systems](#), (Quimper, Fouesnant, October 7-11 2024)
- Laurent Lefevre was co-organizer of the [GreenDays2024@Toulouse: "Exploring multi facets of digital frugality"](#) (Toulouse, France, March 27-28, 2024)

#### 11.1.2 Scientific events: selection

##### Member of the conference program committees

- Yves Caniou was a member of the International Conference on Computational Science and Its Applications.
- Eddy Caron was a member of the program committee for the conferences: CLOSER 2024 and ICCS 2024
- Élise Jeanneau was a member of the program committee for the conference AlgoTel 2024: 26èmes Rencontres Francophones sur les Aspects Algorithmiques des Télécommunications.
- Laurent Lefevre was a member of the program committee for the conference IC2E 2024: 12th IEEE International Conference on Cloud Engineering, Euro-par 2024: 30th International European Conference on Parallel and Distributed Computing, ICDCS 2024: 44th IEEE International Conference on Distributed Computing Systems and ICT4S 2024: The International Conference on ITC for Sustainability.

### 11.1.3 Journal

#### Reviewer - reviewing activities

- Christian Perez was reviewer for the journal Future Generation Computer Systems.
- Eddy Caron was reviewer for:
  - Future Generation Computer Systems,
  - Discrete Mathematics and Theoretical Computer Science
  - Engineering Applications of Artificial Intelligence

### 11.1.4 Invited talks

- Christian Perez gave a talk about SLICES at 8th Symposium on Interoperability of Supercomputing and Cloud Technologies, SC'24, Atlanta, USA, November 22nd, 2024.
- Laurent Lefevre Perez gave the following talks :
  - "Energy efficiency and environmental impacts of large scale parallel and distributed systems - Exploring frugality at large scale", Seminar in University of Tromso, Norway, November 18, 2024
  - "Beaucoup de verrous, peu de leviers : à quand un numérique plus sobre ?", Françoise Berthoud, Laurent Lefevre, NEC2024: Numérique en Commun[s], Chambéry, September 25, 2024
  - "Contre les impacts environnementaux des grandes infrastructures numériques : leviers, IA et cycle de vie", Laurent Lefevre, Adrien Berthelot, Mathilde Jay, Festival du Numérique Responsable, Ecole des Mines de St Etienne (online), March 19, 2024
  - "Accueil de groupes de lycéennes à l'Ecole Normale Supérieure de Lyon", Journée internationale des droits des femmes, École Normale Supérieure de Lyon, March 8, 2024
- Thierry Gautier gave a talk "XKBlas report on AMD multi-GPU and its evolutions" for the MUMPS 2024 annual meeting at Ansys, Villeurbanne, 2024-06-20
- Pierre-Etienne Polet gave a talk "Executing MUMPS & XKBlas on Unified Memory Platforms: Methodology and Initial Experiments on Nvidia Grace Hopper" at EDF, Paris, 2024-11-26

### 11.1.5 Scientific expertise

- Yves Caniou evaluated 2 projets for comité d'évaluation of the «CE46 - Calcul haute performance, Modèles numériques, simulation, applications» and «CE25 - Sciences et génie du logiciel - Réseaux de communication multi-usages, infrastructures numériques».
- Eddy Caron was a member of the ANR CE25 committee and the ANR ASTRID committee.
- Eddy Caron was a member of the committee PHELMA COS MCF 27 for the LIG (Grenoble).
- Élise Jeanneau was a member of the juries for recruiting CRCN candidates in the Inria Saclay and Inria Grenoble centers.
- Christian Perez evaluated 6 projects for the French Direction générale de la Recherche et de l'Innovation. He was a member of the jury for recruiting CRCN candidates in Inria Bordeaux center.

### 11.1.6 Research administration

- Eddy Caron is the leader of the **SkyData** ANR project.
- Élise Jeanneau is a member of the **Inria Evaluation Committee**.
- Laurent Lefevre is the leader of the **FrugalCloud challenge (défi)** between Inria and OVHcloud company. He is a member of the executive board and the sites committee of the **Grid'5000** Scientific Interest Group and member of the executive board and the sites committee of the **SLICES-FR** testbed.
- Christian Perez represents Inria in the overview board of the **France Grilles** Scientific Interest Group. He is a member of the executive board and the sites committee of the **Grid'5000** Scientific Interest Group and member of the executive board of the **SLICES-FR** testbed. He is a member of the **Inria Lyon Strategic Orientation Committee**. He is in charge of organizing scientific collaborations between **Inria** and **SKA France**.

## 11.2 Teaching - Supervision - Juries

### 11.2.1 Teaching

- Licence: Yves Caniou, Algorithmique programmation impérative initiation, 129h, niveau L1, Université Claude Bernard Lyon 1, France.
- Licence: Yves Caniou, Algorithmique et programmation récursive, 54h, niveau L1, Université Claude Bernard Lyon 1, France.
- Licence: Yves Caniou, Programmation Concurrente, 24h and Co-Responsable of UE, niveau L3, Université Claude Bernard Lyon 1, France.
- Licence: Yves Caniou, Réseaux, 12h, niveau L3, Université Claude Bernard Lyon 1, France.
- Licence: Yves Caniou, Systèmes d'information documentaire, 22h, niveau L3, Université Claude Bernard Lyon 1, France.
- Master: Yves Caniou, Responsable of alternance students, 21h, niveau M1, Université Claude Bernard Lyon 1, France.
- Master: Yves Caniou, Sécurité, 24h and Responsable of UE, niveau M2, Université Claude Bernard Lyon 1, France.
- Master: Yves Caniou, Responsable of alternance students, 8h, niveau M2, Université Claude Bernard Lyon 1, France
- Master: Eddy Caron, Distributed System, 20h, M1, École Normale Supérieure de Lyon. France.
- Master: Eddy Caron, Langages, concepts et archi pour les données, 30h, M2, ISFA. Université Claude Bernard Lyon 1
- Master: Eddy Caron, Risques dans les Systèmes et Réseaux - Cloud, 15h, M2, ISFA. Université Claude Bernard Lyon 1.
- Master: Eddy Caron, Service Web et Sécurité, 15h, M2, ISFA. Université Claude Bernard Lyon 1.
- Master: Eddy Caron, Data Mining Avancé: Environnements parallèles et distribués, 12h, M2, ISFA. Université Claude Bernard Lyon 1.
- Licence: Élise Jeanneau, Introduction Réseaux et Web, 36h, niveau L1, Université Lyon 1, France.
- Licence: Élise Jeanneau, Réseaux, 54h, niveau L3, Université Lyon 1, France.
- Licence: Élise Jeanneau, Algorithmique, programmation et structures de données, 24h, niveau L2, Université Lyon 1, France

- Licence: Élise Jeanneau, Architecture des ordinateurs, 24h, niveau L2, Université Lyon 1, France
- Licence: Élise Jeanneau, Réseaux, systèmes et sécurité par la pratique, 23h, niveau L3, Université Lyon 1, France
- Master: Élise Jeanneau, Algorithmes distribués, 45h, niveau M1, Université Lyon 1, France.
- Master: Élise Jeanneau, Réseaux, 6h, niveau M1, Université Lyon 1, France.

### 11.2.2 Supervision

PhD defended:

- Mathilde Jay. "A Versatile Methodology for Assessing the Electricity Consumption and Environmental Footprint of Machine Learning Training: from Supercomputers to Edge Devices", 2021, Laurent Lefevre (co-dir. Inria. Avalon), Denis Trystram (dir. LIG, UGA), defended 2024, October 15th.
- Adrien Berthelot. "The complete environmental footprint of digital usage: a contribution to life-cycle analysis of digital services", 2021, Eddy Caron (dir ENS de Lyon. Inria. Avalon), Laurent Lefevre (co-dir Inria. Avalon), Alexis Nicolas (Octo Technology), defended 2024, November 12th.
- Vladimir Ostapenco. "Modeling, evaluating and orchestrating heterogeneous leverages for large-scale cloud data centers management", 2021, Laurent Lefevre (dir. Inria. Avalon), Anne-Cécile Orgerie (co-dir. Inria. Myriads), Benjamin Fichel (co-dir. OVHcloud), defended 2024, December 18th.

Phd in progress:

- Maxime Agusti. "Observation de plate-formes de co-localisation baremetal, modèles de réduction énergétique et proposition de catalogues", FrugalCloud Inria-OVHCloud collaboration, Feb 2022, Eddy Caron (co-dir. ENS de Lyon. Inria. Avalon), Benjamin Fichel (co-dir. OVHcloud), Laurent Lefevre (dir. Inria. Avalon) et Anne-Cécile Orgerie (co-dir. Inria. Myriads),.
- Simon Lambert. "Forecast and dynamic resource provisioning on a virtualization infrastructure", 2022, Eddy Caron (dir. ENS de Lyon. Inria. Avalon), Laurent Lefevre (co-dir Inria. Avalon), Rémi Grivel (co-dir. Ciril Group).
- Gabriel Suaus. "Résolution de l'équation de transport des neutrons sur des architectures massivement parallèles et hétérogènes : application aux géométries hexagonales. Thierry Gautier (dir. INRIA Avalon), Ansar CALLOO (co-dir. CEA), Romain LE TELLIER (co-dir CEA), Remi LE BARON (co-dir CEA).
- Thomas Stavis. "Replay of environmental leverages in cloud infrastructures and continuums", 2024, Laurent Lefevre (dir. Inria Avalon), Anne-Cécile Orgerie (co-dir CNRS, Magellan team, Irisa Rennes)
- Émile Egreteau-bruet. "Analyzing full life cycle of IoT based 5G solutions for smart agriculture", 2024, Laurent Lefevre (dir. Inria Avalon), Nathalie Mitton (co-dir. Inria Lille) and Doreid Ammar (co-dir. Aivancity Inria Avalon)

Eddy Caron was supervisor or co-supervisor of the following internship:

- Maxime Just (PLR 4th year of ENS). *A study to provide a set of guidelines for calculation and cost-efficient of numerical algorithms in the cloud.*
- Hamza Aabirouche (M2). *Étude visant à établir des recommandations pour le calcul et l'optimisation de la stabilité des algorithmes numériques dans le Cloud*
- Annour Saad Allamine (M2). *SkyISFA : Adaptation d'un modèle spécialisé et mise en œuvre d'un chatbot pour l'actuariat basé sur le Federated Learning.*
- Maxime Just (M2). *Évaluation et analyse du consensus dans un contexte de données autonomes.* Co-supervised with Elise Jeanneau.

### 11.2.3 Juries

- Eddy Caron was member of the HDR defense committee of Antoine Boutet (INSA. Lyon) : "Privacy issues in AI and geolocation : from data protection to user awareness", Lyon, December 10th, 2024.
- Eddy Caron was PhD reviewer and member of the defense committee of Thomas Bouvier (IRISA Rennes): "Distributed Rehearsal Buffers for Continual Learning at Scale", Rennes, November 4th, 2024.
- Eddy Caron was member of the Phd defense committee of Sylvain Lejambre (Université Mont Blanc Savoie. Annecy): "Conception et mise en œuvre d'un système de représentation et de manipulation de la connaissance des Objets Sages (WO) pour la détection de fraudes au sein des systèmes bancaires". Annecy. December 12th, 2024.
- Laurent Lefevre was PhD reviewer of the PhD and member of the defense committee of :
  - Manal Benaissa : "Optimization of the IT and Electrical Sizing Process of a Green, Medium-Sized Datacenter Powered Exclusively by Renewable Energy Sources, Considering Environmental and External Uncertainty", University of Franche-Comté, Besançon, December 10, 2024
  - Vincent Lannurien : "Allocation et placement dynamiques sur ressources hétérogènes pour le cloud serverless", ENSTA Bretagne, Brest, France, November 20, 2024
  - Mohammad Sadegh Aslanpour : "Serverless Edge Computing: Resource Scheduling Algorithms and Software Systems", Monash University, Australia, January 8, 2024
- Laurent Lefevre was member of the defense committee of Klervie Toczé : "Orchestrating a Resource-aware Edge", Linköping University, Sweden, October 4, 2024
- Christian Perez was reviewer and member of the HDR defense committee of Simon Bliudze: "Rigorous Design of Concurrent Component-Based Software and Systems", Lille, March 8th, 2024.
- Christian Perez was PhD reviewer and member of the defense committee of Ophélie Renaud: "Model Based Granularity Optimization for High Performance Computing Systems in Astronomy", Rennes, October 9th, 2024.

## 11.3 Popularization

### 11.3.1 Productions (articles, videos, podcasts, serious games, ...)

- Yves Caniou co-organized the **Campus du libre** with institutes Université Claude Bernard Lyon 1, Université Lyon 2, Université Lyon 3, INSA, Inria.
- Laurent Lefevre participated in the Pop'Sciences project from University of Lyon for the comic book on with middle and high school students from Lyon and Villeurbanne, April 2024

### 11.3.2 Participation in Live events

- Laurent Lefevre was panelist in the GreenTech Forum 2024 event: "Avancées et nouvelles technologies pour des datacenters plus responsables / Advances and new technologies for more responsible data centers", GreenTech Forum, Paris, France, November 6, 2024

## 12 Scientific production

### 12.1 Publications of the year

#### International journals

- [1] A. Berthelot, E. Caron, M. Jay and L. Lefevre. 'Understanding the environmental impact of generative AI services'. In: *Communications of the ACM Special Issue on Sustainability and Computing* (2025). URL: <https://hal.science/hal-04920612>. In press (cit. on p. 12).

- [2] G. Suau, A. Calloo, R. Baron, R. Le Tellier and T. Gautier. 'Efficient sweep kernels on shared-memory architectures for the discrete ordinates neutron transport equation on Cartesian and hexagonal geometries'. In: *EPJ Web of Conferences* 302 (15th Oct. 2024), p. 02009. DOI: [10.1051/epjconf/202430202009](https://doi.org/10.1051/epjconf/202430202009). URL: <https://hal.science/hal-04908404> (cit. on p. 16).

### International peer-reviewed conferences

- [3] M. Agusti, E. Caron, B. Fichel, L. Lefèvre, O. Nicol and A.-C. Orgerie. 'PowerHeat: A non-intrusive approach for estimating the power consumption of bare metal water-cooled servers'. In: 2024 IEEE International Conferences on Internet of Things (iThings) and IEEE Green Computing & Communications (GreenCom) and IEEE Cyber, Physical & Social Computing (CPSCom) and IEEE Smart Data (SmartData) and IEEE Congress on Cybermatics. Copenhagen, Denmark, 2024, pp. 1–7. URL: <https://hal.science/hal-04662683> (cit. on p. 10).
- [4] A. Berthelot, E. Caron, M. Jay and L. Lefèvre. 'Estimating the environmental impact of Generative-AI services using an LCA-based methodology'. In: *Procedia CIRP*. CIRP LCE 2024 - 31st Conference on Life Cycle Engineering. Turin, Italy, 2024, pp. 1–10. URL: <https://inria.hal.science/hal-04346102> (cit. on p. 12).
- [5] F. Berthoud, M. Olivi and L. Lefèvre. 'Beaucoup de verrous, peu de leviers ! Sobriété numérique : le cas est grave mais pas désespéré'. In: *NUDGIS by UBICAST*. JRES 2024 - Journées réseaux de l'enseignement et de la recherche. Rennes, France, 13th Dec. 2024. URL: <https://hal.science/hal-04893986> (cit. on p. 13).
- [6] C. Cérin, M. Jay and L. Lefèvre. 'A Methodology and a Toolbox to Explore Dataset related to the Environmental Impact of HTTP Requests'. In: 2023 IEEE International Conference on Big Data (BigData) - 3rd International Workshop on Big Data Analytics for Sustainability. Sorrento (Naples), Italy: IEEE, 2024, pp. 1–10. DOI: [10.1109/BigData59044.2023.10386275](https://doi.org/10.1109/BigData59044.2023.10386275). URL: <https://inria.hal.science/hal-04386964> (cit. on p. 12).
- [7] H. Hadjur, D. Ammar and L. Lefèvre. 'Deep Reinforcement Learning for Energy-efficient Selection of Embedded Services at the Edge'. In: *2024 IEEE International Conferences on Internet of Things (iThings)*. 2024 IEEE International Conferences on Internet of Things (iThings). Copenhagen, Denmark: IEEE, 19th Aug. 2024, pp. 67–74. DOI: [10.1109/iThings-GreenCom-CPSCom-SmartData-Cybermatics62450.2024.00034](https://doi.org/10.1109/iThings-GreenCom-CPSCom-SmartData-Cybermatics62450.2024.00034). URL: <https://inria.hal.science/hal-04708697> (cit. on p. 11).
- [8] S. Lambert, E. Caron, L. Lefèvre and R. Grivel. 'S-ORCA : A social-based consolidation approach to reduce Cloud infrastructures energy consumption'. In: 15th IEEE International Conference on Cloud Computing Technology and Science - Cloudcom 2024. Abu Dhabi, United Arab Emirates, 9th Dec. 2024. URL: <https://inria.hal.science/hal-04797304> (cit. on p. 11).
- [9] V. Ostapenco, L. Lefèvre, A.-C. Orgerie and B. Fichel. 'Exploring RAPL as a Power Capping Leverage for Power-Constrained Infrastructures'. In: ICA3PP 2024 - 24th International Conference on Algorithms and Architectures for Parallel Processing. Macau SAR, China, 2024, pp. 1–10. URL: <https://hal.science/hal-04742418> (cit. on p. 13).
- [10] R. Pereira, T. Gautier, A. Roussel and P. Carribault. 'Measuring and Interpreting Dependent Task-based Applications Performances'. In: *Parallel Processing and Applied Mathematics 2024*. 15th International Conference on Parallel Processing & Applied Mathematics. Ostrava, Czech Republic, 4th Apr. 2025. URL: <https://hal.science/hal-04767262> (cit. on p. 15).
- [11] R. Pereira, A. Roussel, M. Tsuji, P. Carribault, M. Sato, H. Murai and T. Gautier. 'An Overview on Mixing MPI and OpenMP Dependent Tasking on A64FX'. In: International Workshop on Arm-based HPC 2024 (IWAHPCE-2024). Nagoya, Japan: ACM, 11th Jan. 2024, pp. 7–16. DOI: [10.1145/3636480.3637094](https://doi.org/10.1145/3636480.3637094). URL: <https://hal.science/hal-04370966> (cit. on p. 16).
- [12] R. Pereira, G. Stelle and P. Carribault. 'Taskgrind: Heavyweight Dynamic Binary Instrumentation for Parallel Programs Analysis'. In: *International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*. 8th International Workshop on Software Correctness for HPC Applications (Correctness '24). Atlanta (USA), United States, 17th Nov. 2024. URL: <https://hal.science/hal-04814885> (cit. on p. 16).

- [13] T. Simon, D. Ekchajzer, A. Berthelot, E. Fourboul, S. Rince and R. Rouvoy. 'BoaviztAPI: a bottom-up model to assess the environmental impacts of cloud services'. In: HotCarbon'24 - 3rd Workshop on Sustainable Computer Systems. Santa Cruz, United States, 9th July 2024. URL: <https://hal.science/hal-04621947>.

#### Conferences without proceedings

- [14] S. Bouveret, E. Frenoux, L. Lefèvre, D. Mallarino and D. Trystam. 'IA générative sobre : un oxymore ?' In: JRES (Journées réseaux de l'enseignement et de la recherche) 2024. Rennes, France, 13th Dec. 2024. URL: <https://hal.science/hal-04893830> (cit. on p. 13).

#### Scientific book chapters

- [15] L. Ngale, E. Caron and Y. Zhang. 'Fog-Robotics Infrastructures Simulation-Based Sizing Approach'. In: *Cloud Computing and Services Science : 12th International Conference, CLOSER 2022, Virtual Event, April 27–29, 2022, and 13th International Conference, CLOSER 2023, Prague, Czech Republic, April 26–28, 2023, Revised Selected Papers*. Vol. 1845. Communications in Computer and Information Science. Springer Nature Switzerland, 15th Aug. 2024, pp. 188–211. DOI: [10.1007/978-3-031-68165-3\\_10](https://doi.org/10.1007/978-3-031-68165-3_10). URL: <https://u-picardie.hal.science/hal-04807591>.

#### Doctoral dissertations and habilitation theses

- [16] A. Berthelot. 'The complete environmental footprint of digital usage : a contribution to life-cycle analysis of digital services'. Ecole normale supérieure de lyon - ENS LYON, 12th Nov. 2024. URL: <https://theses.hal.science/tel-04874694> (cit. on p. 12).
- [17] M. Jay. 'A Versatile Methodology for Assessing the Electricity Consumption and Environmental Footprint of Machine Learning Training: from Supercomputers to Edge Devices'. Université Grenoble Alpes, 15th Oct. 2024. URL: <https://theses.hal.science/tel-04907220> (cit. on p. 12).
- [18] V. Ostapenco. 'Modeling, evaluating and orchestrating heterogeneous leverages for large-scale cloud data centers management'. Ecole Normal Supérieure de Lyon - ENS Lyon, 18th Dec. 2024. URL: <https://hal.science/tel-04936934> (cit. on p. 13).
- [19] P.-E. Polet. 'Porting sonar processing chains to heterogeneous architecture : design and evaluation of a moldable task-based programming environment'. Ecole normale supérieure de lyon - ENS LYON, 3rd Apr. 2024. URL: <https://theses.hal.science/tel-04633261>.

#### Reports & preprints

- [20] A. Berthelot, E. Caron, R. de Laage, L. Lefèvre and A. Nicolas. *Fine-grained methodology to assess environmental impact of a set of digital services*. 2024. URL: <https://hal.science/hal-04928998> (cit. on p. 14).

#### Other scientific publications

- [21] A. Berthelot, E. Caron, M. Jay and L. Lefèvre. 'Towards a multi-criteria evaluation of the environmental footprint of generative ai services'. In: ICT4S 2024 - International Conference on Information and Communications Technology for Sustainability. Stockholm, Sweden, 2024, pp. 1–1. URL: <https://inria.hal.science/hal-04586653> (cit. on p. 12).
- [22] P. Carpenter, G. Antoniu, M. Arenaz, O. Aumage, J. Beránek, A. Buttari, A. Costan, S. Happ, V. Kannan, C. Perez, A. Peña, A. Scionti, X. Vigouroux and P. Viviani. *ETP4HPC SRA White Paper - Programming Environment*. Dec. 2024. DOI: [10.5281/zenodo.14446622](https://doi.org/10.5281/zenodo.14446622). URL: <https://hal.science/hal-04905035> (cit. on p. 17).
- [23] S. Lambert, E. Caron, L. Lefèvre and R. Grivel. 'Evaluating the involvement of users for reducing energy consumption of datacenters in a Cloud company'. In: ICT4S 2024 - ICT for Sustainability. Stockholm, Sweden, 24th June 2024. URL: <https://inria.hal.science/hal-04588604> (cit. on p. 11).



## 12.2 Cited publications

- [24] R. Ge, X. Feng, S. Song, H.-C. Chang, D. Li and K. W. Cameron. ‘PowerPack: Energy Profiling and Analysis of High-Performance Systems and Applications’. In: *IEEE Trans. Parallel Distrib. Syst.* 21.5 (May 2010), pp. 658–671. DOI: [10.1109/TPDS.2009.76](https://doi.org/10.1109/TPDS.2009.76). URL: [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=4906989](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4906989) (cit. on p. 4).
- [25] A. Geist and S. Dosanjh. ‘IESP Exascale Challenge: Co-Design of Architectures and Algorithms’. In: *Int. J. High Perform. Comput. Appl.* 23.4 (Nov. 2009), pp. 401–402. DOI: [10.1177/1094342009347766](https://doi.org/10.1177/1094342009347766). URL: <http://dx.doi.org/10.1177/1094342009347766> (cit. on p. 5).
- [26] W. Gropp, S. Huss-Lederman, A. Lumsdaine, E. Lusk, B. Nitzberg, W. Saphir and M. Snir. *MPI: The Complete Reference – The MPI-2 Extensions*. 2nd ed. Vol. 2. ISBN 0-262-57123-4. The MIT Press, Sept. 1998 (cit. on p. 5).
- [27] H. Kimura, T. Imada and M. Sato. ‘Runtime Energy Adaptation with Low-Impact Instrumented Code in a Power-Scalable Cluster System’. In: *Proceedings of the 2010 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing. CCGRID ’10*. Washington, DC, USA: IEEE Computer Society, 2010, pp. 378–387 (cit. on p. 4).
- [28] G. Madec. *NEMO ocean engine*. Note du Pole de modélisation 27. ISSN No 1288-1619. France: Institut Pierre-Simon Laplace (IPSL), 2008 (cit. on p. 7).
- [29] OpenACC. *The OpenACC Application Programming Interface*. Version 1.0. Nov. 2011. URL: <http://www.openacc-standard.org> (cit. on p. 5).
- [30] OpenMP Architecture Review Board. *OpenMP Application Program Interface*. Version 3.1. July 2011. URL: <http://www.openmp.org> (cit. on p. 5).
- [31] B. Rountree, D. K. Lownenthal, B. R. de Supinski, M. Schulz, V. W. Freeh and T. Bletsch. ‘Adagio: Making DVS Practical for Complex HPC Applications’. In: *Proceedings of the 23rd international conference on Supercomputing. ICS ’09*. New York, NY, USA: ACM, 2009, pp. 460–469 (cit. on p. 4).
- [32] C. Szyperski. *Component Software - Beyond Object-Oriented Programming*. 2nd ed. Addison-Wesley / ACM Press, 2002, p. 608 (cit. on p. 5).
- [33] S. Valcke. ‘The OASIS3 coupler: a European climate modelling community software’. In: *Geoscientific Model Development* 6 (2013). doi:10.5194/gmd-6-373-2013, pp. 373–388 (cit. on p. 7).