

RESEARCH CENTRE

**Inria Centre  
at Rennes University**

IN PARTNERSHIP WITH:  
Université de Rennes

2023

ACTIVITY REPORT

Project-Team

PACAP

## **Pushing Architecture and Compilation for Application Performance**

IN COLLABORATION WITH: Institut de recherche en informatique et  
systèmes aléatoires (IRISA)

### **DOMAIN**

**Algorithmics, Programming, Software and  
Architecture**

### **THEME**

**Architecture, Languages and Compilation**

The Inria logo is a stylized, cursive script in red, positioned in the bottom right corner of the page.

# Contents

<b>Project-Team PACAP</b>	<b>1</b>
<b>1 Team members, visitors, external collaborators</b>	<b>2</b>
<b>2 Overall objectives</b>	<b>3</b>
<b>3 Research program</b>	<b>5</b>
3.1 Motivation	5
3.1.1 Technological constraints	5
3.1.2 Evolving community	5
3.1.3 Domain constraints	6
3.2 Research Objectives	6
3.2.1 Static Compilation	6
3.2.2 Software Adaptation	7
3.2.3 Research directions in uniprocessor micro-architecture	7
3.2.4 Towards heterogeneous single-ISA CPU-GPU architectures	9
3.2.5 Real-time systems	9
3.2.6 Power efficiency	9
3.2.7 Security	10
<b>4 Application domains</b>	<b>11</b>
4.1 Domains	11
<b>5 Social and environmental responsibility</b>	<b>11</b>
5.1 Impact of research results	11
<b>6 New software, platforms, open data</b>	<b>11</b>
6.1 New software	11
6.1.1 ATMI	11
6.1.2 HEPTANE	12
6.1.3 tiptop	12
6.1.4 GATO3D	13
6.1.5 OptiPrint	13
6.1.6 SAMVA	13
6.1.7 TimeKlip	13
6.2 New platforms	14
6.3 Open data	14
<b>7 New results</b>	<b>14</b>
7.1 Compilation and Optimization	15
7.1.1 Compilation for Intermittent Systems	15
7.1.2 Dynamic Binary Analysis and Optimization	15
7.1.3 Accurate 3D printing time estimation	15
7.2 Processor Architecture	16
7.2.1 Automatic synthesis of multi-thread pipelines	16
7.2.2 Two-dimensional memory architecture	16
7.3 WCET estimation and optimization	17
7.3.1 Using machine learning for timing analysis of complex processors	17
7.3.2 Static estimation of memory access profiles for real-time multi-core systems	18
7.4 Security	18
7.4.1 Verification of Data Flow Integrity for Real-Time Embedded Systems	18
7.4.2 Multi-nop fault injection attack	18
7.4.3 Gadget chains synthesis driven by SMT Solving for Code-Reuse Attacks	19

<b>8</b>	<b>Bilateral contracts and grants with industry</b>	<b>19</b>
8.1	Bilateral contracts with industry . . . . .	19
<b>9</b>	<b>Partnerships and cooperations</b>	<b>20</b>
9.1	European initiatives . . . . .	20
9.1.1	H2020 projects . . . . .	20
9.2	National initiatives . . . . .	21
<b>10</b>	<b>Dissemination</b>	<b>26</b>
10.1	Promoting scientific activities . . . . .	27
10.1.1	Scientific events: selection . . . . .	27
10.1.2	Journal . . . . .	27
10.1.3	Invited talks . . . . .	27
10.1.4	Leadership within the scientific community . . . . .	27
10.1.5	Research administration . . . . .	27
10.2	Teaching - Supervision - Juries . . . . .	28
10.2.1	Teaching . . . . .	28
10.2.2	Supervision . . . . .	28
10.2.3	Juries . . . . .	29
10.3	Popularization . . . . .	30
10.3.1	Internal or external Inria responsibilities . . . . .	30
10.3.2	Articles and contents . . . . .	30
10.3.3	Education . . . . .	30
10.3.4	Interventions . . . . .	30
<b>11</b>	<b>Scientific production</b>	<b>30</b>
11.1	Major publications . . . . .	30
11.2	Publications of the year . . . . .	31
11.3	Cited publications . . . . .	32

## Project-Team PACAP

*Creation of the Project-Team: 2016 July 01*

### Keywords

#### Computer sciences and digital sciences

- A1.1.1. – Multicore, Manycore
- A1.1.2. – Hardware accelerators (GPGPU, FPGA, etc.)
- A1.1.3. – Memory models
- A1.1.8. – Security of architectures
- A1.1.11. – Quantum architectures
- A1.6. – Green Computing
- A2.2.1. – Static analysis
- A2.2.3. – Memory management
- A2.2.4. – Parallel architectures
- A2.2.5. – Run-time systems
- A2.2.6. – GPGPU, FPGA...
- A2.2.7. – Adaptive compilation
- A2.2.8. – Code generation
- A2.2.9. – Security by compilation
- A2.3. – Embedded and cyber-physical systems
- A2.3.1. – Embedded systems
- A2.3.2. – Cyber-physical systems
- A2.3.3. – Real-time systems
- A4.4. – Security of equipment and software
- A5.10.3. – Planning
- A5.10.5. – Robot interaction (with the environment, humans, other robots)
- A9.2. – Machine learning

#### Other research topics and application domains

- B1. – Life sciences
- B2. – Health
- B3. – Environment and planet
- B4. – Energy
- B5. – Industry of the future
- B5.7. – 3D printing
- B6. – IT and telecom
- B7. – Transport and logistics
- B8. – Smart Cities and Territories
- B9. – Society and Knowledge

# 1 Team members, visitors, external collaborators

## Research Scientists

- Erven Rohou [Team leader, INRIA, Senior Researcher, HDR]
- Caroline Collange [INRIA, Researcher]
- Pierre Michaud [INRIA, Researcher]
- Thomas Rubiano [INRIA, Starting Research Position, from Nov 2023]

## Faculty Members

- Damien Hardy [UNIV RENNES, Associate Professor]
- Isabelle Puaut [UNIV RENNES, Professor, HDR]

## PhD Students

- Abderaouf Nassim Amalou [UNIV RENNES]
- Nicolas Bailluet [UNIV RENNES]
- Hector Chabot [UNIV RENNES, from Sep 2023]
- Antoine Gicquel [INRIA]
- Sara Sadat Hoseininasab [INRIA]
- Anis Peysieux [INRIA, until Mar 2023]
- Aurore Poirier [INRIA]
- Hugo Reymond [INRIA]

## Technical Staff

- Pierre Bedell [INRIA, Engineer]
- Camille Le Bon [INRIA, Engineer]
- Mohammed Mehdi Merah [UNIV RENNES, until Feb 2023]

## Interns and Apprentices

- Valentin Buffa [INRIA, Intern, from May 2023 until Jul 2023]
- Killian Callac [INRIA, Intern, from May 2023 until Jul 2023]
- Hector Chabot [UNIV RENNES, Intern, until Jul 2023]
- Ewen Coquio [UNIV RENNES, Intern, from May 2023 until Aug 2023]
- Lucas Le Guedes [UNIV RENNES, Intern, from May 2023 until Aug 2023]
- Antoine Quere [INRIA, Intern, from May 2023 until Aug 2023]
- Charlotte Thomas [UNIV RENNES, Intern, from Feb 2023 until Jul 2023]

## Administrative Assistant

- Virginie Desroches [INRIA]

## 2 Overall objectives

**Long-Term Goal** In brief, the long-term goal of the PACAP project-team is about *performance*, that is: how fast programs run. We intend to contribute to the ongoing race for exponentially increasing performance and for performance guarantees.

Traditionally, the term “performance” is understood as “how much time is needed to complete execution”. *Latency*-oriented techniques focus on minimizing the average-case execution time (ACET). We are also interested in other definitions of performance. *Throughput*-oriented techniques are concerned with how many units of computation can be completed per unit of time. This is more relevant on manycores and GPUs where many computing nodes are available, and latency is less critical. Finally, we also study worst-case execution time (WCET), which is extremely important for critical real-time systems where designers must guarantee that deadlines are met, in any situation.

Given the complexity of current systems, simply assessing their performance has become a non-trivial task which we also plan to tackle.

We occasionally consider other metrics related to performance, such as power efficiency, total energy, overall complexity, and real-time response guarantee. Our ultimate goal is to propose solutions that make computing systems more efficient, taking into account current and envisioned applications, compilers, runtimes, operating systems, and micro-architectures. And since increased performance often comes at the expense of another metric, identifying the related trade-offs is of interest to PACAP.

The previous decade witnessed the end of the “magically” increasing clock frequency and the introduction of commodity multicore processors. PACAP is experiencing the end of Moore’s law<sup>1</sup>, and the generalization of commodity heterogeneous manycore processors. This impacts how performance is increased and how it can be guaranteed. It is also a time where exogenous parameters should be promoted to first-class citizens:

1. the existence of faults, whose impact is becoming increasingly important when the photo-lithography feature size decreases;
2. the need for security at all levels of computing systems;
3. *green* computing, or the growing concern of power consumption.

**Approach** We strive to address performance in a way that is as transparent as possible to the users. For example, instead of proposing any new language, we consider existing applications (written for example in standard C), and we develop compiler optimizations that immediately benefit programmers; we propose microarchitectural features as opposed to changes in processor instruction sets; we analyze and re-optimize binary programs automatically, without any user intervention.

The perimeter of research directions of the PACAP project-team derives from the intersection of two axes: on the one hand, our high-level research objectives, derived from the overall panorama of computing systems, on the other hand the existing expertise and background of the team members in key technologies (see illustration on Figure 1). Note that it does not imply that we will systematically explore all intersecting points of the figure, yet all correspond to a sensible research direction. These lists are neither exhaustive, nor final. Operating systems in particular constitute a promising operating point for several of the issues we plan to tackle. Other aspects will likely emerge during the lifespan of the project-team.

**Latency-oriented Computing** Improving the ACET of general purpose systems has been the “core business” of PACAP’s ancestors (CAPS and ALF) for two decades. We plan to pursue this line of research, acting at all levels: compilation, dynamic optimizations, and micro-architecture.

**Throughput-Oriented Computing** The goal is to maximize the performance-to-power ratio. We will leverage the execution model of throughput-oriented architectures (such as GPUs) and extend it towards general purpose systems. To address the memory wall issue, we will consider bandwidth saving techniques, such as cache and memory compression.

---

<sup>1</sup>Moore’s law states that the number of transistors in a circuit doubles (approximately) every two years.

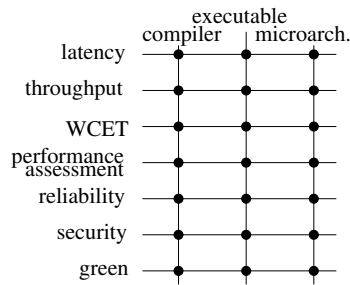


Figure 1: Perimeter of Research Objectives

**Real-Time Systems – WCET** Designers of real-time systems must provide an upper bound of the worst-case execution time of the tasks within their systems. By definition this bound must be safe (i.e., greater than any possible execution time). To be useful, WCET estimates have to be as tight as possible. The process of obtaining a WCET bound consists in analyzing a binary executable, modeling the hardware, and then maximizing an objective function that takes into account all possible flows of execution and their respective execution times. Our research will consider the following directions:

1. better modeling of hardware to either improve tightness, or handle more complex hardware (e.g. multicores);
2. eliminate unfeasible paths from the analysis;
3. consider probabilistic approaches where WCET estimates are provided with a confidence level.

**Performance Assessment** Moore's law drives the complexity of processor micro-architectures, which impacts all other layers: hypervisors, operating systems, compilers and applications follow similar trends. While a small category of experts is able to comprehend (parts of) the behavior of the system, the vast majority of users are only exposed to – and interested in – the bottom line: how fast their applications are actually running. In the presence of virtual machines and cloud computing, multi-programmed workloads add yet another degree of non-determinism to the measure of performance. We plan to research how application performance can be characterized and presented to a final user: behavior of the micro-architecture, relevant metrics, possibly visual rendering. Targeting our own community, we also research techniques appropriate for fast and accurate ways to simulate future architectures, including heterogeneous designs, such as latency/throughput platforms.

Once diagnosed, the way bottlenecks are addressed depends on the level of expertise of users. Experts can typically be left with a diagnostic as they probably know better how to fix the issue. Less knowledgeable users must be guided to a better solution. We plan to rely on iterative compilation to generate multiple versions of critical code regions, to be used in various runtime conditions. To avoid the code bloat resulting from multiversioning, we will leverage split-compilation to embed code generation “recipes” to be applied just-in-time, or even at runtime thanks to dynamic binary translation. Finally, we will explore the applicability of auto-tuning, where programmers expose which parameters of their code can be modified to generate alternate versions of the program (for example trading energy consumption for quality of service) and let a global orchestrator make decisions.

**Dealing with Attacks – Security** Computer systems are under constant attack, from young hackers trying to show their skills, to “professional” criminals stealing credit card information, and even government agencies with virtually unlimited resources. A vast amount of techniques have been proposed in the literature to circumvent attacks. Many of them cause significant slowdowns due to additional checks and countermeasures. Thanks to our expertise in micro-architecture and compilation techniques, we will be able to significantly improve efficiency, robustness and coverage of security mechanisms, as well as to partner with field experts to design innovative solutions.

**Green Computing – Power Concerns** Power consumption has become a major concern of computing systems, at all form factors, ranging from energy-scavenging sensors for IoT, to battery powered embedded systems and laptops, and up to supercomputers operating in the tens of megawatts. Execution time and energy are often related optimization goals. Optimizing for performance under a given power cap, however, introduces new challenges. It also turns out that technologists introduce new solutions (e.g. magnetic RAM) which, in turn, result in new trade-offs and optimization opportunities.

## 3 Research program

### 3.1 Motivation

Our research program is naturally driven by the evolution of our ecosystem. Relevant recent changes can be classified in the following categories: technological constraints, evolving community, and domain constraints. We hereby summarize these evolutions.

#### 3.1.1 Technological constraints

Until recently, binary compatibility guaranteed portability of programs, while increased clock frequency and improved micro-architecture provided increased performance. However, in the last decade, advances in technology and micro-architecture started translating into more parallelism instead. Technology roadmaps even predicted the feasibility of thousands of cores on a chip by the 2020's. Hundreds are already commercially available. Since the vast majority of applications are still sequential, or contain significant sequential sections, such a trend puts an end to the automatic performance improvement enjoyed by developers and users. Many research groups consequently focused on parallel architectures and compiling for parallelism.

Still, the performance of applications will ultimately be driven by the performance of the sequential part. Despite a number of advances (some of them contributed by members of the team), sequential tasks are still a major performance bottleneck. Addressing it is still on the agenda of the PACAP project-team.

In addition, due to power constraints, only part of the billions of transistors of a microprocessor can be operated at any given time (the *dark silicon* paradigm). A sensible approach consists in specializing parts of the silicon area to provide dedicated accelerators (not run simultaneously). This results in diverse and heterogeneous processor cores. Application and compiler designers are thus confronted with a moving target, challenging portability and jeopardizing performance.

*Note on technology.*

Technology also progresses at a fast pace. We do not propose to pursue any research on technology *per se*. Recently proposed paradigms (non-Silicon, brain-inspired) have received lots of attention from the research community. We do *not* intend to invest in those paradigms, but we will continue to investigate compilation and architecture for more conventional programming paradigms. Still, several technological shifts may have consequences for us, and we will closely monitor their developments. They include for example non-volatile memory (impacts security, makes writes longer than loads), 3D-stacking (impacts bandwidth), and photonics (impacts latencies and connection network), quantum computing (impacts the entire software stack).

#### 3.1.2 Evolving community

The PACAP project-team tackles performance-related issues, for conventional programming paradigms. In fact, programming complex environments is no longer the exclusive domain of experts in compilation and architecture. A large community now develops applications for a wide range of targets, including mobile “apps”, cloud, multicore or heterogeneous processors.

This also includes domain scientists (in biology, medicine, but also social sciences) who started relying heavily on computational resources, gathering huge amounts of data, and requiring a considerable amount of processing to analyze them. Our research is motivated by the growing discrepancy between on the one hand, the complexity of the workloads and the computing systems, and on the other hand, the expanding community of developers at large, with limited expertise to optimize and to map efficiently computations to compute nodes.



### 3.1.3 Domain constraints

Mobile, embedded systems have become ubiquitous. Many of them have real-time constraints. For this class of systems, correctness implies not only producing the correct result, but also doing so within specified deadlines. In the presence of heterogeneous, complex and highly dynamic systems, producing a *tight* (i.e., useful) upper bound to the worst-case execution time has become extremely challenging. Our research will aim at improving the tightness as well as enlarging the set of features that can be safely analyzed.

The ever growing dependence of our economy on computing systems also implies that security has become of utmost importance. Many systems are under constant attacks from intruders. Protection has a cost also in terms of performance. We plan to leverage our background to contribute solutions that minimize this impact.

*Note on Applications Domains.*

PACAP works on fundamental technologies for computer science: processor architecture, performance-oriented compilation and guaranteed response time for real-time. The research results may have impact on any application domain that requires high performance execution (telecommunication, multimedia, biology, health, engineering, environment...), but also on many embedded applications that exhibit other constraints such as power consumption, code size and guaranteed response time.

We strive to extract from active domains the fundamental characteristics that are relevant to our research. For example, *big data* is of interest to PACAP because it relates to the study of hardware/software mechanisms to efficiently transfer huge amounts of data to the computing nodes. Similarly, the *Internet of Things* is of interest because it has implications in terms of ultra low-power consumption.

## 3.2 Research Objectives

Processor micro-architecture and compilation have been at the core of the research carried by the members of the project teams for two decades, with undeniable contributions. They continue to be the foundation of PACAP.

Heterogeneity and diversity of processor architectures now require new techniques to guarantee that the hardware is satisfactorily exploited by the software. One of our goals is to devise new static compilation techniques (cf. Section 3.2.1), but also build upon iterative [1] and split [27] compilation to continuously adapt software to its environment (Section 3.2.2). Dynamic binary optimization will also play a key role in delivering adapting software and increased performance.

The end of Moore's law and Dennard's scaling<sup>2</sup> offer an exciting window of opportunity, where performance improvements will no longer derive from additional transistor budget or increased clock frequency, but rather come from breakthroughs in micro-architecture (Section 3.2.3). Reconciling CPU and GPU designs (Section 3.2.4) is one of our objectives.

Heterogeneity and multicores are also major obstacles to determining tight worst-case execution times of real-time systems (Section 3.2.5), which we plan to tackle.

Finally, we also describe how we plan to address transversal aspects such as power efficiency (Section 3.2.6), and security (Section 3.2.7).

### 3.2.1 Static Compilation

Static compilation techniques continue to be relevant in addressing the characteristics of emerging hardware technologies, such as non-volatile memories, 3D-stacking, or novel communication technologies. These techniques expose new characteristics to the software layers. As an example, non-volatile memories typically have asymmetric read-write latencies (writes are much longer than reads) and different power consumption profiles. PACAP studies new optimization opportunities and develops tailored compilation techniques for upcoming compute nodes. New technologies may also be coupled with traditional solutions to offer new trade-offs. We study how programs can adequately exploit the specific features of the proposed heterogeneous compute nodes.

<sup>2</sup>According to Dennard scaling, as transistors get smaller the power density remains constant, and the consumed power remains proportional to the area.

We propose to build upon iterative compilation [1] to explore how applications perform on different configurations. When possible, Pareto points are related to application characteristics. The best configuration, however, may actually depend on runtime information, such as input data, dynamic events, or properties that are available only at runtime. Unfortunately a runtime system has little time and means to determine the best configuration. For these reasons, we also leverage split-compilation [27]: the idea consists in pre-computing alternatives, and embedding in the program enough information to assist and drive a runtime system towards to the best solution.

### 3.2.2 Software Adaptation

More than ever, software needs to adapt to its environment. In most cases, this environment remains unknown until runtime. This is already the case when one deploys an application to a cloud, or an “app” to mobile devices. The dilemma is the following: for maximum portability, developers should target the most general device; but for performance they would like to exploit the most recent and advanced hardware features. JIT compilers can handle the situation to some extent, but binary deployment requires dynamic binary rewriting. Our work has shown how SIMD instructions can be upgraded from SSE to AVX transparently [2]. Many more opportunities will appear with diverse and heterogeneous processors, featuring various kinds of accelerators.

On shared hardware, the environment is also defined by other applications competing for the same computational resources. It becomes increasingly important to adapt to changing runtime conditions, such as the contention of the cache memories, available bandwidth, or hardware faults. Fortunately, optimizing at runtime is also an opportunity, because this is the first time the program is visible as a whole: executable and libraries (including library versions). Optimizers may also rely on dynamic information, such as actual input data, parameter values, etc. We have already developed software platforms [37, 33] to analyze and optimize programs at runtime, and we started working on automatic dynamic parallelization of sequential code, and dynamic specialization.

We addressed some of these challenges in previous projects such as Nano2017 PSAIC Collaborative research program with STMicroelectronics, as well as within the Inria Project Lab MULTICORE. The H2020 FET HPC project ANTAREX also addressed these challenges from the energy perspective, while the ANR Continuum project and the Inria Challenge ZEP focused on opportunities brought by non-volatile memories. We further leverage our platform and initial results to address other adaptation opportunities. Efficient software adaptation requires expertise from all domains tackled by PACAP, and strong interaction between all team members is expected.

### 3.2.3 Research directions in uniprocessor micro-architecture

Achieving high single-thread performance remains a major challenge even in the multicore era (Amdahl’s law). The members of the PACAP project-team have been conducting research in uniprocessor micro-architecture research for about 25 years covering major topics including caches, instruction front-end, branch prediction, out-of-order core pipeline, and value prediction. In particular, in recent years they have been recognized as world leaders in branch prediction [39] [34] and in cache prefetching [5] and they have revived the forgotten concept of value prediction [8][7]. This research was supported by the ERC Advanced grant DAL (2011-2016) and also by Intel. We pursue research on achieving ultimate uncore performance. Below are several non-orthogonal directions that we have identified for mid-term research:

1. management of the memory hierarchy (particularly the hardware prefetching);
2. practical design of very wide issue execution cores;
3. speculative execution.

#### *Memory design issues:*

Performance of many applications is highly impacted by the memory hierarchy behavior. The interactions between the different components in the memory hierarchy and the out-of-order execution engine have high impact on performance.

The *Data Prefetching Contest* held with ISCA 2015 has illustrated that achieving high prefetching efficiency is still a challenge for wide-issue superscalar processors, particularly those featuring a very large

instruction window. The large instruction window enables an implicit data prefetcher. The interaction between this implicit hardware prefetcher and the explicit hardware prefetcher is still relatively mysterious as illustrated by Pierre Michaud's BO prefetcher (winner of DPC2) [5]. The first research objective is to better understand how the implicit prefetching enabled by the large instruction window interacts with the L2 prefetcher and then to understand how explicit prefetching on the L1 also interacts with the L2 prefetcher.

The second research objective is related to the interaction of prefetching and virtual/physical memory. On real hardware, prefetching is stopped by page frontiers. The interaction between TLB prefetching (and on which level) and cache prefetching must be analyzed.

The prefetcher is not the only actor in the hierarchy that must be carefully controlled. Significant benefits can also be achieved through careful management of memory access bandwidth, particularly the management of spatial locality on memory accesses, both for reads and writes. The exploitation of this locality is traditionally handled in the memory controller. However, it could be better handled if larger temporal granularity was available. Finally, we also intend to continue to explore the promising avenue of compressed caches. In particular we proposed the skewed compressed cache [11]. It offers new possibilities for efficient compression schemes.

#### *Ultra wide-issue superscalar.*

To effectively leverage memory level parallelism, one requires huge out-of-order execution structures as well as very wide issue superscalar processors. For the two past decades, implementing ever wider issue superscalar processors has been challenging. The objective of our research on the execution core is to explore (and revisit) directions that allow the design of a very wide-issue (8-to-16 way) out-of-order execution core while mastering its complexity (silicon area, hardware logic complexity, power/energy consumption).

The first direction that we are exploring is the use of clustered architectures [6]. Symmetric clustered organization allows to benefit from a simpler bypass network, but induce large complexity on the issue queue. One remarkable finding of our study [6] is that, when considering two large clusters (e.g. 8-wide), steering large groups of consecutive instructions (e.g. 64  $\mu$ ops) to the same cluster is quite efficient. This opens opportunities to limit the complexity of the issue queues (monitoring fewer buses) and register files (fewer ports and physical registers) in the clusters, since not all results have to be forwarded to the other cluster.

The second direction that we are exploring is associated with the approach that we developed with Sembrant et al. [38]. It reduces the number of instructions waiting in the instruction queues for the applications benefiting from very large instruction windows. Instructions are dynamically classified as ready (independent from any long latency instruction) or non-ready, and as urgent (part of a dependency chain leading to a long latency instruction) or non-urgent. Non-ready non-urgent instructions can be delayed until the long latency instruction has been executed; this allows to reduce the pressure on the issue queue. This proposition opens the opportunity to consider an asymmetric micro-architecture with a cluster dedicated to the execution of urgent instructions and a second cluster executing the non-urgent instructions. The micro-architecture of this second cluster could be optimized to reduce complexity and power consumption (smaller instruction queue, less aggressive scheduling...)

#### *Speculative execution.*

Out-of-order (OoO) execution relies on speculative execution that requires predictions of all sorts: branch, memory dependency, value...

The PACAP members have been major actors of branch prediction research for the last 25 years; and their proposals have influenced the design of most of the hardware branch predictors in current microprocessors. We will continue to steadily explore new branch predictor designs, as for instance [40].

In speculative execution, we have recently revisited value prediction (VP) which was a hot research topic between 1996 and 2002. However it was considered until recently that value prediction would lead to a huge increase in complexity and power consumption in every stage of the pipeline. Fortunately, we have recently shown that complexity usually introduced by value prediction in the OoO engine can be overcome [8][7] [39] [34]. First, very high accuracy can be enforced at reasonable cost in coverage and minimal complexity [8]. Thus, both prediction validation and recovery by squashing can be done outside the out-of-order engine, at commit time. Furthermore, we propose a new pipeline organization, EOLE ({Early | Out-of-order | Late} Execution), that leverages VP with validation at commit to execute many instructions outside the OoO core, in-order [7]. With EOLE, the issue-width in OoO core can be reduced

without sacrificing performance, thus benefiting the performance of VP without a significant cost in silicon area and/or energy. In the near future, we will explore new avenues related to value prediction. These directions include register equality prediction and compatibility of value prediction with weak memory models in multiprocessors.

### 3.2.4 Towards heterogeneous single-ISA CPU-GPU architectures

Heterogeneous single-ISA architectures have been proposed in the literature during the 2000's [32] and are now widely used in the industry (Arm big.LITTLE, NVIDIA 4+1, Intel Alder Lake...) as a way to improve power-efficiency in mobile processors. These architectures include multiple cores whose respective micro-architectures offer different trade-offs between performance and energy efficiency, or between latency and throughput, while offering the same interface to software. Dynamic task migration policies leverage the heterogeneity of the platform by using the most suitable core for each application, or even each phase of processing. However, these works only tune cores by changing their complexity. Energy-optimized cores are either identical cores implemented in a low-power process technology, or simplified in-order superscalar cores, which are far from state-of-the-art throughput-oriented architectures such as GPUs.

We investigate the convergence of CPU and GPU at both architecture and compiler levels.

#### *Architecture.*

The architecture convergence between Single Instruction Multiple Threads (SIMT) GPUs and multicore processors that we have been pursuing [15] opens the way for heterogeneous architectures including latency-optimized superscalar cores and throughput-optimized GPU-style cores, which all share the same instruction set. Using SIMT cores in place of superscalar cores will enable the highest energy efficiency on regular sections of applications. As with existing single-ISA heterogeneous architectures, task migration will not necessitate any software rewrite and will accelerate existing applications.

#### *Compilers for emerging heterogeneous architectures.*

Single-ISA CPU+GPU architectures will provide the necessary substrate to enable efficient heterogeneous processing. However, it will also introduce substantial challenges at the software and firmware level. Task placement and migration will require advanced policies that leverage both static information at compile time and dynamic information at run-time. We are tackling the heterogeneous task scheduling problem at the compiler level.

### 3.2.5 Real-time systems

Safety-critical systems (e.g. avionics, medical devices, automotive...) have so far used simple uncore hardware systems as a way to control their predictability, in order to meet timing constraints. Still, many critical embedded systems have increasing demand in computing power, and simple uncore processors are not sufficient anymore. General-purpose multicore processors are not suitable for safety-critical real-time systems, because they include complex micro-architectural elements (cache hierarchies, branch, stride and value predictors) meant to improve average-case performance, and for which worst-case performance is difficult to predict. The prerequisite for calculating tight WCET is a deterministic hardware system that avoids dynamic, time-unpredictable calculations at run-time.

Even for multi and manycore systems designed with time-predictability in mind (Kalray MPPA manycore architecture or the Recore manycore hardware) calculating WCETs is still challenging. The following two challenges will be addressed in the mid-term:

1. definition of methods to estimate WCETs tightly on manycores, that smartly analyze and/or control shared resources such as buses, NoCs or caches;
2. methods to improve the programmability of real-time applications through automatic parallelization and optimizations from model-based designs.

### 3.2.6 Power efficiency

PACAP addresses power-efficiency at several levels. First, we design static and split compilation techniques to contribute to the race for Exascale computing (the general goal is to reach  $10^{18}$  FLOP/s at less

than 20 MW). Second, we focus on high-performance low-power embedded compute nodes. Within the ANR project Continuum, in collaboration with architecture and technology experts from LIRMM and the SME Cortus, we researched new static and dynamic compilation techniques that fully exploit emerging memory and NoC technologies. Finally, in collaboration with the TARAN project-team, we investigate the synergy of reconfigurable computing and dynamic code generation.

*Green and heterogeneous high-performance computing.*

Concerning HPC systems, our approach consists in mapping, runtime managing and autotuning applications for green and heterogeneous High-Performance Computing systems up to the Exascale level. One key innovation of the proposed approach consists in introducing a separation of concerns (where self-adaptivity and energy efficient strategies are specified aside to application functionalities) promoted by the definition of a Domain Specific Language (DSL) inspired by aspect-oriented programming concepts for heterogeneous systems. The new DSL will be introduced for expressing adaptivity/energy/performance strategies and to enforce at runtime application autotuning and resource and power management. The goal is to support the parallelism, scalability and adaptability of a dynamic workload by exploiting the full system capabilities (including energy management) for emerging large-scale and extreme-scale systems, while reducing the Total Cost of Ownership (TCO) for companies and public organizations.

*High-performance low-power embedded compute nodes.*

We will address the design of next generation energy-efficient high-performance embedded compute nodes. We focus at the same time on software, architecture and emerging memory and communication technologies in order to synergistically exploit their corresponding features. The approach of the project is organized around three complementary topics: 1) compilation techniques; 2) multicore architectures; 3) emerging memory and communication technologies. PACAP will focus on the compilation aspects, taking as input the software-visible characteristics of the proposed emerging technology, and making the best possible use of the new features (non-volatility, density, endurance, low-power).

*Hardware Accelerated JIT Compilation.*

Reconfigurable hardware offers the opportunity to limit power consumption by dynamically adjusting the number of available resources to the requirements of the running software. In particular, VLIW processors can adjust the number of available issue lanes. Unfortunately, changing the processor width often requires recompiling the application, and VLIW processors are highly dependent of the quality of the compilation, mainly because of the instruction scheduling phase performed by the compiler. Another challenge lies in the high constraints of the embedded system: the energy and execution time overhead due to the JIT compilation must be carefully kept under control.

We started exploring ways to reduce the cost of JIT compilation targeting VLIW-based heterogeneous manycore systems. Our approach relies on a hardware/software JIT compiler framework. While basic optimizations and JIT management are performed in software, the compilation back-end is implemented by means of specialized hardware. This back-end involves both instruction scheduling and register allocation, which are known to be the most time-consuming stages of such a compiler.

### 3.2.7 Security

Security is a mandatory concern of any modern computing system. Various threat models have led to a multitude of protection solutions. Members of PACAP already contributed in the past, thanks to the HAVEGE [41] random number generator, and code obfuscating techniques (the obfuscating just-in-time compiler [30], or thread-based control flow mangling [35]). Still, security is not core competence of PACAP members.

Our strategy consists in partnering with security experts who can provide intuition, know-how and expertise, in particular in defining threat models, and assessing the quality of the solutions. Our expertise in compilation and architecture helps design more efficient and less expensive protection mechanisms.

Examples of collaborations so far include the following:

**Compilation:** We partnered with experts in security and codes to prototype a platform that demonstrates resilient software. They designed and proposed advanced masking techniques to hide sensitive data in application memory. PACAP's expertise is key to select and tune the protection mechanisms developed within the project, and to propose safe, yet cost-effective solutions from an implementation point of view.

**Dynamic Binary Rewriting:** Our expertise in dynamic binary rewriting combines well with the expertise of the CIDRE team in protecting application. Security has a high cost in terms of performance, and static insertion of countermeasures cannot take into account the current threat level. In collaboration with CIDRE, we proposed an adaptive insertion/removal of countermeasures in a running application based of dynamic assessment of the threat level.

**WCET Analysis:** Designing real-time systems requires computing an upper bound of the worst-case execution time. Knowledge of this timing information opens an opportunity to detect attacks on the control flow of programs. In collaboration with CIDRE, we developed a technique to detect such attacks thanks to a hardware monitor that makes sure that statically computed time information is preserved (TARAN is also involved in the definition of the hardware component).

## 4 Application domains

### 4.1 Domains

The PACAP team is working on fundamental technologies for computer science: processor architecture, performance-oriented compilation and guaranteed response time for real-time. The research results may have impact on any application domain that requires high performance execution (telecommunication, multimedia, biology, health, engineering, environment...), but also on many embedded applications that exhibit other constraints such as power consumption, code size and guaranteed response time. Our research activity implies the development of software prototypes.

## 5 Social and environmental responsibility

### 5.1 Impact of research results

For a few years now, the PACAP team has been contributing to the transition from traditional IoT networks to battery-less networks. The increasing number of IoT devices led to a proliferation of batteries in the environment, associated with their well-known ecological and social footprint.

In a effort to reduce this footprint, PACAP provides compiler building blocks to support intermittent computing, i.e. the execution of programs on battery-less devices, powered by energy harvesting. This supports allow the devices to endure frequent power failures.

This work has been presented and discussed in events on sustainable development [21] and main conference [20].

## 6 New software, platforms, open data

### 6.1 New software

#### 6.1.1 ATMI

**Keywords:** Analytic model, Chip design, Temperature

**Scientific Description:** Research on temperature-aware computer architecture requires a chip temperature model. General-purpose models based on classical numerical methods like finite differences or finite elements are not appropriate for such research, because they are generally too slow for modeling the time-varying thermal behavior of a processing chip.

ATMI (Analytical model of Temperature in MICroprocessors) is an ad hoc temperature model for studying thermal behaviors over a time scale ranging from microseconds to several minutes. ATMI is based on an explicit solution to the heat equation and on the principle of superposition. ATMI can model any power density map that can be described as a superposition of rectangle sources, which is appropriate for modeling the microarchitectural units of a microprocessor.



**Functional Description:** ATMI is a library for modelling steady-state and time-varying temperature in microprocessors. ATMI uses a simplified representation of microprocessor packaging.

**URL:** <https://team.inria.fr/pacap/software/atmi/>

**Contact:** Pierre Michaud

**Participant:** Pierre Michaud

### 6.1.2 HEPTANE

**Keywords:** IPET, WCET, Performance, Real time, Static analysis, Worst Case Execution Time

**Scientific Description:** WCET estimation

The aim of Heptane is to produce upper bounds of the execution times of applications. It is targeted at applications with hard real-time requirements (automotive, railway, aerospace domains). Heptane computes WCETs using static analysis at the binary code level. It includes static analyses of microarchitectural elements such as caches and cache hierarchies.

**Functional Description:** In a hard real-time system, it is essential to comply with timing constraints, and Worst Case Execution Time (WCET) in particular. Timing analysis is performed at two levels: analysis of the WCET for each task in isolation taking account of the hardware architecture, and schedulability analysis of all the tasks in the system. Heptane is a static WCET analyser designed to address the first issue.

**URL:** <https://team.inria.fr/pacap/software/heptane/>

**Contact:** Isabelle Puaut

**Participants:** Benjamin Lesage, Loïc Besnard, Damien Hardy, François Joulaud, Isabelle Puaut, Thomas Piquet

**Partner:** Université de Rennes 1

### 6.1.3 tiptop

**Keywords:** Instructions, Cycles, Cache, CPU, Performance, HPC, Branch predictor

**Scientific Description:** Tiptop is a simple and flexible user-level tool that collects hardware counter data on Linux platforms (version 2.6.31+) and displays them in a way simple to the Linux "top" utility. The goal is to make the collection of performance and bottleneck data as simple as possible, including simple installation and usage. Unless the system administrator has restricted access to performance counters, no privilege is required, any user can run tiptop.

Tiptop is written in C. It can take advantage of libncurses when available for pseudo-graphic display. Installation is only a matter of compiling the source code. No patching of the Linux kernel is needed, and no special-purpose module needs to be loaded.

Current version is 2.3.2, released December 2023. Tiptop has been integrated in major Linux distributions, such as Fedora, Debian, Ubuntu, CentOS.

**Functional Description:** Today's microprocessors have become extremely complex. To better understand the multitude of internal events, manufacturers have integrated many monitoring counters. Tiptop can be used to collect and display the values from these performance counters very easily. Tiptop may be of interest to anyone who wants to optimize the performance of their HPC applications.

**URL:** <https://team.inria.fr/pacap/software/tiptop/>

**Contact:** Erven Rohou

**Participant:** Erven Rohou

#### 6.1.4 GATO3D

**Keywords:** Code optimisation, 3D printing

**Functional Description:** GATO3D stands for "G-code Analysis Transformation and Optimization". It is a library that provides an abstraction of the G-code, the language interpreted by 3D printers, as well as an API to manipulate it easily. First, GATO3D reads a file in G-code format and builds its representation in memory. This representation can be transcribed into a G-code file at the end of the manipulation. The software also contains client codes for the computation of G-code properties, the optimization of displacements, and a graphical rendering.

**Authors:** Damien Hardy, Erven Rohou

**Contact:** Erven Rohou

#### 6.1.5 OptiPrint

**Keywords:** 3D printing, Planning, Optimization

**Functional Description:** OptiPrint is a software library dedicated to print time optimization for fused filament deposition (FDM) printers. This library is integrated to the Gato3D compiler. Its role is to allow the optimization of the printing time by reordering / filtering the G-code sent to a 3D printer. The optimization is fully configurable. It adapts to the characteristics of the printers (type of nozzle, speed of movement of the nozzle). It also allows to describe scheduling constraints allowing to make a compromise between printing quality and optimization.

**Contact:** Fabrice Lamarche

#### 6.1.6 SAMVA

**Keywords:** Static analysis, Fault injection

**Functional Description:** SAMVA is a software package for determining attack paths in the context of precise, multiple fault injection attacks. It is framework for efficiently searching vulnerabilities of applications in presence of multiple instruction-skip faults with various widths. SAMVA relies solely on static analysis to determine attack paths in a binary code. It is configurable with the fault injection capacity of the attacker and the attacker's objective

**Authors:** Antoine Gicquel, Erven Rohou, Karine Heydemann, Damien Hardy

**Contact:** Erven Rohou

#### 6.1.7 TimeKlip

**Keywords:** Simulator, 3D printing

**Functional Description:** 3D printing simulator calculating the printing time of a G-code file. It is able to give timing information for each instruction in the file. The simulator does not require a printer to run, only configuration files. It is also slicer agnostic.

The simulator takes the form of a module integrated into the Klipper firmware.

**Authors:** Damien Hardy, Camille Le Bon

**Contact:** Damien Hardy



## 6.2 New platforms

**Participants:** Pierre Bedell, Damien Hardy.

In the context of the Inria exploratory action Ofast3D, a 3D printing platform for research experiments is under construction. At this stage, it is composed of 11 printers and 4 test benches. This allows to evaluate optimizations and time prediction on different kinematics and configurations as well as different firmwares. Furthermore, air quality sensors are under deployment to evaluate the impact of 3D printing materials.

This platform is used by other teams in particular: MimeTIC, Rainbow, LACODAM, TARAN and LogicA.

## 6.3 Open data

**Participants:** Abderaouf Nassim Amalou, Isabelle Puaut.

**Pre-training and fine-tuning dataset for transformers consisting of basic blocks and their execution times (average, minimum, and maximum) along with the execution context of these blocks, for various Cortex processors M7, M4, A53, and A72** We are releasing the dataset used for training CAWET [16], a tool for estimating the Worst-Case Execution Time (WCET) of basic blocks using the Transformer XL model. CAWET leverages the Transformer architecture for accurate WCET predictions, and its training involves two main phases: self-supervised pre-training and fine-tuning.

CAWET undergoes a pre-training process on a substantial corpus of basic blocks to enable the Transformer to grasp the intricacies of the assembly language in focus. For this, we utilized CodeNet [36], a comprehensive collection of publicly submitted solutions to competitive programming challenges, comprising roughly 900,000 C programs. These programs were cross-compiled to the target architecture and subsequently disassembled using GNU binary utilities with objdump. The textual output from objdump, after a series of basic parsing operations (e.g., address extraction, separation of basic blocks), serves as the foundation for an extensive pre-training dataset. We employed this dataset to develop a vocabulary model utilizing sentence piece [31]. Following the completion of the sentence piece model's training, it becomes ready for use in tokenizing any binary programs written in the target instruction set.

The fine-tuning phase of CAWET involves its adaptation to basic blocks along with their contextual information. Here, we used a varied and openly accessible collection of programs, namely, The Algorithms (accessible at: <https://github.com/TheAlgorithms/C>), MiBench [28], and Polybench [42].

Data is available at: <https://zenodo.org/records/10043908>.

**A dataset of synthetically generated code blocks for the learning of WCET on Cortex A53** We are releasing the data related to our RTCSA publication [26]: WE-HML Dataset WCET code block with different pollution factors on data cache for Cortex A53.

Data is available at: <https://zenodo.org/records/10043653>.

## 7 New results

**Participants:** Abderaouf Nassim Amalou, Nicolas Bellec, Caroline Collange, Antoine Gicquel, Damien Hardy, Camille Le Bon, Pierre Michaud, Valentin Pasquale, Anis Peysieux, Isabelle Puaut, Erven Rohou.

## 7.1 Compilation and Optimization

**Participants:** Pierre Bedell, Damien Hardy, Camille Le Bon, Mohammed Mehdi Merah, Aurore Poirier, Isabelle Puaut, Hugo Reymond, Erven Rohou.

### 7.1.1 Compilation for Intermittent Systems

**Participants:** Isabelle Puaut, Hugo Reymond, Erven Rohou

**Context:** CominLabs project NOP

**External collaborators:** Sébastien Faucou, Mikaël Briday, Jean-Luc Béchenec, LS2N Nantes

Battery-free devices enable sensing in hard-to-access locations, opening up new opportunities in various fields such as healthcare, space, or civil engineering. Such devices harvest ambient energy and store it in a capacitor. Due to the unpredictable nature of the harvested energy, a power failure can occur at any time, resulting in a loss of all non-persistent information (e.g., processor registers, data stored in volatile memory). Checkpointing volatile data in non-volatile memory allows the system to recover after a power failure, but raises two issues: (i) spatial and temporal placement of checkpoints; (ii) memory allocation of variables between volatile and non-volatile memory, with the overall objective of using energy as efficiently as possible. While many techniques rely on the developer to address these issues, we present SCHEMATIC [20], a compiler technique that automates checkpoint placement and memory allocation to minimize the overall energy consumption. SCHEMATIC ensures that programs will eventually terminate (forward progress property). Moreover, checkpoint placement and memory allocation adapt to the size of the energy buffer and the capacity of volatile memory. SCHEMATIC takes advantage of volatile memory (VM) to reduce the energy consumed, by automatically placing the most used variables in VM. We tested SCHEMATIC for different experimental settings (size of volatile memory and capacitor) and results show an average energy reduction of 51 % compared to related techniques.

### 7.1.2 Dynamic Binary Analysis and Optimization

**Participants:** Aurore Poirier, Erven Rohou

**Context:** Exploratory Action AoT.js

**External collaborators:** Manuel Serrano, INDES/SPLiTS team (Sophia)

Since the creation of JavaScript, web browsers embedding it have been relying on JIT compilers or even interpreters to execute it, a choice that seems natural because of the dynamic aspect of the language. However, researchers have wondered whether it is possible to produce ahead-of-time compiled code that allows dynamic languages to be used in contexts that do not allow the use of JIT compilers (due to limited resources, security policies, response time...) In order to implement performance measurement methods that feed back information from the processor to the high-level languages and design new ways to compile JavaScript, we explored ways to exploit low-level performance measurement from hardware counters.

### 7.1.3 Accurate 3D printing time estimation

**Participants:** Pierre Bedell, Damien Hardy, Camille Le Bon, Mohammed Mehdi Merah

**Context:** Inria Exploratory Action Ofast3D

**External collaborators:** MimeTIC and MFX (Nancy) teams.

Fused deposition modeling 3D printing is a process that requires hours or even days to print a 3D model. To assess the benefits of optimizations, it is mandatory to have a fast 3D printing time estimator to avoid waste of materials and a very long validation process. Furthermore, the estimation must be accurate [29]. To reach that goal, we have modified an existing 3D printer firmwares Klipper in simulation mode to determine the timing per G-code instruction (the language interpreted by 3D printers). This extension named TimeKlip 6.1.7 is printer and slicer agnostic. We conduct an extensive study to highlight the precision and versatility of our simulator on 3D printers with different kinematics, using different slicers. We show that our simulator can be up to 145 times faster than an actual print. Its average error,

without requiring any calibration, is 0.04 % on a total of 66 printed models representing more than 133 hours of print.

See also GATO3D 6.1.4 and OptiPrint 6.1.5.

## 7.2 Processor Architecture

**Participants:** Caroline Collange, Erven Rohou, Sara Sadat Hoseininasab, Pierre Michaud.

### 7.2.1 Automatic synthesis of multi-thread pipelines

**Participants:** Sara Sadat Hoseininasab, Caroline Collange

**Context:** ANR Project DYVE

**External collaborator:** Steven Derrien, TARAN team.

Register-Transfer Level (RTL) design has been a traditional approach in hardware design for several decades. However, with the growing complexity of designs and the need for fast time-to-market, the design and verification process at the RTL level can become impractical. This has motivated for raising the abstraction level in hardware design. High-Level Synthesis (HLS) provides higher-level abstraction by automatically transforming a behavioral specification of a circuit into a low-level RTL, making it easier to design, simulate and verify complex digital systems. HLS relies on statically scheduled data paths which can limit its effectiveness. This limitation makes it difficult to design the micro-architectural features of processors from an Instruction Set Architecture described in high-level languages.

This project aims to demonstrate how the available features of HLS can be deployed in designing various pipelined processors micro-architecture. Our approach takes advantage of the capabilities of HLS and employs multi-threading and dynamic scheduling techniques to overcome the limitation of HLS in pipelining a processor from an Instruction Set Simulator written in C [19].

### 7.2.2 Two-dimensional memory architecture

**Participants:** Pierre Michaud, Erven Rohou

**Context:** ANR Project Maplurinum.

The Maplurinum project revisits the foundations of computer architecture and operating systems for cloud computing and high-performance computing, to better manage the growing heterogeneity of hardware components. A 128-bit address space is considered as one of the possible solutions to achieve this goal. As a participant of the Maplurinum project, the PACAP team explores some of the architectural implications of a 128-bit address space.

Most general-purpose computing systems today implement a virtual memory where the memory addresses seen by the programmer, called virtual addresses, are distinct from the physical memory addresses used by the hardware. During execution of a program, the hardware and the operating system collaborate to translate virtual addresses automatically into physical addresses. A virtual memory provides several advantages, in particular with respect to program portability and process isolation.

As memory needs keep growing, the 64-bit virtual address space of current architecture will eventually be insufficient, and the transition to 128-bit addresses will have to be considered [24, 25]. 128-bit addresses will be beneficial to applications needing a huge amount of memory. However, doubling the virtual address length has a cost in hardware complexity. The transition from 32-bit to 64-bit architectures happened at a time when the hardware cost could be partially hidden by Dennard scaling, making the transition relatively smooth as far as the hardware is concerned, with little impact on clock frequency, silicon area and energy consumption. However, Dennard scaling stopped around 2005, and the hardware cost of transitioning to 128-bit will be a factor to take into account. If the applications that do not benefit from 128-bit addresses still consume more energy than on 64-bit architectures, this will be a problem. The transition will be easier if it benefits a large class of applications, not just the ones requiring a huge amount of memory, and if the hardware cost is small.

We consider the possibility to make the transition to 128-bit addresses coincide with a fundamental architectural change: instead of representing a virtual address as a single 128-bit value, we represent

it as two independent 64-bit coordinates  $X$  and  $Y$ . In other words, we consider replacing the classical one-dimensional (1D) virtual address space with a two-dimensional (2D) one.

Conventional architectures have a 1D virtual address space because it is the simplest option and because the decisions made several decades ago, when simplicity was required, tend to last (e.g., backward binary compatibility). Besides these reasons, the choice of a 1D address space is somewhat arbitrary, as the data manipulated by programs are often inherently multi-dimensional.

We have proposed a 2D-VM architecture called XYA and a programming language called XYZ, similar to the C language but exposing the 2D memory to the programmer. XYA partitions the address space into regions called *books*, each book corresponding to a different page aspect ratio. The compiler or the programmer can select the book where an array is mapped so as to maximize page locality. We have written a C++ library called MXY allowing to write programs in pseudo-XYZ whose execution generates a 2D address trace. With MXY, we did a case study of a 2D fast Fourier transform and have shown that XYA offers new degrees of freedom to programmers for maximizing performance. Besides a variable page aspect ratio, XYA offers other advantages compared to a conventional 128-bit architecture: a 64-bit integer data path (registers, ALUs), efficient 2D array slicing, and efficient 2D array indexing (no multiplication needed). A paper describing XYA is under submission at the ISCA conference.

### 7.3 WCET estimation and optimization

**Participants:** Abderaouf Nassim Amalou, Hector Chabot, Isabelle Puaut.

#### 7.3.1 Using machine learning for timing analysis of complex processors

**Participants:** Abderaouf Nassim Amalou, Isabelle Puaut

**External collaborators:** Elisa Fromont, LACODAM team

Modern processors raise a challenge for timing estimation in general, and WCET estimation in particular, since detailed knowledge of the processor microarchitecture is not available. Our recent work has explored the use of the advanced language processing technique Transformer-XL to estimate the timing of machine code, running on embedded processors. The most recent results in 2023 have targeted worst-case performance (WCET estimation) and average-case performance.

Regarding worst-case performance, we have introduced CAWET, a hybrid worst-case program timing estimation technique [16]. CAWET identifies the longest execution path using static techniques, whereas the worst-case execution time (WCET) of basic blocks is predicted using Transformer-XL. By employing Transformers-XL in CAWET, the execution context formed by previously executed basic blocks is taken into account, allowing for consideration of the micro-architecture of the processor pipeline without explicit modeling. Through a series of experiments on the TacleBench benchmarks, using different target processors (Arm Cortex M4, M7, and A53), our method is demonstrated to never underestimate WCETs and is shown to be less pessimistic than its competitors.

As far as average-case execution time is concerned, we have designed ORXESTRA, a context-aware execution time prediction model based on Transformers XL, specifically designed to accurately estimate performance in embedded system applications. Unlike traditional machine learning models that often overlook contextual information, resulting in biased predictions for individual basic blocks, ORXESTRA overcomes this limitation by incorporating execution context awareness. By doing so, ORXESTRA effectively accounts for the processor micro-architecture without explicitly modeling micro-architectural elements such as caches, pipelines, and branch predictors. Our evaluations demonstrate ORXESTRA's ability to provide precise timing estimations for different ARM targets (Cortex M4, M7, A53, and A72), surpassing existing machine learning-based approaches in both prediction accuracy and prediction speed [17].

This work is done in collaboration with Elisa Fromont, from the LACODAM team, who co-supervizes the PhD thesis of Abderaouf Nassim Amalou, defended in December 2023 [22].

### 7.3.2 Static estimation of memory access profiles for real-time multi-core systems

**Participants:** Hector Chabot, Isabelle Puaut

**External collaborators:** Hugues Cassé, Thomas Carle, IRIT Toulouse

Real-time systems are defined by their need for guarantees of timing properties to ensure a correct behavior. These properties call for the computation of the upper-bound of the Worst-Case Execution Time (WCET) of tasks. The calculated upper-bound needs to be both *safe* and *tight*: the result must be an over-approximation without being too far off from the ground truth. WCET analysis for single-core systems have been well studied in the past thirty years. However, the spread of multi-core usage in real-time systems raises new issues.

In a multi-core context, resource sharing is often used to decrease the overall cost of the platform. This leads to conflicts on concurrent access to shared resource. These conflicts, named *interferences*, causes delays which increase the execution time of tasks. Accounting for interferences delay in WCET computation may rely on interference avoidance, through hardware or software-enforced policies, or on calculation of interference delays on an un-modified multi-core system. Both techniques rely on the knowledge of the shared resources accesses for a given task to find possible interferences. We propose a new technique that extracts memory access profiles for real-time software. In contrast to state-of-the-art studies, the constructed profiles are both *fine-grain* for a more precise calculation of interference delays and *code-based* to allow dynamic schedule reconfiguration for on-line interference reduction.

This work is part of the PhD thesis of Hector Chabot, who is co-supervised by Hugues Cassé and Thomas Carle from IRIT, Toulouse.

## 7.4 Security

**Participants:** Nicolas Bailluet, Nicolas Bellec, Antoine Gicquel, Damien Hardy, Isabelle Puaut, Erven Rohou.

### 7.4.1 Verification of Data Flow Integrity for Real-Time Embedded Systems

**Participants:** Nicolas Bellec (former PhD student in PACAP), Isabelle Puaut

**External collaborators:** Guillaume Hiet, Frédéric Tronel, CIDRE team and Simon Rokicki, TARAN team

Real-time systems have more and more ways to communicate wirelessly with external users. These same means can be hijacked to attack these systems, breaking their guarantees and potentially causing accidents. To protect real-time systems against these new attacks, it is necessary to develop new protections taking into account the specificities of these systems. In this work, we seek to improve the security of real-time systems against so-called memory corruption attacks. These attacks use a bad memory management in a program to modify its behavior. We are particularly interested in a defense called Data Flow Integrity, which can protect against a large class of memory corruption attacks. We adapt this protection to the context of real-time systems by optimizing the worst-case execution time, a fundamental metric to ensure the correct execution of these systems.

This work is done in collaboration with the CIDRE and TARAN teams, who co-supervised the thesis of Nicolas Bellec, defended in May 2023 [23].

### 7.4.2 Multi-nop fault injection attack

**Participants:** Antoine Gicquel, Damien Hardy, Erven Rohou

**External collaborators:** CIDRE and TARAN team.

Multi-fault injection attacks are powerful since they allow to bypass software security mechanisms of embedded devices. Assessing the vulnerability of an application while considering multiple faults with various effects is an open problem due to the size of the fault space to explore. We propose SAMVA [18] (see Section 6.1.6), a framework for efficiently searching vulnerabilities of applications in presence of multiple instruction-skip faults with various widths. SAMVA relies solely on static analysis to determine attack paths in a binary code. It is configurable with the fault injection capacity of the attacker and the attacker's objective. We evaluate the proposed approach on eight PIN verification programs containing

various software countermeasures. Our framework finds numerous attack paths, even for the most hardened version, in very limited time.

### 7.4.3 Gadget chains synthesis driven by SMT Solving for Code-Reuse Attacks

**Participants:** Nicolas Bailluet, Isabelle Puaut, Erven Rohou

**External collaborators:** Emmanuel Fleury, LaBRI Bordeaux.

The final objective of this work is to develop compiler approaches, based on binary code modifications, to protect programs against attacks such as *Return-Oriented Programming* (ROP) or *Jump-Oriented Programming* (JOP).

As a first step towards this ambitious goal, we have designed a new technique for the automatic chaining of gadgets. Performing complex code-reuse attacks require discovering small code snippets (gadgets) and chaining them together to reach the attacker's goal. We have designed a new method to synthesize gadget chains using SMT solving. The proposed method addresses three challenges, yet not solved by other related works: (i) it is able to build gadget chains for arbitrary exploitation contexts (whether the stack is controlled or not); (ii) it guarantees that data that are out of attacker's control do not interfere with the generated gadget chains (robust reachability property); (iii) it is able to leverage multi-path gadgets to synthesize multi-behavior chains that contain conditions and can behave differently based on their execution context. Experiments show the benefit of the robust reachability property and thoroughly evaluate the quality of the approach in terms of expressivity and performance compared to other related chain synthesis techniques.

## 8 Bilateral contracts and grants with industry

**Participants:** Pierre Bedell, Damien Hardy, Camille Le Bon, Pierre Michaud, Erven Rohou.

### 8.1 Bilateral contracts with industry

**Ampere Computing:**

**Participants:** Pierre Michaud.

- Duration: 2023
- Local coordinator: Pierre Michaud
- Collaboration between the PACAP team and Ampere Computing on features of the microarchitecture of next generation CPUs.

**Ofast3D:**

**Participants:** Pierre Bedell, Damien Hardy, Camille Le Bon, Erven Rohou.

- Duration: 2023-2024
- Local coordinator: Damien Hardy
- Collaboration between the PACAP team in the context of the Ofast3D Inria exploratory action and the following companies: 3D News Tech, Cosmyx and Pollen AM.



## 9 Partnerships and cooperations

### 9.1 European initiatives

#### 9.1.1 H2020 projects

**HPCQS** [HPCQS project on cordis.europa.eu](https://cordis.europa.eu)

**Title:** High Performance Computer and Quantum Simulator hybrid

**Duration:** From December 1, 2021 to November 30, 2025

**Partners:**

- INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET AUTOMATIQUE (INRIA), France
- GRAND ÉQUIPEMENT NATIONAL DE CALCUL INTENSIF (GENCI), France
- NATIONAL UNIVERSITY OF IRELAND GALWAY (NUI GALWAY), Ireland
- FORSCHUNGSZENTRUM JÜLICH GMBH (FZJ), Germany
- PARITY QUANTUM COMPUTING GMBH (ParityQC), Austria
- FRAUNHOFER GESELLSCHAFT ZUR FORDERUNG DER ANGEWANDTEN FORSCHUNG EV (Fraunhofer), Germany
- COMMISSARIAT À L'ÉNERGIE ATOMIQUE ET AUX ÉNERGIES ALTERNATIVES (CEA), France
- EURICE EUROPEAN RESEARCH AND PROJECT OFFICE GMBH, Germany
- CONSIGLIO NAZIONALE DELLE RICERCHE (CNR), Italy
- BULL SAS (BULL), France
- FLYSIGHT SRL, Italy
- PARTEC AG (PARTEC), Germany
- UNIVERSITAET INNSBRUCK (UIBK), Austria
- CINECA CONSORZIO INTERUNIVERSITARIO (CINECA), Italy
- CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE (CNRS), France
- CENTRALESUPELEC, France
- BARCELONA SUPERCOMPUTING CENTER CENTRO NACIONAL DE SUPERCOMPUTACION (BSC CNS), Spain
- SORBONNE UNIVERSITÉ, France

**Inria contact:** Luc Giraud

**Coordinator:** Forschungszentrum Jülich GmbH

**Summary:** The aim of HPCQS is to prepare European research, industry and society for the use and federal operation of quantum computers and simulators. These are future computing technologies that are promising to overcome the most difficult computational challenges. HPCQS is developing the programming platform for the quantum simulator, which is based on the European ATOS Quantum Learning Machine (QLM), and the deep, low-latency integration into modular HPC systems based on ParTec's European modular supercomputing concept. A twin pilot system, developed as a prototype by the European company Pasqal, will be implemented and integrated at CEA/TGCC (France) and FZJ/JSC (Germany), both hosts of European Tier-0 HPC systems. The pre-exascale sites BSC (Spain) and CINECA (Italy) as well as ICECH (Ireland) will be connected to the TGCC and JSC via the European data infrastructure FENIX. It is planned to offer quantum HPC hybrid resources to the public via the access channels of PRACE. To achieve these goals, HPCQS brings together leading quantum and supercomputer experts from science and industry, thus creating an incubator for practical quantum HPC hybrid computing that is unique in the world.

The HPC-QS technology will be developed in a co-design process together with selected exemplary use cases from chemistry, physics, optimization and machine learning suitable for quantum HPC hybrid calculations. HPCQS fits squarely to the challenges and scope of the call by acquiring a quantum device with two times 100+ neutral atoms. HPCQS develops the connection between the classical supercomputer and the quantum simulator by deep integration in the modular supercomputing architecture and will provide cloud access and middleware for programming and execution of applications on the quantum simulator through the QLM, as well as a Jupyter-Hub platform with safe access guarantee through the European UNICORE system to its ecosystem of quantum programming facilities and application libraries.

## 9.2 National initiatives

### EPIQ: Study of the quantum stack: Algorithm, models, and simulation for quantum computing

**Participants:** Caroline Collange, Audrey Fauveau.

- Funding: PEPR
- Duration: 2022-2027
- Local coordinator: Caroline Collange
- Partners: CNRS, Inria, CEA
- The EPIQ project aims at developing algorithmic techniques for both noisy quantum machines (NISQ) and fault-tolerant ones so as to facilitate their practical implementation. To this end, a first Work Package (WP) is dedicated to algorithmic techniques, a second one focuses on computational models and languages so as to facilitate the programming of quantum machines and to optimize the code execution steps. Lastly, the third WP aims at developing the simulation techniques of quantum computers.

### ARSENE: Secure architectures for embedded digital systems (ARchitectures SEcurisées pour le Numérique Embarqué)

**Participants:** Damien Hardy, Erven Rohou, Thomas Rubiano.

- Funding: PEPR
- Duration: 2022-2027
- Local coordinator: Ronan Lashermes
- Partners: CNRS, Inria, CEA, UGA, IMT
- The security of communicating objects and the components they integrate is of growing importance in the cybersecurity arena. To address those challenges, the already-rich French research community in embedded systems security is joining forces within the ARSENE project in order to accelerate research & development in this field in a coordinated and structured way to achieve secure solutions. The main objectives of the project are to allow the French community to make significant advances in the field to strengthen the community's expertise and visibility on the international stage. The first part of the ARSENE project is on the study and implementation of two families of RISC-V processors: 32-bit RISC-V for low power secure circuits against physical attacks for IoT applications and 64-bit RISC-V secure circuits against micro-architectural attacks for rich applications. The second aspect of the project pertains to the secure integration of such



new generations of secure processors into System of Chips, to the research and development of secure building blocks for such SoCs like secure and robust Random Number Generators, memory blocks secured against physical attacks, memories instrumented for security and agile hardware accelerators for next generation of cryptography. This work on hardware security is completed by studies on software tools for dynamic annotation of code for next generation of secure embedded software, by the implementation of a secure kernel for an embedded OS and by research work on the dynamic embedded supervision of the system. A last, but very significant, aspect of this project is the implementation of FPGA and ASIC demonstrators integrating the components developed in this project. Those demonstrators shall offer a unique opportunity to showcase the results of the project. This ambitious project will result in increasing the scientific visibility of the research teams involved on the international level, but also in the regional, national and international ecosystems. This project shall trigger a durable, lifelong, cooperation among the main French research teams of the field, not only in terms of scientific achievements, but also for building new collaborative projects on the EU level or other national projects involving industrial partners.

### **EQIP: Engineering for Quantum Information Processors**

**Participants:** Caroline Collange.

- Funding: Inria Challenge project
- Duration: 2021-2024
- Local coordinator: Caroline Collange
- Partners: COSMIQ, CAGE, CASCADE, DEDUCTEAM, GRACE, HIEPACS, MATHERIALS, MOCQUA, PACAP, PARSYS, QUANTIC, STORM, and ATOS Quantum
- Building a functional quantum computer is one of the grand scientific challenges of the 21st century. This formidable task is the object of Quantum Engineering, a new and very active field of research at the interface between physics, computer science and mathematics. EQIP brings together all the competences already present in the institute, to turn Inria into a major international actor in quantum engineering, including both software and hardware aspects of quantum computing.
- website: [project.inria.fr/eqip](http://project.inria.fr/eqip)

### **DYVE: Dynamic vectorization for heterogeneous multi-core processors with single instruction set**

**Participants:** Caroline Collange, Sara Sadat Hoseininasab.

- Funding: ANR, JCJC
- Duration: 2020-2023
- Local coordinator: Caroline Collange
- Most of today's computer systems have CPU cores and GPU cores on the same chip. Though both are general-purpose, CPUs and GPUs still have fundamentally different software stacks and programming models, starting from the instruction set architecture. Indeed, GPUs rely on static vectorization of parallel applications, which demands vector instruction sets instead of CPU scalar instruction sets. In the DYVE project, we advocate a disruptive change in both CPU and GPU architecture by introducing Dynamic Vectorization at the hardware level.

Dynamic Vectorization aims to combine the efficiency of GPUs with the programmability and compatibility of CPUs by bringing them together into heterogeneous general-purpose multicores.

It will enable processor architectures of the next decades to provide (1) high performance on sequential program sections thanks to latency-optimized cores, (2) energy-efficiency on parallel sections thanks to throughput-optimized cores, (3) programmability, binary compatibility and portability.

### **NOP: Safe and Efficient Intermittent Computing for a Batteryless IoT**

**Participants:** Isabelle Puaut, Hugo Reymond, Erven Rohou.

- Funding: LabEx CominLabs (50 %)
- Duration: 2021-2024
- Local coordinator: Erven Rohou
- Partners: IRISA/Granit Lannion, LS2N/STR Nantes, IETR/Syscom Nantes
- Intermittent computing is an emerging paradigm for batteryless IoT nodes powered by harvesting ambient energy. It intends to provide transparent support for power losses so that complex computations can be distributed over several power cycles. It aims at significantly increasing the complexity of software running on these nodes, and thus at reducing the volume of outgoing data, which improves the overall energy efficiency of the whole processing chain, reduces reaction latencies, and, by limiting data movements, preserves anonymity and privacy.

NOP aims at improving the efficiency and usability of intermittent computing, based on consolidated theoretical foundations and a detailed understanding of energy flows within systems. For this, it brings together specialists in system architecture, energy-harvesting IoT systems, compilation, and real-time computing, to address the following scientific challenges:

1. develop sound formal foundations for intermittent systems,
  2. develop precise predictive energy models of a whole node (including both harvesting and consumption) usable for online decision making,
  3. significantly improve the energy efficiency of run-time support for intermittency,
  4. develop techniques to provide formal guarantee through static analysis of the systems behavior (forward progress),
  5. develop a proof of concept: an intermittent system for bird recognition by their songs, to assess the costs and benefits of the proposed solutions.
- website: [project.inria.fr/nopcl/](http://project.inria.fr/nopcl/)

### **CAOTIC: Collaborative Action on Timing Interference**

**Participants:** Hector Chabot, Isabelle Puaut.

- Funding: ANR
- Duration: 2022-2026
- Local coordinator: Isabelle Puaut
- Partners: CEA List, Inria, Univ Rennes/IRISA, IRIT, IRT Saint Exupery, LS2N, LTCI, Verimag (Project Coordinator)

- Project CAOTIC is an ambitious initiative aimed at pooling and coordinating the efforts of major French research teams working on the timing analysis of multicore real-time systems, with a focus on interference due to shared resources. The objective is to enable the efficient use of multicore in critical systems. Based on a better understanding of timing anomalies and interference, taking into account the specificities of applications (structural properties and execution model), and revisiting the links between timing analysis and synthesis processes (code generation, mapping, scheduling), significant progress is targeted in timing analysis models and techniques for critical systems, as well as in methodologies for their application in industry.

In this context, the originality and strength of the CAOTIC project resides in the complementarity of the approaches proposed by the project members to address the same set of scientific challenges: (i) build a consistent and comprehensive set of methods to quantify and control the timing interferences and their impact on the execution time of programs; (ii) define interference-aware timing analysis and real-time scheduling techniques suitable for modern multi-core real-time systems; (iii) consolidate these methods and techniques in order to facilitate their transfer to industry.

- website: [anr-caotic.imag.fr/](http://anr-caotic.imag.fr/)

### **OWL: Operating Within Limits**

**Participants:** Erven Rohou, Isabelle Puaut.

- Funding: ANR
- Duration: 2023-2027
- Local coordinator: Erven Rohou
- Partners: IRISA/Granit Lannion, LS2N/STR Nantes (Project Coordinator), LS2N/SIMS Nantes
- Project OWL proposes a new model of computation for more frugal intelligent autonomous sensors: circadian artificial intelligence (AI). The targeted applications are in the field of environmental monitoring, especially bioacoustic and its application to conservation ecology. This model is particularly well suited for sensors without batteries that are intermittently powered by ambient energy. The great promises of these systems is the extension of their lifetime without the need for human intervention allowing for long-term biostatistics observation missions, and a lower impact on the environment thanks to the absence of battery.

Circadian AI is interested in observing phenomena that have a period of one day, such as the activity of birds or the pollution associated with traffic in a metropolis. It exploits the fact that this period is shared with the availability of solar energy, which is used to power the sensors. This correlation allows the systems to temporally shift the costly computations required to perform the AI functions to times when the observed phenomenon is at rest and energy is abundant.

The project proposes two main contributions. The first is to design new algorithms for circadian AI that allow for this temporal shift in computation. The second is to provide the software and hardware infrastructure necessary to run circadian AI on intermittently powered sensors.

The work done in the project will be based as much as possible on open source / open hardware technologies. Those built during the project (dataset, software, hardware design) will all be freely distributed.

### **LOTR: Lord Of The RISCs**

**Participants:** Isabelle Puaut.

- Funding: ANR

- Duration: 2023-2027
- Local coordinator: Steven Derrien (Univ. Rennes/IRISA)
- Partners: CEA List, Univ. Rennes/IRISA (coordinator)
- Lord Of The RISCs (LOTR) is a novel flow for designing highly customized RISC-V processor microarchitectures for embedded and IoT platforms. The LOTR flow operates on a description of the processor Instruction Set Architecture (ISA). It can automatically infer synthesizable Register Transfer Level (RTL) descriptions of a large number of microarchitecture variants with different performance/cost trade-offs. In addition, the flow integrates two domain-specific toolboxes dedicated to the support of timing predictability (for safety-critical systems) and security (through hardware protection mechanisms)

### **AIxIA (Artificial Intelligence for Interference Analysis)**

**Participants:** Isabelle Puaut.

- Funding: FRAE (Fondation de Recherche pour l'Aéronautique et l'Espace) AIRSTRIP (L'intelligence Artificielle au service de l'Ingénierie des Systèmes aéronautiques et spatiaux) project
- Duration: 2024-2026
- Local coordinator: Isabelle Puaut
- Partners: IRT Saint Exupéry, INRIA Bordeaux, IRIT, Univ. Rennes/IRISA
- Demonstrating the satisfaction of temporal performance in an embedded software with the required level of confidence is a difficult and costly task. One of the main issues is accounting for temporal interference phenomena that occur between software applications sharing elements of the execution structure (e.g., cores, GPU, etc.). In this context, the AIxIA project aims to study the contribution of artificial intelligence techniques to identifying these interferences and analyzing their effects. The project will apply artificial intelligence techniques to three dimensions of the problem: (i) identifying sources of interference, (ii) quantifying and predicting their effects, and (iii) avoidance.

### **Maplurinum (Machinae pluribus unum): (make) one machine out of many**

**Participants:** Pierre Michaud, Erven Rohou.

- Funding: ANR, PRC
- Duration: 2021-2024
- Local coordinator: Pierre Michaud
- Partners: Télécom Sud Paris/PDS, CEA List, Université Grenoble Alpes/TIMA
- Cloud and high-performance architectures are increasingly heterogeneous and often incorporate specialized hardware. We have first seen the generalization of GPUs in the most powerful machines, followed a few years later by the introduction of FPGAs. More recently we have seen nascence of many other accelerators such as tensor processor units (TPUs) for DNNs or variable precision FPUs. Recent hardware manufacturing trends make it very likely that specialization will not only persist, but increase in future supercomputers. Because manually managing this heterogeneity in each application is complex and not maintainable, we propose in this project to revisit how we design both hardware and operating systems in order to better hide the heterogeneity to supercomputer users.

- website: [project.inria.fr/maplurinum/](https://project.inria.fr/maplurinum/)

### Ofast3D

**Participants:** Pierre Bedell, Damien Hardy, Camille Le Bon, Erven Rohou.

- Funding: Inria Exploratory Action
- Duration: 2021-2024
- Local coordinator: Damien Hardy
- Partners: MimeTIC (Rennes) and MFX (Nancy)
- The goal of Ofast3D is to increase the production capacity of fused deposition modeling 3D printing, without requiring any modification of existing production infrastructures. Ofast3D aims to reduce printing time without impacting the print quality by optimizing the code interpreted by 3D printers during its generation by taking into account the geometry of 3D models. Ofast3D is complementary to methods aiming either at improving printers or at optimizing 3D models.
- website: [project.inria.fr/ofast3d](https://project.inria.fr/ofast3d)

### AoT.js

**Participants:** Aurore Poirier, Erven Rohou.

- Funding: Inria Exploratory Action
- Duration: 2022-2025
- Local coordinator: Erven Rohou
- Partners: INDES/SPLiTS (Sophia)
- JavaScript programs are typically executed by a JIT compiler, able to handle efficiently the dynamic aspects of the language. However, JIT compilers are not always viable or sensible (*e.g.*, on constrained IoT systems, due to secured read-only memory (W $\oplus$ X), or because of the energy spent recompiling again and again). We propose to rely on ahead-of-time compilation, and achieve performance thanks to optimistic compilation, and detailed analysis of the behavior of the processor, thus requiring a wide range of expertise from high-level dynamic languages to microarchitecture.

## 10 Dissemination

**Participants:** Abderaouf Nassim Amalou, Nicolas Bailluet, Pierre Bedell, Hector Chabot, Caroline Collange, Antoine Gicquel, Damien Hardy, Sara Sadat Hoseininasab, Pierre Michaud, Anis Peysieux, Aurore Poirier, Isabelle Puaut, Hugo Raymond, Erven Rohou.

## 10.1 Promoting scientific activities

### 10.1.1 Scientific events: selection

#### Member of the conference program committees

- E. Rohou was a member of the program committees of the CC 2024 conference.
- P. Michaud was a member of the program committee of the ISCA 2023 conference.
- I. Puaut is member of the program committee of Euromicro Conference on Real Time Systems (ECRTS) 2024.
- I. Puaut is member of the program committee of IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS), 2024.
- I. Puaut was member of the program committee of Conference on Real-Time Networks and Systems (RTNS) 2023.
- D. Hardy was member of the program committee of the International Symposium on Advanced Security on Software and Systems (ASSS) 2023.
- C. Collange was member of the external review committee of the ASPLOS 2023 conference.

**Reviewer** Members of PACAP routinely review submissions to international conferences and events.

### 10.1.2 Journal

**Member of the editorial boards** I. Puaut is associate editor of the Springer International Journal of Time-Critical Computing Systems.

**Reviewer - reviewing activities** Members of PACAP routinely review submissions to international journals.

### 10.1.3 Invited talks

- A. Gicquel was invited to give a talk at the “École de Printemps Recherche” of the CyberSchool in Apr 2023.
- H. Reymond was invited to present our results entitled “Towards Sustainable IoT Nodes” [21] at GreenDays 2023 in Lyon, France, on Mar 27th, 2023. In collaboration with LS2N, Nantes.

### 10.1.4 Leadership within the scientific community

E. Rohou is a member of the organization committee of focused days within the GdR SoC2.

### 10.1.5 Research administration

- E. Rohou is the contact for international relations for Inria Rennes Bretagne Atlantique (for scientific matters). He co-organized the 2023 edition of the French-American Doctoral Exchange Seminar (FADEx).
- E. Rohou is a member of the steering committee of the high security research laboratory (LHS).
- E. Rohou is a member of the steering committee of the Inria Rennes computer grid “igrida”.
- I. Puaut is member of the Advisory board of the Euromicro Conference on Real Time Systems (ECRTS).
- I. Puaut is member of the steering committee of the WCET workshop, satellite workshop to ECRTS.

- Since Nov 2023, I. Puaut is elected member of section 27 of CNU (*Conseil National des Universités* – National Council of Universities). The CNU is a national consultative and decision-making body. It makes decisions regarding the career progression of assistant professors and professors in institutions under the jurisdiction of the Ministry of Higher Education and Research (MESR).
- Since Feb 2023, I. Puaut is a member of the thesis committee (*comité des thèses*) at the Matisse doctoral school. The committee is responsible for reviewing thesis registration applications and forming juries. The thesis committee oversees the 250 doctoral students hosted at IRISA.

## 10.2 Teaching - Supervision - Juries

### 10.2.1 Teaching

- Master: C. Collange, GPU programming, 20 hours, M1, Université de Rennes, France
- Licence: D. Hardy, Real-time systems, 95 hours, L3, Université de Rennes, France
- Master: D. Hardy, Operating systems, 59 hours, M1, Université de Rennes, France
- Master: D. Hardy, Students project, 30 hours, M1, Université de Rennes, France
- Master: I. Puaut, Advanced Operating Systems (SEA), 100 hours, M1, Université de Rennes
- Master: I. Puaut, Low Level Programming (LLP), 40 hours, Université de Rennes
- Master: I. Puaut, Writing of scientific publications, 9 hours, M2 and PhD students, Université de Rennes
- Licence: N. Bailluet, Programmation C et réseau, 20 hours, L3, ENS Rennes, France
- Licence: N. Bailluet, Programmation fonctionnelle, 20 hours, L1, Université de Rennes, France
- Master: N. Bailluet, Noyaux de systèmes d'exploitation, 24 hours, M1, Université de Rennes, France
- Licence: A. Poirier, ALG2 (Algorithmique 2), 18 hours, L2, Université de Rennes, France
- Master: A. Poirier, NFS (Network For Security), 15 hours, M1, Université de Rennes, France
- Master: H. Reymond, OS and Real-Time (SEL, ITR), 32 hours, M1, Université de Rennes, France
- Licence: H. Reymond, OS and Architecture (SIN), 30 hours, L3, Université de Rennes, France
- Licence: H. Chabot, PO (Programmation Objet), 15 hours, L2, Université de Rennes, France
- Licence: H. Chabot, SYRE (Base de système et réseaux), 18 hours, L3, Université de Rennes, France

### 10.2.2 Supervision

- PhD: Nicolas Bellec, *Security in real-time embedded systems* [23], May 2023, advisors I. Puaut (50 %), G. Hiet from CIDRE (25 %), F. Tronel from CIDRE (25 %)
- PhD: Abderaouf Nassim Amalou, *Machine learning for performance prediction* [22], Dec 2023, advisors I. Puaut (75 %), E. Fromont from LACODAM (25 %)
- PhD in progress: Hector Chabot, *Fine grain software modeling and analysis for interference management in multi-core real-time systems*, started Sep 2023, advisors I. Puaut (50 %), H. Cassé and T. Carle (IRIT, Toulouse, 25 % each)
- PhD in progress: Anis Peysieux, *Towards simple and highly efficient execution cores*, started Jan 2020, advisor since Jan 2021: P. Michaud (A. Seznec was advisor from Jan 2020 until Dec 2020)
- PhD in progress: Hugo Reymond, *Energy-aware execution model in intermittent systems*, started Oct 2021, advisors I. Puaut, E. Rohou, S. Faucou (LS2N Nantes), J.-L. Béchenec (LS2N Nantes)

- PhD in progress: Antoine Gicquel, *Étude de vulnérabilité d'un programme au format binaire en présence de fautes précises et nombreuses : métriques et contremesures*, started Sep 2021, advisors D. Hardy, E. Rohou, K. Heydemann (Sorbonne Université)
- PhD in progress: Sara Hoseininasab, *Automatic synthesis of multi-thread pipelines*, started Nov 2021, advisors C. Collange (70 %) and S. Derrien (30 %, TARAN)
- PhD in progress: Aurore Poirier, *Profile-Guided optimization for Dynamic Languages*, started Oct 2022, advisors E. Rohou (50 %) and M. Serrano (50 %, Inria Sophia)
- PhD in progress: Nicolas Bailluet, *Approches par modification de code machine pour la défense contre les attaques ROP et JOP*, started Sep 2022, advisors I. Puaut (50 %) and E. Rohou (50 %)

### 10.2.3 Juries

E. Rohou was member of the following hiring committees:

- Professor, ENS Lyon
- “commission moyens incitatifs”, Centre Inria de l’Université de Rennes
- “comité recrutement” for a team assistant, Centre Inria de l’Université de Rennes

I. Puaut was member of the hiring committee for a Junior Professor Position (CPI) on Artificial Intelligence, University of Rennes, 2023

E. Rohou was a member of the following PhD thesis committees:

- Nicolas Derumigny, *Throughput Optimization Techniques for Heterogeneous Architectures*, Dec 2023 (reviewer)

I. Puaut was a member of the following PhD thesis committees:

- Samira Ait Bensaid, Formal Semantics of Hardware Compilation Framework, Université Paris Saclay, Nov 2023 (examiner)
- Sandro Grebant, Efficient tree-based symbolic WCET computation, Université de Lille, Nov 2023, (reviewer)
- Mathieu Leonel Mba, génération automatique de plateformes matérielles distribuées pour des applications de traitement du signal, Sorbonne Université, Sep 2023, (examiner, president of committee)
- Jatin Arora, Shared resource contention-aware schedulability analysis of hard real-time systems, University of Porto, Portugal, Sep 2023 (external examiner)
- Zhenyu Bai, modélisation du comportement temporel du pipeline pour le calcul du WCET, Université de Toulouse 3 – Paul Sabatier, Toulouse, May 2023 (reviewer)
- Michael Platzer, predictable and performant computer architectures for time-critical systems, Technische Universität Wien, Vienna, Austria, Feb 2023 (reviewer)

C. Collange was a member of the following PhD thesis committees:

- Kevin Mambu, Modèle de programmation bas niveau pour architecture de calcul proche mémoire, CEA / Grenoble Université, Mar 2023 (examiner)
- Alexandre Clément, Graphical Languages for Quantum Control and Linear Optics, LORIA / Université de Lorraine, May 2023 (examinee)

E. Rohou was a member of the CSI of Bruno Mateu, Jean-Loup Hatchikian-Houdot, Aaron Randriana, Ikram Dendani. I. Puaut was member of the CSI of Jean-Michel Gorius and Constance Bocquillon.



## 10.3 Popularization

### 10.3.1 Internal or external Inria responsibilities

E. Rohou is a member of the Inria Rennes working group on sustainable development.

### 10.3.2 Articles and contents

D. Hardy presented the Ofast3D Inria exploratory action in *Émergences*.

### 10.3.3 Education

E. Rohou was invited to present the job of a researcher to secondary-school students (classe de 4e) at Collège de Bourgchevreuil, Cesson-Sévigné.

### 10.3.4 Interventions

- E. Rohou contributed to the program “1 scientifique, 1 classe : Chiche !”, Lycée Sainte-Thérèse, Quimper.
- P. Bedell and D. Hardy presented and demonstrate the Ofast3D Inria exploratory action in the 3D Print Congress & Exhibition in Paris.

## 11 Scientific production

### 11.1 Major publications

- [1] F. Bodin, T. Kisuki, P. M. W. Knijnenburg, M. F. P. O’Boyle and E. Rohou. ‘Iterative Compilation in a Non-Linear Optimisation Space’. In: *Workshop on Profile and Feedback-Directed Compilation (FDO-1), in conjunction with PACT’98*. Paris, France, Oct. 1998.
- [2] N. Hallou, E. Rohou, P. Clauss and A. Ketterlin. ‘Dynamic Re-Vectorization of Binary Code’. In: *SAMOS*. July 2015. URL: <https://hal.inria.fr/hal-01155207>.
- [3] D. Hardy and I. Puaut. ‘Static probabilistic Worst Case Execution Time Estimation for architectures with Faulty Instruction Caches’. In: *21st International Conference on Real-Time Networks and Systems*. Sophia Antipolis, France, Oct. 2013. DOI: [10.1145/2516821.2516842](https://doi.org/10.1145/2516821.2516842). URL: <https://hal.inria.fr/hal-00862604>.
- [4] D. Hardy, I. Sideris, N. Ladas and Y. Sazeides. ‘The performance vulnerability of architectural and non-architectural arrays to permanent faults’. In: *MICRO 45*. Vancouver, Canada, Dec. 2012. URL: <https://hal.inria.fr/hal-00747488>.
- [5] P. Michaud. ‘Best-Offset Hardware Prefetching’. In: *International Symposium on High-Performance Computer Architecture*. Barcelona, Spain, Mar. 2016. DOI: [10.1109/HPCA.2016.7446087](https://doi.org/10.1109/HPCA.2016.7446087). URL: <https://hal.inria.fr/hal-01254863>.
- [6] P. Michaud, A. Mondelli and A. Sez nec. ‘Revisiting Clustered Microarchitecture for Future Superscalar Cores: A Case for Wide Issue Clusters’. In: *ACM Transactions on Architecture and Code Optimization (TACO)* 13.3 (Aug. 2015), p. 22. DOI: [10.1145/2800787](https://doi.org/10.1145/2800787). URL: <https://hal.inria.fr/hal-01193178>.
- [7] A. Perais and A. Sez nec. ‘EOLE: Paving the Way for an Effective Implementation of Value Prediction’. In: *International Symposium on Computer Architecture*. Vol. 42. ACM/IEEE. Minneapolis, MN, United States, June 2014, pp. 481–492. DOI: [10.1109/ISCA.2014.6853205](https://doi.org/10.1109/ISCA.2014.6853205). URL: <https://hal.inria.fr/hal-01088130>.
- [8] A. Perais and A. Sez nec. ‘Practical data value speculation for future high-end processors’. In: *International Symposium on High Performance Computer Architecture*. IEEE. Orlando, FL, United States, Feb. 2014, pp. 428–439. DOI: [10.1109/HPCA.2014.6835952](https://doi.org/10.1109/HPCA.2014.6835952). URL: <https://hal.inria.fr/hal-01088116>.

- [9] E. Rohou, B. Narasimha Swamy and A. Sez nec. ‘Branch Prediction and the Performance of Interpreters - Don’t Trust Folklore’. In: *International Symposium on Code Generation and Optimization*. Burlingame, United States, Feb. 2015. URL: <https://hal.inria.fr/hal-01100647>.
- [10] D. Sampaio, R. M. De Souza, C. Collange and F. M. Quintão Pereira. ‘Divergence Analysis’. In: *ACM Transactions on Programming Languages and Systems (TOPLAS)* 35.4 (Nov. 2013), 13:1–13:36. DOI: [10.1145/2523815](https://doi.org/10.1145/2523815). URL: <https://hal.inria.fr/hal-00909072>.
- [11] S. Sardashti, A. Sez nec and D. A. Wood. ‘Skewed Compressed Caches’. In: *47th Annual IEEE/ACM International Symposium on Microarchitecture, 2014*. Minneapolis, United States, Dec. 2014. URL: <https://hal.inria.fr/hal-01088050>.
- [12] S. Sardashti, A. Sez nec and D. A. Wood. ‘Yet Another Compressed Cache: a Low Cost Yet Effective Compressed Cache’. In: *ACM Transactions on Architecture and Code Optimization* (Sept. 2016), p. 25. URL: <https://hal.inria.fr/hal-01354248>.
- [13] A. Sez nec and P. Michaud. ‘A case for (partially)-tagged geometric history length branch prediction’. In: *Journal of Instruction Level Parallelism* (Feb. 2006). URL: <http://www.jilp.org/vol18>.
- [14] M. Y. Siraichi, V. F. d. Santos, C. Collange and F. M. Quintão Pereira. ‘Qubit allocation as a combination of subgraph isomorphism and token swapping’. In: *OOPSLA*. Vol. 3. Athens, Greece, 10th Oct. 2019, pp. 1–29. DOI: [10.1145/3360546](https://doi.org/10.1145/3360546). URL: <https://hal.inria.fr/hal-02316820>.
- [15] A. Tino, C. Collange and A. Sez nec. ‘SIMT-X: Extending Single-Instruction Multi-Threading to Out-of-Order Cores’. In: *ACM Transactions on Architecture and Code Optimization* 17.2 (May 2020), p. 15. DOI: [10.1145/3392032](https://doi.org/10.1145/3392032). URL: <https://hal.inria.fr/hal-02542333>.

## 11.2 Publications of the year

### International peer-reviewed conferences

- [16] A. N. Amalou, E. Fromont and I. Puaut. ‘CAWET: Context-Aware Worst-Case Execution Time Estimation Using Transformers’. In: *Leibniz International Proceedings in Informatics (LIPIcs), Volume 262, pp. 7:1-7:20, Schloss Dagstuhl - Leibniz-Zentrum für Informatik*. ECRTS 2023 - 35th Euromicro Conference on Real-Time Systems. Vol. 262. Vienne, Austria: Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 3rd July 2023, 7:1–7:20. DOI: [10.4230/LIPIcs.ECRTS.2023.7](https://doi.org/10.4230/LIPIcs.ECRTS.2023.7). URL: <https://hal.science/hal-04148587>.
- [17] A. N. Amalou, E. Fromont and I. Puaut. ‘Fast and Accurate Context-Aware Basic Block Timing Prediction using Transformers’. In: *Proceedings of the ACM SIGPLAN 2024 International Conference on Compiler Construction*. ACM SIGPLAN 2024 International Conference on Compiler Construction. Edimbourg, United Kingdom, 2nd Mar. 2024. DOI: [10.1145/nnnnnnnn.nnnnnnnn](https://doi.org/10.1145/nnnnnnnn.nnnnnnnn). URL: <https://hal.science/hal-04406073>.
- [18] A. Gicquel, D. Hardy, K. Heydemann and E. Rohou. ‘SAMVA: Static Analysis for Multi-Fault Attack Paths Determination’. In: *Constructive Side-Channel Analysis and Secure Design*. COSADE 2023 - 14th International Workshop on Constructive Side-Channel Analysis and Secure Design. Munich (Allemagne), Germany: Springer, 23rd Mar. 2023, pp. 3–22. DOI: [10.1007/978-3-031-29497-6\\_1](https://doi.org/10.1007/978-3-031-29497-6_1). URL: <https://hal.science/hal-03980128>.
- [19] S. S. Hoseininasab, C. Collange and S. Derrien. ‘Rapid Prototyping of Complex Micro-architectures Through High-Level Synthesis’. In: *Applied Reconfigurable Computing*. ARC 2023 - 19th International Symposium on Applied Reconfigurable Computing. Vol. 14251. Lecture Notes in Computer Science. Cottbus, Germany: Springer Nature Switzerland, 2023, pp. 19–34. DOI: [10.1007/978-3-031-42921-7\\_2](https://doi.org/10.1007/978-3-031-42921-7_2). URL: <https://hal.science/hal-04225360>.
- [20] H. Reymond, J.-L. Béchenec, M. Briday, S. Faucou, I. Puaut and E. Rohou. ‘SCHEMATIC: Compile-time checkpoint placement and memory allocation for intermittent systems’. In: *Proceedings of the IEEE/ACM International Symposium on Code Generation and Optimization (CGO’24)*. IEEE/ACM International Symposium on Code Generation and Optimization (CGO’24). Edinburgh, United Kingdom, 2nd Mar. 2024. URL: <https://hal.science/hal-04345348>.

### Conferences without proceedings

- [21] A. Bernabeu and H. Reymond. ‘Towards Sustainable IoT Nodes’. In: GreenDays 2023 - Efficacité énergétique, impacts environnementaux du numérique, sobriété et frugalité numérique : une vision décloisonnée ! Lyon, France, 2023, pp. 1–39. URL: <https://hal.science/hal-04385269>.

### Doctoral dissertations and habilitation theses

- [22] A. N. Amalou. ‘Machine learning for timing estimation’. Université de Rennes, 12th Dec. 2023. URL: <https://hal.science/tel-04406029>.
- [23] N. Bellec. ‘Security enhancement in embedded hard real-time systems’. Université de Rennes, 23rd May 2023. URL: <https://inria.hal.science/tel-04219240>.

### Other scientific publications

- [24] M. Bacou, A. Chader, C. S. Deshpande, C. Fabre, C. Fuguet, P. Michaud, A. Perais, F. Pétrot, G. Thomas and E. Tomasi Ribeiro. ‘128-bit addresses for the masses (of memory and devices)’. In: HotInfra 2023 - Workshop on Hot Topics in System Infrastructure. Orlando, United States, 2023. URL: <https://hal.science/hal-04161640>.
- [25] M. Bacou, A. Chader, C. S. Deshpande, C. Fabre, C. Fuguet, P. Michaud, A. Perais, F. Pétrot, G. Thomas and E. Tomasi Ribeiro. ‘We had 64-bit, yes. What about second 64-bit?’ In: RISC-V Summit Europe 2023. Barcelona, Spain, 2023. URL: <https://hal.science/hal-04161612>.

## 11.3 Cited publications

- [26] A. N. Amalou, I. Puaut and G. Muller. ‘WE-HML: hybrid WCET estimation using machine learning for architectures with caches’. In: *RTCSA 2021 - 27th IEEE International Conference on Embedded Real-Time Computing Systems and Applications*. Online Virtual Conference, France: IEEE, Aug. 2021, pp. 1–10. URL: <https://inria.hal.science/hal-03280177>.
- [27] A. Cohen and E. Rohou. ‘Processor Virtualization and Split Compilation for Heterogeneous Multicore Embedded Systems’. In: *DAC*. Anaheim, CA, USA, June 2010, pp. 102–107.
- [28] M. R. Guthaus, J. S. Ringenberg, D. Ernst, T. M. Austin, T. Mudge and R. B. Brown. ‘MiBench: A free, commercially representative embedded benchmark suite’. In: *4th IEEE international workshop on workload characterization*. 2001.
- [29] D. Hardy. *Ofast3D - Étude de faisabilité*. Technical Report RT-0511. Inria Rennes - Bretagne Atlantique ; IRISA, Dec. 2020, p. 18. URL: <https://hal.inria.fr/hal-03093905>.
- [30] M. Hataba, A. El-Mahdy and E. Rohou. ‘OJIT: A Novel Obfuscation Approach Using Standard Just-In-Time Compiler Transformations’. In: *International Workshop on Dynamic Compilation Everywhere*. Jan. 2015.
- [31] T. Kudo and J. Richardson. ‘Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing’. In: *arXiv preprint arXiv:1808.06226* (2018).
- [32] R. Kumar, D. M. Tullsen, N. P. Jouppi and P. Ranganathan. ‘Heterogeneous chip multiprocessors’. In: *IEEE Computer* 38.11 (Nov. 2005), pp. 32–38.
- [33] C. Le Bon. ‘Analyse et optimisation dynamiques de programmes au format binaire pour la cybersécurité’. Theses. Université Rennes 1, July 2022. URL: <https://theses.hal.science/tel-03906421>.
- [34] P. Michaud and A. Sez nec. ‘Pushing the branch predictability limits with the multi-poTAGE+SC predictor : **Champion in the unlimited category**’. In: *4th JILP Workshop on Computer Architecture Competitions (JWAC-4): Championship Branch Prediction (CBP-4)*. Minneapolis, United States, June 2014. URL: <https://hal.archives-ouvertes.fr/hal-01087719>.
- [35] R. Omar, A. El-Mahdy and E. Rohou. ‘Arbitrary control-flow embedding into multiple threads for obfuscation: a preliminary complexity and performance analysis’. In: *Proceedings of the 2nd international workshop on Security in cloud computing*. ACM. 2014, pp. 51–58.

- [36] R. Puri, D. S. Kung, G. Janssen, W. Zhang, G. Domeniconi, V. Zolotov, J. Dolby, J. Chen, M. Choudhury, L. Decker et al. 'CodeNet: A large-scale AI for code dataset for learning a diversity of coding tasks'. In: *arXiv preprint arXiv:2105.12655* (2021).
- [37] E. Riou, E. Rohou, P. Clauss, N. Hallou and A. Ketterlin. 'PADRONE: a Platform for Online Profiling, Analysis, and Optimization'. In: *Dynamic Compilation Everywhere*. Vienna, Austria, Jan. 2014.
- [38] A. Sembrant, T. Carlson, E. Hagersten, D. Black-Shaffer, A. Perais, A. Seznec and P. Michaud. 'Long Term Parking (LTP): Criticality-aware Resource Allocation in OOO Processors'. In: *International Symposium on Microarchitecture, Micro 2015*. Proceeding of the International Symposium on Microarchitecture, Micro 2015. Honolulu, United States: ACM, Dec. 2015. URL: <https://hal.inria.fr/hal-01225019>.
- [39] A. Seznec. 'TAGE-SC-L Branch Predictors: **Champion in 32Kbits and 256 Kbits category**'. In: *JILP - Championship Branch Prediction*. Minneapolis, United States, June 2014. URL: <https://hal.inria.fr/hal-01086920>.
- [40] A. Seznec, J. San Miguel and J. Albericio. 'The Inner Most Loop Iteration counter: a new dimension in branch history'. In: *48th International Symposium On Microarchitecture*. Honolulu, United States: ACM, Dec. 2015, p. 11. URL: <https://hal.inria.fr/hal-01208347>.
- [41] A. Seznec and N. Sendrier. 'HAVEGE: A user-level software heuristic for generating empirically strong random numbers'. In: *ACM Transactions on Modeling and Computer Simulation (TOMACS)* 13.4 (2003), pp. 334–346.
- [42] T. Yuki. 'Understanding polybench/C 3.2 kernels'. In: *International workshop on polyhedral compilation techniques (IMPACT)*. 2014, pp. 1–5.