# Activity Report 2015

# **Project-Team COMPSYS**

# Compilation and embedded computing systems

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

# Table of contents

<div align="center">

**Project-Team COMPSYS**

</div>

*Creation of the Project-Team: 2004 January 01, updated into Team: 2016 January 01*

**Keywords:**

### Computer Science and Digital Science:

- 2.1.1. - Semantics of programming languages
- 2.1.10. - Domain-specific languages
- 2.1.6. - Concurrent programming
- 2.2.1. - Static analysis
- 2.2.5. - GPGPU, FPGA, etc.
- 2.4.1. - Analysis
- 6.2.6. - Optimization
- 6.2.7. - High performance computing
- 7.2. - Discrete mathematics, combinatorics

### Other Research Topics and Application Domains:

- 6.6. - Embedded systems
- 9.4.1. - Computer science

*Compsys is located at Ecole normale supérieure de Lyon.*

# 1. Members

**Research Scientists**

Alain Darte [Team leader, CNRS, Senior Researcher, HdR]
Christophe Alias [Inria, Researcher, until Sep. 2015]
Tomofumi Yuki [Inria, Researcher]

**Faculty Members**

Paul Feautrier [ENS Lyon, Emeritus Professor, HdR]
Laure Gonnord [Univ. Lyon I, Associate Professor, until Sep. 2015]

**PhD Students**

Guillaume Iooss [PhD student, ENS-Lyon/Colorado State Univ.]
Alexandre Isoard [PhD student, ENS-Lyon]
Maroua Maalej [PhD student, Univ. Lyon, until Sep. 2015]

**Visiting Scientists**

Fernando Magno Quintao Pereira [ENS Lyon, Jan. 2015–Feb. 2015]
Tristan Dubois [Master 1 internship, Lyon 1, Jan. 2015–Feb. 2015]
Adilla Susungi [Master 2 internship, Inria, Mar. 2015–Jul. 2015]
Marc Vincenti [Master 1 internship, Lyon 1, Jan. 2015–Feb. 2015]

**Administrative Assistant**

Evelyne Blesle [Inria]

# 2. Overall Objectives

## 2.1. Introduction

> Keywords: Compilation, code analysis, code optimization, memory optimization, combinatorial optimization, algorithmics, polyhedral optimization, hardware accelerators, high-level synthesis, high-performance computing.

Compsys develops compilation techniques, more precisely code analysis and code optimization techniques, to help programming or designing "embedded computing systems" or platforms for "small" HPC (High-Performance Computing). The team focuses on both low-level (back-end) optimizations and high-level (front-end, mainly source-to-source) transformations, for specialized processors, programmable hardware accelerators (GPU, multicores), and FPGA platforms (high-level synthesis). Recent activities include a shift towards the analysis of parallel languages, and links with abstract interpretation and program termination. The main characteristic of Compsys is its use of algorithmic and formal methods (with graph algorithms, linear programming, polyhedral optimizations) to address code analysis and optimization problems (e.g., termination, register allocation, memory optimizations, scheduling, automatic generation of interfaces) and the validation of these techniques through the development of compilation tools.

Compsys started as an Inria project in 2004, after 2 years of maturation. This first period of Compsys, Compsys I, was positively evaluated in Spring 2007 after its first 4 years period (2004-2007). It was again evaluated by AERES in 2009, as part of the general evaluation of Lip, and got the best possible mark, A+. The second period (2007-2012), Compsys II, was again evaluated positively by Inria in Spring 2012 and formally prolonged into Compsys III at the very end of 2012. In 2013, Fabrice Rastello moved to Grenoble first to expand the activities of Compsys in the context of Giant, a R&D technology center with several industrial and academic actors. He left officially the team in 2014 to work on his own. The research directions of Compsys III followed the lines presented in the synthesis report provided for the 2012 evaluation [1], including a shift towards the compilation of streaming programming, the analysis and optimizations of parallel languages, and an even stronger focus on polyhedral optimizations and their extensions. Christophe Alias was mostly involved in the development of the Zettice/XTREMLOGIC start-up. The hiring of Laure Gonnord (in 2013) and Tomofumi Yuki (in 2014) added new forces on the code analysis research aspects and on HPC polyhedral-related topics. However, Christophe Alias and Laure Gonnord left the team in Sep. 2015. The project-team itself ended officially in Dec. 2015, reaching the limit of 12 years. Nevertheless, it will be evaluated again by Inria in Spring 2016. It has been extended as an Inria team until Dec. 2016.

Section 2.2 defines the general context of the team's activities. Section 2.3 presents the research objectives and main achievements in Compsys I, i.e., until 2007, and how its research directions were modified for Compsys II. Section 2.4 briefly presents the main achievements of Compsys II and of the first years of Compsys III, referring to the annual reports from 2008 to 2014 for details. As for the highlights of the past year, i.e., 2015, they are given in Section 5.1.

## 2.2. General Presentation

Classically, an embedded computer is a digital system that is part of a larger system and that is not directly accessible to the user. Examples are appliances like phones, TV sets, washing machines, game platforms, or even larger systems like radars and sonars. In particular, this computer is not programmable in the usual way. Its program, if it exists, is supplied as part of the manufacturing process and is seldom (or ever) modified thereafter. As the embedded systems market grows and evolves, this view of embedded systems is becoming obsolete and tends to be too restrictive. Many aspects of general-purpose computers apply to modern embedded platforms. Nevertheless, embedded systems remain characterized by a set of specialized application domains, rigid constraints (cost, power, efficiency, heterogeneity), and its market structure. The term *embedded system* has been used for naming a wide variety of objects. More precisely, there are two categories of so-called *embedded systems*: a) control-oriented and hard real-time embedded systems (automotive, plant control,

---

[1] See http://www.ens-lyon.fr/LIP/COMPSYS/wordpress/wp-content/uploads/2013/09/ficheSynthese.pdf

airplanes, etc.); b) compute-intensive embedded systems (signal processing, multi-media, stream processing) processing large data sets with parallel and/or pipelined execution. Compsys is primarily concerned with this second type of embedded systems, now referred to as *embedded computing systems*.

Today, the industry sells many more embedded processors than general-purpose processors; the field of embedded systems is one of the few segments of the computer market where the European industry still has a substantial share, hence the importance of embedded system research in the European research initiatives. Our priority towards embedded software is motivated by the following observations: a) the embedded system market is expanding, among many factors, one can quote pervasive digitalization, low-cost products, appliances, etc.; b) research on software for embedded systems is poorly developed in France, especially if one considers the importance of actors like Alcatel, STMicroelectronics, Matra, Thales, etc.; c) since embedded systems increase in complexity, new problems are emerging: computer-aided design, shorter time-to-market, better reliability, modular design, and component reuse.

A specific aspect of embedded computing systems is the use of various kinds of processors, with many particularities (instruction sets, registers, data and instruction caches, now multiple cores) and constraints (code size, performance, storage). The development of *compilers* is crucial for this industry, as selling a platform without its programming environment and compiler would not be acceptable. To cope with such a range of different processors, the development of robust, generic (retargetable), though efficient compilers is mandatory. Unlike standard compilers for general-purpose processors, compilers for embedded processors and hardware accelerators can be more aggressive (i.e., take more time to optimize) for optimizing some important parts of applications. This opens a new range of optimizations. Another interesting aspect is the introduction of platform-independent intermediate languages, such as Java bytecode, that is compiled dynamically at runtime (aka just-in-time). Extreme lightweight compilation mechanisms that run faster and consume less memory have to be developed. The introduction of intermediate languages such as OpenCL is also a sign of the need for portability (as well as productivity) across diverse (if not heterogeneous) platforms. One of the initial objectives of Compsys was thus to revisit existing compilation techniques in the context of such embedded computing systems, to deconstruct some of these techniques, to improve them, and to develop new techniques taking constraints of embedded processors and platforms into account.

As for *high-level synthesis* (HLS), several compilers/systems have appeared, after some first unsuccessful industrial attempts in the past. These tools are mostly based on C or C++ as for example SystemC, VCC, CatapultC, Altera C2H, Pico-Express, Vivado HLS. Academic projects also exist (or existed) such as Flex and Raw at MIT, Piperench at Carnegie-Mellon University, Compaan at the University of Leiden, Ugh/Disydent at LIP6 (Paris), Gaut at Lester (Bretagne), MMAlpha (Insa-Lyon), and others. In general, the support for parallelism in HLS tools is minimal, especially in industrial tools. Also, the basic problem that these projects have to face is that the definition of performance is more complex than in classical systems. In fact, it is a multi-criteria optimization problem and one has to take into account the execution time, the size of the program, the size of the data structures, the power consumption, the manufacturing cost, etc. The impact of the compiler on these costs is difficult to assess and control. Success will be the consequence of a detailed knowledge of all steps of the design process, from a high-level specification to the chip layout. A strong cooperation of the compilation and chip design communities is needed. The main expertise in Compsys for this aspect is in the *parallelization* and optimization of *regular computations*. Hence, we target applications with a large potential parallelism, but we attempt to integrate our solutions into the big picture of CAD environments.

More generally, the aims of Compsys are to develop new compilation and optimization techniques for the field of embedded computing system design. This field is large, and Compsys does not intend to cover it in its entirety. As previously mentioned, we are mostly interested in the automatic design of accelerators, for example designing a VLSI or FPGA circuit for a digital filter, and in the development of new back-end compilation strategies for embedded processors. We study code transformations that optimize features such as execution time, power consumption, code and die size, memory constraints, and compiler reliability. These features are related to embedded systems but some are not specific to them. The code transformations we develop are both at source level and at assembly level. A specificity of Compsys is to mix a solid theoretical basis for all code optimizations we introduce with algorithmic/software developments. Within Inria, our

project is related to the "architecture and compilation" theme, more precisely code optimization, as some of the research conducted in Parkas (previously known as Alchemy), Alf (previously known as Caps), Camus, and to high-level architectural synthesis, as some of the research in Cairn.

Most french researchers working on high-performance computing (automatic parallelization, languages, operating systems, networks) moved to grid computing at the end of the 90s. We thought that applications, industrial needs, and research problems were more interesting in the design of embedded platforms. Furthermore, we were convinced that our expertise on high-level code transformations could be more useful in this field. This is the reason why Tanguy Risset came to Lyon in 2002 to create the Compsys team with Anne Mignotte and Alain Darte, before Paul Feautrier, Antoine Fraboulet, Fabrice Rastello, and finally Christophe Alias joined the group. Then, Tanguy Risset left Compsys to become a professor at INSA Lyon, and Antoine Fraboulet and Anne Mignotte moved to other fields of research. As for Laure Gonnord, after a post-doc in Compsys, she obtained an assistant professor position in Lille but remained external collaborator of the team for the period 2009-2013 and finally obtained an assistant professor position in Lyon, and integrated officially the team. About the same time, Fabrice Rastello left while Tomofumi Yuki was hired as Inria researcher in 2014.

All present and past members of Compsys have a background in automatic parallelization and high-level program analyses and transformations. Paul Feautrier was the initiator of the polyhedral model for program transformations around 1990 and, before coming to Lyon, started to be more interested in programming models and optimizations for embedded applications, in particular through collaborations with Philips. Alain Darte worked on mathematical tools and algorithmic issues for parallelism extraction in programs. He became interested in the automatic generation of hardware accelerators, thanks to his stay at HP Labs in the Pico project in 2001. Antoine Fraboulet did a PhD with Anne Mignotte – who was working on high-level synthesis (HLS) – on code and memory optimizations for embedded applications. Fabrice Rastello did a PhD on tiling transformations for parallel machines, then was hired by STMicroelectronics where he worked on assembly code optimizations for embedded processors. Tanguy Risset worked for a long time on the synthesis of systolic arrays, being the main architect of the HLS tool MMAlpha. Christophe Alias did a PhD on algorithm recognition for program optimizations and parallelization. He first spent a year in Compsys working on array contraction, where he started to develop the tool Bee, then a year at Ohio State University with Prof. P. Sadayappan on memory optimizations. Laure Gonnord did a PhD on invariant generation and program analysis and became interested on application on compilation and code generation since her postdoc in the team. Finally, Tomofumi Yuki did a PhD on polyhedral programming environments and optimizations (in Colorado State University, with Prof. S. Rajopadhye) before a post-doc on polyhedral HLS in the Cairn team (Rennes).

To understand why we think automation in our field is highly important, it may be worth to quote Bob Rau and his colleagues (IEEE Computer, Sep. 2002):

*"Engineering disciplines tend to go through fairly predictable phases: ad hoc, formal and rigorous, and automation. When the discipline is in its infancy and designers do not yet fully understand its potential problems and solutions, a rich diversity of poorly understood design techniques tends to flourish. As understanding grows, designers sacrifice the flexibility of wild and woolly design for more stylized and restrictive methodologies that have underpinnings in formalism and rigorous theory. Once the formalism and theory mature, the designers can automate the design process. This life cycle has played itself out in disciplines as diverse as PC board and chip layout and routing, machine language parsing, and logic synthesis.*

*We believe that the computer architecture discipline is ready to enter the automation phase. Although the gratification of inventing brave new architectures will always tempt us, for the most part the focus will shift to the automatic and speedy design of highly customized computer systems using well-understood architecture and compiler technologies."*

We share this view of the future of architecture and compilation. Without targeting too ambitious objectives, we were convinced of two complementary facts: a) the mathematical tools developed in the past for manipulating programs in automatic parallelization were lacking in high-level synthesis and embedded computing optimizations and, even more, they started to be rediscovered frequently in less mature forms, b) before being

able to really use these techniques in HLS and embedded program optimizations, we needed to learn a lot from the application side, from the electrical engineering side, and from the embedded architecture side. Our primary goal was thus twofold: to increase our knowledge of embedded computing systems and to adapt/extend code optimization techniques, primarily designed for high performance computing, to the special case of embedded computing systems. In the initial Compsys proposal, we proposed four research directions, centered on compilation methods for embedded applications, both for software and accelerators design:

- Code optimization for specific processors (mainly DSP and VLIW processors);
- Platform-independent loop transformations (including memory optimization);
- Silicon compilation and hardware/software codesign;
- Development of polyhedral (but not only) optimization tools.

These research activities were primarily supported by a marked investment in polyhedra manipulation tools and, more generally, solid mathematical and algorithmic studies, with the aim of constructing operational software tools, not just theoretical results. Hence the fourth research theme was centered on the development of these tools.

## 2.3. Summary of Compsys I Achievements

The Compsys team has been evaluated by Inria for the first time in April 2007. The evaluation, conducted by Erik Hagersted (Uppsala University), Vinod Kathail (Synfora, inc), J. (Ram) Ramanujam (Baton Rouge University) was positive. Compsys I thus continued into Compsys II for 4-5 years but in a new configuration as Tanguy Risset and Antoine Fraboulet left the project to follow research directions closer to their host laboratory at Insa-Lyon. The main achievements of Compsys I, for this period, were the following:

- The development of a strong collaboration with the compilation group at STMicroelectronics, with important results in aggressive optimizations for instruction cache and register allocation.
- New results on the foundation of high-level program transformations, including scheduling techniques for process networks and a general technique for array contraction (memory reuse) based on the theory of lattices.
- Many original contributions with partners closer to hardware constraints, including CEA, related to SoC simulation, hardware/software interfaces, power models, and simulators.

Due to Compsys size reduction (from 5 permanent researchers to 3 in 2008, then 4 again in 2009), the team then focused, in Compsys II, on two research directions only:

- Code generation for embedded processors, on the two opposite, though connected, aspects: aggressive compilation and just-in-time compilation.
- High-level program analysis and transformations for high-level synthesis tools.

## 2.4. Quick View of Compsys II Achievements and Directions for Compsys III

The main achievements of Compsys II were:

- the great success of the collaboration with STMicroelectronics with many deep results on SSA (Static Single Assignment), register allocation, and intermediate program representations;
- the design of high-level program analysis, optimizations, and tools, mainly related to high-level synthesis, some leading to the development of the Zettice start-up.

For more details on the past years of Compsys II, see the previous annual reports from 2008 to 2012. Compsys II was positively evaluated in Spring 2012 by Inria. The evaluation committee members were Walid Najjar (University of California Riverside), Paolo Faraboschi (HP Labs), Scott Mahlke (University of Michigan), Pedro Diniz (University of Southern California), Peter Marwedel (TU Dortmund), and Pierre Paulin (STMicroelectronics, Canada), the last three assigned specifically to Compsys.

For Compsys III, the changes in the permanent members (departure of Fabrice Rastello and arrival of Laure Gonnord, while she was only external collaborator of Compsys until Sep. 2013) reduced the forces on back-end code optimizations, and in particular dynamic compilation, but increased the forces on program analysis. In this context, Compsys III has continued to develop fundamental concepts or techniques whose applicability should go beyond a particular architectural or language trend, as well as stand-alone tools (either as proofs of concepts or to be used as basic blocks in larger tools/compilers developed by others) and our own experimental prototypes. One of the main objectives of Compsys III has been to try to push the polyhedral model beyond its present limits both in terms of analysis techniques (possibly integrating approximation and runtime support) and of applicability (e.g., analysis of parallel or streaming languages, program verification, compilation towards accelerators such as GPU or multicores). The hiring of Tomofumi Yuki supported this new direction. A summary of 2013 and 2014 activities are given in the corresponding annual reports, while new results for 2015 are provided in this document, in Section 5.1 (highlights) and from Section 7.1 to 7.15 (new results).

# 3. Research Program

## 3.1. Architecture and Compilation Trends

The embedded system design community is facing two challenges:

- The complexity of embedded applications is increasing at a rapid rate.
- The needed increase in processing power is no longer obtained by increases in the clock frequency, but by increased parallelism.

While, in the past, each type of embedded application was implemented in a separate appliance, the present tendency is toward a universal hand-held object, which must serve as a cell-phone, as a personal digital assistant, as a game console, as a camera, as a Web access point, and much more. One may say that embedded applications are of the same level of complexity as those running on a PC, but they must use a more constrained platform in terms of processing power, memory size, and energy consumption. Furthermore, most of them depend on international standards (e.g., in the field of radio digital communication), which are evolving rapidly. Lastly, since ease of use is at a premium for portable devices, these applications must be integrated seamlessly to a degree that is unheard of in standard computers.

All of this dictates that modern embedded systems retain some form of programmability. For increased designer productivity and reduced time-to-market, programming must be done in some high-level language, with appropriate tools for compilation, run-time support, and debugging. This does not mean however that all embedded systems (or all of an embedded system) must be processor based. Another solution is the use of field programmable gate arrays (FPGA), which may be programmed at a much finer grain than a processor, although the process of FPGA "programming" is less well understood than software generation. Processors are better than application-specific circuits at handling complicated control and unexpected events. On the other hand, FPGAs may be tailored to just meet the needs of their application, resulting in better energy and silicon area usage. It is expected that most embedded systems will use a combination of general-purpose processors, specific processors like DSPs, and FPGA accelerators (or even low-power GPUs). Such a combination DSP+FPGA is already present in recent versions of the Atom Intel processor.

As a consequence, parallel programming, which has long been confined to the high-performance community, must become the common place rather than the exception. In the same way that sequential programming moved from assembly code to high-level languages at the price of a slight loss in performance, parallel programming must move from low-level tools, like OpenMP or even MPI, to higher-level programming environments. While fully-automatic parallelization is a Holy Grail that will probably never be reached in our lifetimes, it will remain as a component in a comprehensive environment, including general-purpose parallel programming languages, domain-specific parallelizers, parallel libraries and run-time systems, back-end compilation, dynamic parallelization. The landscape of embedded systems is indeed very diverse and

many design flows and code optimization techniques must be considered. For example, embedded processors (micro-controllers, DSP, VLIW) require powerful back-end optimizations that can take into account hardware specificities, such as special instructions and particular organizations of registers and memories. FPGA and hardware accelerators, to be used as small components in a larger embedded platform, require "hardware compilation", i.e., design flows and code generation mechanisms to generate non-programmable circuits. For the design of a complete system-on-chip platform, architecture models, simulators, debuggers are required. The same is true for multicores of any kind, GPGPU ("general-purpose" graphical processing units), CGRA (coarse-grain reconfigurable architectures), which require specific methodologies and optimizations, although all these techniques converge or have connections. In other words, embedded systems need all usual aspects of the process that transforms some specification down to an executable, software or hardware. In this wide range of topics, Compsys concentrates on the code optimizations aspects (and the associated analysis) in this transformation chain, restricting to compilation (transforming a program to a program) for embedded processors and programmable accelerators, and to high-level synthesis (transforming a program into a circuit description) for FPGAs.

Actually, it is not a surprise to see compilation and high-level synthesis getting closer (in the last 10 years now). Now that high-level synthesis has grown up sufficiently to be able to rely on place-and-route tools, or even to synthesize C-like languages, standard techniques for back-end code generation (register allocation, instruction selection, instruction scheduling, software pipelining) are used in HLS tools. At the higher level, programming languages for programmable parallel platforms share many aspects with high-level specification languages for HLS, for example, the description and manipulations of nested loops, or the model of computation/communication (e.g., Kahn process networks and its many "streaming" variants). In all aspects, the frontier between software and hardware is vanishing. For example, in terms of architecture, customized processors (with processor extension as first proposed by Tensilica) share features with both general-purpose processors and hardware accelerators. FPGAs are both hardware and software as they are fed with "programs" representing their hardware configurations.

In other words, this convergence in code optimizations explains why Compsys studies both program compilation and high-level synthesis, and at both front-end and back-end levels, the first one acting more at the granularity of memories, transfers, and multiple cores, the second one more at the granularity of registers, system calls, and single core. Both levels must be considered as they interact with each other. Front-end optimizations must be aware of what back-end optimizations will do, as single core performance remain the basis for good parallel performances. Some front-end optimizations even act directly on back-end features, for example register tiling considered as a source-level transformation. Also, from a conceptual point of view, the polyhedral techniques developed by Compsys are actually the symbolic front-end counterpart, for structured loops, of back-end analysis and optimizations of unstructured programs (through control-flow graphs), such as dependence analysis, scheduling, lifetime analysis, register allocation, etc. A strength of Compsys so far was to juggle with both aspects, one more on graph theory with SSA-type optimizations, the other with polyhedra representing loops, and to exploit the correspondence between both. This has still to be exploited, for applying polyhedral techniques to more irregular programs.

Besides, Compsys has a tradition of building free software tools for linear programming and optimization in general, and will continue it, as needed for our current research.

### 3.1.1. Compilation and Languages Issues in the Context of Embedded Processors, "Embedded Systems", and Programmable Accelerators

Compilation is an old activity, in particular back-end code optimizations. The development of embedded systems was one of the reasons for the revival of compilation activities as a research topic. Applications for embedded computing systems generate complex programs and need more and more processing power. This evolution is driven, among others, by the increasing impact of digital television, the first instances of UMTS networks, and the increasing size of digital supports, like recordable DVD, and even Internet applications. Furthermore, standards are evolving very rapidly (see for instance the successive versions of MPEG). As a consequence, the industry has focused on programmable structures, whose flexibility more

than compensates for their larger size and power consumption. The appliance provider has a choice between hard-wired structures (Asic), special-purpose processors (Asip), (quasi) general-purpose processors (DSP for multimedia applications), and now hardware accelerators (dedicated platforms – such as those developed by Thales or the CEA –, or more general-purpose accelerators such as GPUs or even multicores, even if these are closer to small HPC platforms than truly embedded systems). Our cooperation with STMicroelectronics, until 2012, focused on investigating the compilation for specialized processors, such as the ST100 (DSP processor) and the ST200 (VLIW DSP processor) family. Even for this restricted class of processors, the diversity is large, and the potential for instruction level parallelism (SIMD, MMX), the limited number of registers and the small size of the memory, the use of direct-mapped instruction caches, of predication, generate many open problems. Our goal was to contribute to their understanding and their solutions.

An important concept to cope with the diversity of platforms is the concept of *virtualization*, which is a key for more portability, more simplicity, more reliability, and of course more security. This concept – implemented at low level through binary translation and just-in-time (JIT) compilation [2] – consists in hiding the architecture-dependent features as long as possible during the compilation process. It has been used for a while for servers such as HotSpot, a bit more recently for workstations, and now for embedded computing. The same needs drive the development of intermediate languages such as OpenCL to, not necessarily hide, but at least make more uniform, the different facets of the underlying architectures. The challenge is then to design and compile high-productivity and high-performance languages [3] (coping with parallelism and heterogeneity) that can be ported to such intermediate languages, or to architecture-dependent runtime systems. The offloading of computation kernels, through source-to-source compilation, targeting back-end C dialects, has the same goals: to automate application porting to the variety of accelerators.

For JIT compilation, the compactness of the information representation, and thus its pertinence, is an important criterion for such late compilation phases. Indeed, the intermediate representation (IR) is evolving not only from a target-independent description to a target-dependent one, but also from a situation where the compilation time is almost unlimited (cross-compilation) to one where any type of resource is limited. This is one of the reasons why static single assignment (SSA), a sparse compact representation of liveness information, became popular in embedded compilation. If time constraints are common to all JIT compilers (not only for embedded computing), the benefit of using SSA is also in terms of its good ratio pertinence/storage of information. It also enables to simplify algorithms, which is also important for increasing the reliability of the compiler. In this context, our aim has been, in particular, to develop exact or heuristic solutions to *combinatorial* problems that arise in compilation for VLIW and DSP processors, and to integrate these methods into industrial compilers for DSP processors (mainly ST100, ST200, Strong ARM). Such combinatorial problems can be found in register allocation, opcode selection, code placement, when removing the SSA multiplexer functions (known as $\phi$ functions). These optimizations are usually done in the last phases of the compiler, using an assembly-level intermediate representation. As mentioned in Sections 2.3 and 2.4, we made a lot of progress in this area in our past collaborations with STMicroelectronics (see also previous activity reports). Through the Sceptre and Mediacom projects, we first revisited, in the light of SSA, some code optimizations in an aggressive context, to develop better strategies, without eliminating too quickly solutions that may have been considered as too expensive in the past. Then we exploited the new concepts introduced in the aggressive context to design better algorithms in a JIT context, focusing on the speed of algorithms and their memory footprint, without compromising too much on the quality of the generated code.

---

[2]*Aggressive compilation* consists in allowing more time to implement more complete and costly solutions: the compiled program is loaded in permanent memory (ROM, flash, etc.) and its compilation time is less relevant than the execution time, size, and energy consumption of the produced code, which can have a critical impact on the cost and quality of the final product. Hence, the application is cross-compiled, i.e., compiled on a powerful platform distinct from the target processor. *Just-in-time compilation*, on the other hand, corresponds to compiling applets on demand on the target processor. For compatibility and compactness, the source languages are CIL or Java bytecode. The code can be uploaded or sold separately on a flash memory. Compilation is performed at load time and even dynamically during execution. The optimization heuristics, constrained by time and limited resources, are far from being aggressive. They must be fast but smart enough.

[3]For examples of such languages, see the keynotes event we organized in 2013: http://labexcompilation.ens-lyon.fr/hpc-languages.

Our research directions are currently more focused on programmable accelerators, such as GPU and multi-cores, but still considering *static* compilation and without forgetting the link between high-level (in general at source-code level) and low-level (i.e., at assembly-code level) optimizations. They concern program analysis (of both sequential and parallel specifications), program optimizations (for memory hierarchies, parallelism, streaming, etc.), and also the link with applications and between compilers and users (programmers). Polyhedral techniques play an important role in these directions, even if control-flow-based techniques remain in the background and may come back at any time in the foreground. This is also the case for high-level synthesis, as exposed in the next section.

### 3.1.2. *Context of High-Level Synthesis and FPGA Platforms*

High-level synthesis has become a necessity, mainly because the exponential increase in the number of gates per chip far outstrips the productivity of human designers. Besides, applications that need hardware accelerators usually belong to domains, like telecommunications and game platforms, where fast turn-around and time-to-market minimization are paramount. When Compsys started, we were convinced that our expertise in compilation and automatic parallelization could contribute to the development of the needed tools.

Today, synthesis tools for FPGAs or ASICs come in many shapes. At the lowest level, there are proprietary Boolean, layout, and place-and-route tools, whose input is a VHDL or Verilog specification at the structural or register-transfer level (RTL). Direct use of these tools is difficult, for several reasons:

- A structural description is completely different from an usual algorithmic language description, as it is written in term of interconnected basic operators. One may say that it has a spatial orientation, in place of the familiar temporal orientation of algorithmic languages.

- The basic operators are extracted from a library, which poses problems of selection, similar to the instruction selection problem in ordinary compilation.

- Since there is no accepted standard for VHDL synthesis, each tool has its own idiosyncrasies and reports its results in a different format. This makes it difficult to build portable HLS tools.

- HLS tools have trouble handling loops. This is particularly true for logic synthesis systems, where loops are systematically unrolled (or considered as sequential) before synthesis. An efficient treatment of loops needs the polyhedral model. This is where past results from the automatic parallelization community are useful.

- More generally, a VHDL specification is too low level to allow the designer to perform, easily, higher-level code optimizations, especially on multi-dimensional loops and arrays, which are of paramount importance to exploit parallelism, pipelining, and perform communication and memory optimizations.

Some intermediate tools were proposed that generate VHDL from a specification in restricted C, both in academia (such as SPARK, Gaut, UGH, CloogVHDL), and in industry (such as C2H, CatapultC, Pico-Express, Vivado HLS). All these tools use only the most elementary form of parallelization, equivalent to instruction-level parallelism in ordinary compilers, with some limited form of block pipelining, and communication through FIFOs. Targeting one of these tools for low-level code generation, while we concentrate on exploiting loop parallelism, might be a more fruitful approach than directly generating VHDL. However, it may be that the restrictions they impose preclude efficient use of the underlying hardware. Our first experiments with these HLS tools reveal two important issues. First, they are, of course, limited to certain types of input programs so as to make their design flows successful, even if, over the years, they become more and more mature. But it remains a painful and tricky task for the user to transform the program so that it fits these constraints and to tune it to get good results. Automatic or semi-automatic program transformations can help the user achieve this task. Second, users, even expert users, have only a very limited understanding of what back-end compilers do and why they do not lead to the expected results. An effort must be done to analyze the different design flows of HLS tools, to explain what to expect from them, and how to use them to get a good quality of results. Our first goal is thus to develop high-level techniques that, used in front of existing HLS tools, improve their utilization. This should also give us directions on how to modify them or to design new tools from scratch.

More generally, we want to consider HLS as a more global parallelization process. So far, no HLS tools is capable of generating designs with communicating *parallel* accelerators, even if, in theory, at least for the scheduling part, a tool such as Pico-Express could have such capabilities. The reason is that it is, for example, very hard to automatically design parallel memories and to decide the distribution of array elements in memory banks to get the desired performances with parallel accesses. Also, how to express communicating processes at the language level? How to express constraints, pipeline behavior, communication media, etc.? To better exploit parallelism, a first solution is to extend the source language with parallel constructs, as in all derivations of the Kahn process networks model, including communicating regular processes (CRP, see later). The other solution is a form of automatic parallelization. However, classical methods, which are mostly based on scheduling, need to be revisited, to pay more attention to locality, process streaming, and low-level pipelining, which are of paramount importance in hardware. Besides, classical methods mostly rely on the runtime system to tailor the parallelism degree to the available resources. Obviously, there is no runtime system in hardware. The real challenge is thus to invent new scheduling algorithms that take resource, locality, and pipelining into account, and then to infer the necessary hardware from the schedule. This is probably possible only for programs that fit into the polyhedral model, or in an incrementally-extended model.

Our research activities on polyhedral code analysis and optimizations directly target these HLS challenges. But they are not limited to the automatic generation of hardware as can be seen from our different contributions on X10, OpenStream, parametric tiling, etc. The same underlying concepts also arise when optimizing codes for GPUs and multicores. In this context of polyhedral analysis and optimizations, we will focus on three aspects:

- developing high-level transformations, especially for loops and memory/communication optimizations, that can be used in front of HLS tools so as to improve their use, as well as for hardware accelerators;

- developing concepts and techniques in a more global view of high-level synthesis and high-level parallel programming, starting from specification languages down to hardware implementation;

- developing more general code analysis so as to extract more information from codes as well as to extend the programs that can be handled.

## 3.2. Code Analysis, Code Transformations, Code Optimizations

Embedded systems generated new problems in code analysis and optimization both for optimizing embedded software (compilation) and hardware (HLS). We now give a bit more details on some general challenges for program analysis, optimizations, and transformations, induced by this context, and on our methodology, in particular our development and use of polyhedral optimizations and its extensions.

### 3.2.1. Processes, Scheduling, Mapping, Communications, etc.

Before mapping an application to an architecture, one has to decide which execution model is targeted and where to intervene in the design flow. Then one has to solve scheduling, placement, and memory management problems. These three aspects should be handled as a whole, but present state of the art dictates that they be treated separately. One of our aims will be to find more comprehensive solutions. The last task is code generation, both for the processing elements and the interfaces processors/accelerators.

There are basically two execution models for embedded systems: one is the classical accelerator model, in which data is deposited in the memory of the accelerator, which then does its job, and returns the results. In the streaming model, computations are done on the fly, as data items flow from an input channel to the output. Here, the data are never stored in (addressable) memory. Other models are special cases, or sometimes compositions of the basic models. For instance, a systolic array follows the streaming model, and sometimes extends it to higher dimensions. Software radio modems follow the streaming model in the large, and the accelerator model in detail. The use of first-in first-out queues (FIFO) in hardware design is an application of the streaming model. Experience shows that designs based on the streaming model are more efficient that those based on memory, for such applications. One of the point to be investigated is whether it is general enough to handle arbitrary (regular) programs. The answer is probably negative. One possible implementation of the streaming model is as a network of communicating processes either as Kahn process networks (FIFO

based) or as our more recent model of communicating regular processes (memory based, see for example CRP below). It is an interesting fact that several researchers have investigated translation from process networks [21] and to process networks [32], [33]. Streaming languages such as StreamIt and OpenStream have also been developed.

Kahn process networks (KPN) were introduced 30 years ago as a notation for representing parallel programs. Such a network is built from processes that communicate via perfect FIFO channels. Because the channel histories are deterministic, one can define a semantics and talk meaningfully about the equivalence of two implementations. As a bonus, the dataflow diagrams used by signal processing specialists can be translated on-the-fly into process networks. The problem with KPNs is that they rely on an asynchronous execution model, while VLIW processors and FPGAs are synchronous or partially synchronous. Thus, there is a need for a tool for synchronizing KPNs. This can be done by computing a schedule that has to satisfy data dependences within each process, a causality condition for each channel (a message cannot be received before it is sent), and real-time constraints. However, there is a difficulty in writing the channel constraints because one has to count messages in order to establish the send/receive correspondence and, in multi-dimensional loop nests, the counting functions may not be affine. Recent developments on the theory of polynomials (see Section 7.11) may offer a solution to this problem. One can also define another model, *communicating regular processes* (CRP), in which channels are represented as write-once/read-many arrays. One can then dispense with counting functions and prove that the determinacy property still holds. As an added benefit, a communication system in which the receive operation is not destructive is closer to the expectations of system designers.

The main difficulty with this approach is that ordinary programs are usually not constructed as process networks. One needs automatic or semi-automatic tools for converting sequential programs into process networks. One possibility is to start from array dataflow analysis [23] or variants. Another approach attempts to construct threads, i.e., pieces of sequential code with the smallest possible interactions. In favorable cases, one may even find outermost parallelism, i.e., threads with no interactions whatsoever. Tiling mechanisms can also be used to define atomic processes that can be pipelined as we proposed initially for FPGA [17].

Whatever the chosen solution (FIFO or addressable memory) for communicating between two accelerators or between the host processor and an accelerator, the problems of optimizing communication between processes and of optimizing buffers have to be addressed. Many local memory optimization problems have already been solved theoretically. Some examples are loop fusion and loop alignment for array contraction, techniques for data allocation in scratch-pad memory, or techniques for folding multi-dimensional arrays [20]. Nevertheless, the problem is still largely open. Some questions are: how to schedule a loop sequence (or even a process network) for minimal scratch-pad memory size? How is the problem modified when one introduces unlimited and/or bounded parallelism (same questions for analyzing explicitly-parallel programs)? How does one take into account latency or throughput constraints, bandwidth constraints for input and output channels, memory hierarchies? All loop transformations are useful in this context, in particular loop tiling, and may be applied either as source-to-source transformations (when used in front of HLS or C-level compilers) or to generate directly VHDL or lower-level C-dialects such as OpenCL. One should keep in mind that theory will not be sufficient to solve these problems. Experiments are required to check the relevance of the various models (computation model, memory model, power consumption model) and to select the most important factors according to the architecture. Besides, optimizations do interact: for instance, reducing memory size and increasing parallelism are often antagonistic. Experiments will be needed to find a global compromise between local optimizations. In particular, the design of cost models remain a fundamental challenge.

Finally, there remains the problem of code generation for accelerators. It is a well-known fact that modern methods for program optimization and parallelization do not generate a new program, but just deliver blueprints for program generation, in the form, e.g., of schedules, placement functions, or new array subscripting functions. A separate code generation phase must be crafted with care, as a too naive implementation may destroy the benefits of high-level optimization. There are two possibilities here as suggested before; one may target another high-level synthesis or compilation tool, or one may target directly VHDL or low-level code. Each approach has its advantages and drawbacks. However, both situations require that the input program re-

spects some strong constraints on the code shape, array accesses, memory accesses, communication protocols, etc. Furthermore, to get the compilers do what the user wants requires a lot of program tuning, i.e., of program rewriting or of program annotations. What can be automated in this rewriting process? Semi-automated?

In other words, we still need to address scheduling, memory, communication, and code generation issues, in the light of the developments of new languages and architectures, pushing the limits of such an automation.

### 3.2.2. *Beyond Static Control Programs*

With the advent of parallelism in supercomputers, the bulk of research in code transformation resulted in (semi-)automatic parallelization, with many techniques (analysis, scheduling, code generation, etc.) based on the description and manipulation of nested loops with polyhedra. Compsys has always taken an active part in the development of these so-called "polyhedral techniques". Historically, these analysis were (wrongly) understood to be limited to static control programs.

Actually, the polyhedral model is neither a programming language nor an execution model rather an intermediate representation. As such, it can be generated from imperative sequential languages like C or Fortran, streaming languages like CRP, or equational languages like Alpha. While the structure of the model is the same in all three cases, it may enjoy different properties, e.g., a schedule always exists in the first case, not in the two others. The import of the polyhedral model is that many questions relative to the analysis of a program and the applicability of transformations can be answered precisely and efficiently by applying well-known mathematical results to the model.

For irregular programs, the basic idea is to construct a polyhedral over-approximation, i.e., a program which has more operations, a larger memory footprint, and more dependences than the original. One can then parallelize the approximated program using polyhedral tools, and then return to the original, either by introducing guards, or by insuring that approximations are harmless. This technique is the standard way of dealing with approximated dependences. We already started to study the impact of approximations in our kernel offloading technique, for optimizing remote communications [3]. It is clear however that this method will apply only to mildly non-polyhedral programs. The restriction to arrays as the only data structure is still present. Its advantage is that it will be able to subsume in a coherent framework many disparate tricks: the extraction of SCoPs, induction variable detection, the omission of non-affine subscripts, or the conversion of control dependences into data dependences. The link with the techniques developed in the PIPS compiler (based on array region analysis) is strong and will have to be explored.

Such over-approximations can be found by mean of abstract interpretation, a general framework to develop static analysis on real-life programs. However, they were designed mainly for verification purposes, thus precision was the main issue before scalability. Although many efforts were made in designing specialized analyses (pointers, data structures, arrays), these approaches still suffer from a lack of experimental evidence concerning their applicability for code optimization. Following our experience and work on termination analysis (that connects the work on back-end CFG-like and front-end polyhedral-like optimizations), and our work on range analysis of numerical variables and on the memory footprint on real-world C programs [29], our objective is to bridge the gap between abstract interpretation and compilation, by designing cheaper analyses that scale well, mainly based on compact representations derived from variants of static single assignment (SSA). We will focus on complex control, and complex data structures (pointers, lists) that still suffer from complexity issues in the area of optimization.

Another possibility is to rely on application specific knowledge to guide compiler decisions, as it is impossible for a compiler alone to fully exploit such pieces of information. A possible approach to better utilize such knowledge is to put the programmers "in the loop". Expert parallel programmers often have a good idea about coarse-grain parallelism and locality that they want to use for an application. On the other hand, fine-grain parallelism (e.g., ILP, SIMD) is tedious and specific to each underlying architecture, and is best left to the compiler. Furthermore, approximations will have opportunities to be refined using programmer knowledge. The key challenge is to create a programming environment where compiler techniques and programmer knowledge can be combined effectively. One of the difficulties is to design a common language between the compiler and the programmer.

# 3.3. Mathematical Tools

All compilers have to deal with *sets* and relations. In classical compilers, these sets are finite: the set of statements of a program, the set of its variables, its abstract syntax tree (AST), its control-flow graph (CFG), and many others. It is only in the first phase of compilation, parsing, that one has to deal with infinite objects, regular and context-free languages, and those are represented by finite grammars, and are processed by a symbolic algorithm, `yacc` or one of its clones.

When tackling parallel programs and parallel compilation, it was soon realized that this position was no longer tenable. Since it makes no sense to ask whether a statement can be executed in parallel with itself, one has to consider sets of operations, which may be so large as to forbid an extensive representation, or even be infinite. The same is true for dependence sets, for memory cells, for communication sets, and for many other objects a parallel compiler has to consider. The representation is to be *symbolic*, and all necessary algorithms have to be promoted to symbolic versions.

Such symbolic representations have to be efficient – the formula representing a set has to be much smaller than the set itself – and effective – the operations one needs, union, intersection, emptiness tests and many others – have to be feasible and fast. As an aside, note that progress in algorithm design has blurred the distinction between polynomially-solvable and NP-complete problems, and between decidable and undecidable questions. For instance SAT, SMT, and ILP software tools solve efficiently many NP-complete problems, and the Z3 tool is able to "solve" many instances of the undecidable Hilbert's 10th problem.

Since the times of Pip and of the Polylib, Compsys has been active in the implementation of basic mathematical tools for program analysis and synthesis. Pip is still developed by Paul Feautrier and Cédric Bastoul, while the Polylib is now taken care of by the Inria Camus project, which introduced Ehrhart polynomials. These tools are still in use world-wide and they also have been reimplemented many times with (sometimes slight) improvements, e.g., as part of the Parma Polylib, of Sven Verdoolaege's Isl and Barvinok libraries, or of the Jollylib of Reservoir Labs. Other groups also made a lot of efforts towards the democratization of the use of polyhedral techniques, in particular the Alchemy Inria project, with Cloog and the development of Graphite in GCC, and Sadayappan's group in the USA, with the development of U. Bondhugula's Pluto prototype compiler. The same effort is made through the PPCG prototype compiler (for GPU) and Pencil (directives-based language on top of PPCG).

After 2009, Compsys continued to focus on the introduction of concepts and techniques to extend the polytope model, with a shift toward tools that may prepare the future. For instance, PoCo and C2fsm are able to parse general programs, not just SCoPs (static control programs), while the efficient handling of Boolean affine formulas [22] is a prerequisite for the construction of non-convex approximations. Euclidean lattices provide an efficient abstraction for the representation of spatial phenomena, and the construction of *critical lattices* as embedded in the tool Cl@k is a first step towards memory optimization in stream languages and may be useful in other situations. Our work on Chuba introduced a new element-wise array reuse analysis and the possibility of handling approximations. Our work on the analysis of while loops is both an extension of the polytope model itself (i.e., beyond SCoPs) and of its applications, here links with program termination and worst-case execution time (WCET) tools.

A recent example of the same approach is the proposal by Paul Feautrier to use polynomials for program analysis and optimization [6]. The associated tools are based on Handelman and Schweighofer theorems, the polynomial analogue of Farkas lemma. While this is definitely work in progress, with many unsolved questions, it has the potential of greatly enlarging the set of tractable programs.

As a last remark, observe that a common motif of these development is the transformation of finite algorithms into symbolic algorithms, able to solve very large or even infinite instances. For instance, PIP is a symbolic extension of the Simplex; our work on memory allocation is a symbolic extension of the familiar register allocation problem; loop scheduling extends DAG scheduling. Many other algorithms await their symbolic transformation: a case in point is resource-constrained scheduling.

# 4. Application Domains

## 4.1. Compilers for Embedded Computing Systems

The previous sections described our main activities in terms of research directions, but also places Compsys within the embedded computing systems domain, especially in Europe. We will therefore not come back here to the importance, for industry, of compilation and embedded computing systems design.

In terms of application domain, the embedded computing systems we consider are mostly used for multimedia: phones, TV sets, game platforms, etc. But, more than the final applications developed as programs, our main application is the computer itself: how the system is organized (architecture) and designed, how it is programmed (software), how programs are mapped to it (compilation and high-level synthesis).

The industry that can be impacted by our research is thus all the companies that develop embedded processors, hardware accelerators (programmable or not), embedded systems, and those (the same plus other) that need software tools to map applications to these platforms, i.e., that need to use or even develop programming languages, program optimization techniques, compilers, operating systems. Compsys do not focus on all these critical parts, but our activities are connected to them.

## 4.2. Users of HPC Platforms and Scientific Computing

The convergence between embedded computing systems and high-performance computing (HPC) technologies offers new computing platforms and tools for the users of scientific computing (e.g., people working in numerical analysis, in simulation, modeling, etc.). The proliferation of "cheap" hardware accelerators and multicores makes the "small HPC" (as opposed to computing centers with more powerful computers, grid computing, and exascale computing) accessible to a larger number of users, even though it is still difficult to exploit, due to the complexity of parallel programming, code tuning, interaction with compilers, which result from the multiple levels of parallelism and of memories in the recent architectures. The link between compiler and code optimization research (as in Compsys) and such users are still to be reinforced, both to guarantee the relevance of compiler research efforts with respect to application needs, and to help users better interact with compiler choices and understand performance issues.

The support of Labex MILYON (through its thematic quarters, such as the thematic quarter on compilation we organized in 2013 [4], or the upcoming 2016 thematic quarter on high-performance computing) and the activities of the LyonCalcul initiative [5] are means to get closer to users of scientific computing, even if it is too early to know if Compsys will indeed be directly helpful to them.

# 5. Highlights of the Year

## 5.1. Highlights of the Year

### Scientific Results

2015 showed good successes, in terms of scientific results, with respect to the objectives we fixed for Compsys III, i.e., pushing static compilation beyond its present limits, both in terms of techniques and applications, bridging the gap between polyhedral techniques and abstract interpretation, sequential codes and parallel specifications, back-end and front-end techniques. Important advances in 2015 are as follows:

- **Towards a polynomial model** We developed new techniques to handle polynomials (see Section 7.11) and thereby generalizing polyhedral (e.g., affine) techniques, with applications to the analysis of the OpenStream parallel language (see Section 7.10).

---

- **Handling parallel specifications** In complement to our current studies of parallel languages such as X10 (see Sections 7.8 and 7.9) and OpenStream (see Section 7.10), and kernel offloading with pipelined specifications (see Section 7.7), we succeeded to extend liveness analysis (see Section 7.12) and array contraction (see Section 7.13) to parallel specifications.

- **Enhancing interactions between programmer and compiler** This is an important challenge for the expansion of the applicability of our techniques. The work exposed in Sections 7.9 and 7.15 (effort for collecting and analyzing real applications), as well as the interaction with users of HPC, including the organization a joint spring school in 2016, are important steps in this direction.

- **Links with abstract interpretation and SMT solvers** The extension of our previous work on loop termination, with an iterative technique relying on SMT solvers for exhibiting counter-examples (see Section 7.4), is an interesting combination of polyhedral and abstract interpretation techniques. This is the case also for the array analysis of Section 7.3.

- **Back-end analysis** Considering back-end optimizations remains important, as complementary to front-end optimizations. See the results on register spilling (Section 7.1), pointer analysis (Section 7.2), liveness analysis (Section 7.12), the latter exploiting the fact that a polyhedral representation of arrays and loops is a symbolic unrolled view of registers and traces.

### Awards

The CC'15 paper on parametric tiling [3] was nominated as a best paper candidate for the group of conferences ETAPS'15 where, unfortunately, CC papers never finally got an award.

### End of Compsys

Compsys exists since 2012 as an Inria team. It has been created in 2004 as an Inria project-team, and evaluated by Inria first in 2007, then in 2012. It will again be evaluated in March 2016, which will be its final evaluation as an Inria project-team is limited to 12 years. The construction of a new project is thus necessary. The research directions of Compsys III were already a shift towards this future project. A few tentative research directions may be:

- Shift the application domain from embedded systems to high performance computing (HPC) but at small scale (desktop HPC: FPGA, GPU, multicores). In fact, the two ecosystems are nowadays slowly converging.

- A stronger attention to real HPC users and real HPC applications may lead to better programming models ("putting the programmer in the loop").

- Design new models of programs. The polynomial model is but an example.

- Explore the synergy between parallel programming and program verification and certification; in particular, import approximation methods from one field to the other. Abstract interpretation is a case in point.

However, while its field of expertise, compilation for parallel and heterogeneous systems, is still of crucial importance, the unexpected departure in Sep. 2015 of two of its staff members makes it difficult to have a clear view of the future.

# 6. New Software and Platforms

## 6.1. Aspic

Accelerated Symbolic Polyhedral Invariant Generation
KEYWORDS: Abstract Interpretation - Invariant Generation
FUNCTIONAL DESCRIPTION

Aspic is an invariant generator for general counter automata. Combined with C2fsm (a tool developed by P. Feautrier in Compsys), it can be used to derive invariants for numerical C programs, and also to prove safety. It is also part of the WTC toolsuite (see http://compsys-tools.ens-lyon.fr/wtc/index.html), a tool chain to compute worse-case time complexity of a given sequential program.

Aspic implements the theoretical results of Laure Gonnord's PhD thesis on acceleration techniques and has been maintained since 2007.

- Participant: Laure Gonnord
- Contact: Laure Gonnord
- URL: http://laure.gonnord.org/pro/aspic/aspic.html

## 6.2. DCC

DPN C Compiler
KEYWORDS: Polyhedral compilation - Automatic parallelization - High-level synthesis
FUNCTIONAL DESCRIPTION

Dcc (Data-aware process network C compiler) analyzes a sequential regular program written in C and generates an equivalent architecture of parallel computer as a communicating process network (Data-aware Process Network, DPN). Internal communications (channels) and external communications (external memory) are automatically handled while fitting optimally the characteristics of the global memory (latency and throughput). The parallelism can be tuned. Dcc has been registered at the APP ("Agence de protection des programmes") and transferred to the XtremLogic start-up under an Inria license.

- Participants: Christophe Alias and Alexandru Plesco
- Contact: Christophe Alias

## 6.3. Lattifold

Lattice-based Memory Folding
KEYWORDS: Polyhedral compilation - Euclidean Lattices
FUNCTIONAL DESCRIPTION

Implements advanced lattice-based memory folding techniques. The idea is to reduce memory footprint of multidimensional arrays by reducing the size of each dimension. Given a relation denoting conflicting array cells, it produces a new mapping based on affine functions bounded by moduli. The moduli induces memory reuse and bound memory accesses to a tighter area, allowing to reduce the array size without loss of correctness.

- Partner: ENS Lyon
- Contact: Alexandre Isoard

## 6.4. OpenOrdo

OpenStream scheduler
FUNCTIONAL DESCRIPTION

Finding polynomial schedules for the streaming language OpenStream. Main use: detecting deadlocks.

- Contact: Paul Feautrier

## 6.5. PoCo

Polyhedral Compilation library
KEYWORDS: Polyhedral compilation - Automatic parallelization
FUNCTIONAL DESCRIPTION

PoCo (Polyhedral Compilation library) is a compilation framework allowing to develop parallelizing compilers for regular programs. PoCo features many state-of-the-art polyhedral program analysis (dependences, affine scheduling, code generation) and a symbolic calculator on execution traces (represented as convex polyhedra). PoCo has been registered at the APP ("agence de protection des programmes") and transferred to the XtremLogic start-up under an Inria license.

- Participant: Christophe Alias
- Contact: Christophe Alias

## 6.6. PolyOrdo

Polynomial Scheduler
FUNCTIONAL DESCRIPTION

Computes a polynomial schedule for a sequential polyhedral program having no affine schedule. Uses algorithms for finding positive polynomials in semi-algebraic sets. Status: proof of concept software.

- Contact: Paul Feautrier

## 6.7. PPCG-ParamTiling

Parametric Tiling Extension for PPCG
KEYWORDS: Source-to-source compiler - Polyhedral compilation
FUNCTIONAL DESCRIPTION

PPCG is a source-to-source compiler, based on polyhedral techniques, targeting GPU architectures. It involves automatic parallelization and tiling using polyhedral techniques. This version replaces the static tiling of PPCG by a fully parametric tiling and code generator. It allows to choose tile sizes at run time when the memory size is known. It also provides a symbolic expression of memory usage depending on the problem size and the tile sizes.

- Partner: ENS Lyon
- Contact: Alexandre Isoard

## 6.8. Termite

Termination of C programs
KEYWORDS: Abstract Interpretation - Termination
FUNCTIONAL DESCRIPTION

TERMITE is the implementation of our new algorithm "Counter-example based generation of ranking functions" (see Section 7.4). Based on LLVM and Pagai (a tool that generates invariants), the tool automatically generates a ranking function for each *head of loop*.

TERMITE represents 3000 lines of OCaml and is now available via the opam installer.

- Participants: Laure Gonnord, Gabriel Radanne (PPS, Univ Paris 7), David Monniaux (CNRS/Verimag).
- Contact: Laure Gonnord
- URL: https://termite-analyser.github.io/

## 6.9. Vaphor

Validation of C programs with arrays with Horn Clauses

KEYWORDS: Abstract Interpretation - Safety - Array Programs
FUNCTIONAL DESCRIPTION

VAPHOR (Validation of Programs with Horn Clauses) is the implementation of our new algorithm "An encoding of array verification problems into array-free Horn clauses" (see Section 7.3). The tool implements a translation from a C-like imperative language into Horn clauses in the SMT-lib Format.

VAPHOR represents 2000 lines of OCaml and its development is under consolidation.

- Participants: Laure Gonnord, David Monniaux (CNRS/Verimag).
- Contact: Laure Gonnord
- URL: not yet published, under consolidation.

# 7. New Results

## 7.1. Studying Optimal Spilling in the Light of SSA

**Participants:** Florian Brandner [ENSTA ParisTech, previously Compsys], Quentin Colombet [Apple, previously Compsys], Alain Darte.

Recent developments in register allocation, mostly linked to static single assignment (SSA) form, have shown the benefits of decoupling the problem in two phases: a first spilling phase places load and store instructions so that the register pressure at all program points is small enough, and a second assignment and coalescing phase maps the variables to physical registers and reduces the number of move instructions among registers. We focused on the first phase, for which many open questions remain: in particular, we studied the notion of optimal spilling (what can be expressed?) and the impact of SSA form (does it help?).

To identify the important features for optimal spilling on load-store architectures, we developed a new integer linear programming formulation, more accurate and expressive than previous approaches. Among other features, we can express SSA $\phi$-functions, memory-to-memory copies, and the fact that a value can be stored simultaneously in a register and in memory. Based on this formulation, we presented a thorough analysis of the results obtained for the SPECINT 2000 and EEMBC 1.1 benchmarks, from which we have drawn, among others, the following conclusions: (1) rematerialization is extremely important; (2) SSA complicates the formulation of optimal spilling, especially because of memory coalescing when the code is not in conventional SSA (CSSA); (3) micro-architectural features are significant and thus have to be accounted for; and (4) significant savings can be obtained in terms of static spill costs, cache miss rates, and dynamic instruction counts.

Parts of this work were published at CASES 2011 [18]. The journal publication [1] contains more detailed discussions, more examples illustrating new concepts and existing approaches, and additional experiments covering the observed worst-case behavior, a new post-latency heuristic, and empiric evidence showing why static spill costs are a poor metric. Three configurations were added: Appel and George under SSA, Koes and Goldstein, and the heuristic of Braun and Hack.

## 7.2. Symbolic Range of Pointers in C programs

**Participants:** Vitor Paisante [Univ. Mineas Gerais, Brazil], Maroua Maalej, Leonardo Barbosa [Univ. Mineas Gerais, Brazil], Laure Gonnord, Fernando Pereira [Univ. Mineas Gerais, Brazil].

Alias analysis is one of the most fundamental techniques that compilers use to optimize languages with pointers. However, in spite of all the attention that this topic has received, the current state-of-the-art approaches inside compilers still face challenges regarding precision and speed. In particular, pointer arithmetic, a key feature in C and C++, is yet to be handled satisfactorily. We designed a new alias analysis algorithm to solve this problem. The key insight of our approach is to combine alias analysis with symbolic range analysis. This combination lets us disambiguate fields within arrays and structs, effectively achieving more precision than traditional algorithms. To validate our technique, we have implemented it on top of the LLVM compiler. Tests on a vast suite of benchmarks show that we can disambiguate several kinds of C idioms that current state-of-the-art analyses cannot deal with. In particular, we can disambiguate 1.35x more queries than the alias analysis currently available in LLVM. Furthermore, our analysis is very fast: we can go over one million assembly instructions in 10 seconds.

This work has been accepted at CGO'16 [30]. An extended version of the related work is available as an Inria research report [27] and will be the basis of a journal submission.

## 7.3. Analyzing C Programs with Arrays

**Participants:** Laure Gonnord, David Monniaux [CNRS/VERIMAG].

Automatically verifying safety properties of programs is hard, and it is even harder if the program acts upon arrays or other forms of maps. Many approaches exist for verifying programs operating upon Boolean and integer values (e.g., abstract interpretation, counter-examples guided abstraction refinement using interpolants), but transposing them to array properties has been fraught with difficulties.

In contrast to most preceding approaches, we do not introduce a new abstract domain or a new interpolation procedure for arrays. Instead, we generate an abstraction as a scalar problem and feed it to a preexisting solver. The intuition is that if there is a proof of safety of the program, it is likely that it can be expressed by elementary steps between properties involving only a small (tunable) number $N$ of cells from the array.

Our transformed problem is expressed using Horn clauses over scalar variables, a common format with clear and unambiguous logical semantics, for which there exist several solvers. In contrast, solvers directly operating over Horn clauses with arrays are still very immature.

An important characteristic of our encoding is that it creates a non-linear Horn problem, with tree unfoldings, contrary to the linear problems obtained by flatly encoding the control-graph structure. Our encoding thus cannot be expressed by encoding into another control-flow graph problem, and truly leverages the Horn clause format.

Experiments with our prototype VAPHOR (see Section 6.9) show that this approach can prove automatically the functional correctness of several classical examples of the literature, including *selection sort*, *bubble sort*, *insertion sort*, as well as examples from previous articles on array analysis.

This work is presented in a research report [28] and is currently under submission.

## 7.4. Termination of C Programs

**Participants:** Laure Gonnord, David Monniaux [CNRS/VERIMAG], Gabriel Radanne [Univ Paris 7/ PPS].

The work of Compsys on the generation of multi-dimensional ranking functions [15], through a mix of polyhedral and abstract interpretation techniques, and its implementation in the tool RanK [16], was continued by Laure Gonnord in collaboration with D. Monniaux. A complete method for synthesizing lexicographic linear ranking functions (and thus proving termination), supported by inductive invariants, was designed in the case where the transition relation of the program includes disjunctions and existentials (large block encoding of control flow).

Previous work would either synthesize a ranking function at every basic block head, not just loop headers, which reduces the scope of programs that may be proved to be terminating, or expand large block transitions including tests into (exponentially many) elementary transitions, prior to computing the ranking function, resulting in a very large global constraint system. In contrast, the new algorithm incrementally refines a global linear constraint system according to extremal counterexamples: only constraints that exclude spurious solutions are included.

Experiments with our tool Termite 6.8 show marked performance and scalability improvements compared to other systems.

This work has been published at the PLDI'15 conference [7].

## 7.5. Data-aware Process Networks

**Participants:** Christophe Alias, Alexandru Plesco [XtremLogic SAS].

High-level circuit synthesis (HLS, high-level synthesis) consists in compiling a program described in a high-level programming language (as C) to a circuit. The circuit must be as efficient as possible while using properly the resources (power consumption, silicon area, FPGA elementary units, memory accesses, etc). Although a lot of progress was achieved on the back-end (low-level) aspects (pipeline generation, place/route), the front-end aspects (parallelism, I/O) are still rudimentary compared to the techniques developed by the HPC community, notably the analysis stemming from the *polyhedral model*.

We introduced data-aware process networks (DPN), a parallel execution model adapted to the hardware constraints of high-level synthesis, where the data transfers are made explicit. We have shown that the DPN model is consistent in the sense that any translation of a sequential program produces an equivalent DPN without deadlocks. Finally, we show how to compile a sequential program to a DPN and how to optimize the input/output and the parallelism.

This work has been published as an Inria research report [9] and will be submitted to a journal.

## 7.6. Mono-parametric Tiling

**Participants:** Guillaume Iooss, Sanjay Rajopadhye [Colorado State University], Christophe Alias, Yun Zou [Colorado State University].

Tiling is a crucial program transformation with many benefits. It improves locality, exposes parallelism, allows for adjusting the ops-to-bytes balance of codes, and can be applied at multiple levels. Allowing tile sizes to be symbolic parameters at compile time has many benefits, including efficient auto-tuning, and run-time adaptability to system variations. For polyhedral programs, parametric tiling in its full generality is known to be non-linear, breaking the mathematical closure properties of the polyhedral model. Most compilation tools therefore either avoid it by only performing fixed size tiling, or apply it only in the final, code generation step. Both strategies have limitations.

We first introduced mono-parametric partitioning, a restricted parametric, tiling-like transformation that can be used to express a tiling. We showed that, despite being parametric, it is a polyhedral transformation. We first proved that applying mono-parametric partitioning (i) to a polyhedron yields a union of polyhedra, and (ii) to an affine function produces a piecewise-affine function. We then used these properties to show how to partition an entire polyhedral program, including one with reductions. Next, we generalized this transformation to tiles with arbitrary tile shapes that can tessellate the iteration space (e.g., hexagonal, trapezoidal, etc). We showed how mono-parametric tiling can be applied at multiple levels, and how it enables a wide range of polyhedral analyses and transformations to be applied.

This work has been published as an Inria research report [14] and will be submitted to a journal. It is the extended version of the work published at IMPACT'14 [26].

## 7.7. Exact and Approximated Data-Reuse Optimizations for Tiling with Parametric Sizes

**Participants:** Alain Darte, Alexandre Isoard.

As mentioned in Section 7.6, loop tiling is a loop transformation widely used to improve spatial and temporal data locality, to increase computation granularity, and to enable blocking algorithms, which are particularly useful when offloading kernels on computing units with smaller memories. When caches are not available or used, data transfers and local storage must be software-managed, and some useless remote communications can be avoided by exploiting data reuse between tiles. An important parameter of tiling is the sizes of the tiles, which impact the size of the required local memory. However, for most analyses involving several tiles, which is the case for inter-tile data reuse, the tile sizes induce non-linear constraints, unless they are numerical constants. This complicates or prevents a parametric analysis with polyhedral optimization techniques.

We showed that, when tiles are executed in sequence along tile axes, the parametric (with respect to tile sizes) analysis for inter-tile data reuse is nevertheless possible, i.e., one can determine, at compile-time and in a parametric fashion, the copy-in and copy-out data sets for all tiles, with inter-tile reuse, as well as sizes for the induced local memories (this is also connected to the liveness analysis described in Section 7.12). When approximations of transfers are performed, the situation is much more complex, and involves a careful analysis to guarantee correctness when data are both read and written. We provide the mathematical foundations to make such approximations possible, thanks to the introduction of the concept of *pointwise functions*. Combined with hierarchical tiling, this result opens perspectives for the automatic generation of blocking algorithms, guided by parametric cost models, where blocks can be pipelined and/or can contain parallelism. Previous work on FPGAs and GPUs already showed the interest and feasibility of such automation with tiling, but in a non-parametric fashion.

Our method is currently implemented with the `iscc` calculator of ISL, a library for the manipulation of integer sets defined with Presburger arithmetic, a complete implementation within the PPCG compiler is in progress (see also Section 6.7).

We believe that our approximation technique can be used for other applications linked to the extension of the polyhedral model as it turns out to be fairly powerful. Our future work will be to derive efficient approximation techniques, either because the program cannot be fully analyzable, or because approximations can speed-up or simplify the results of the analysis without losing much in terms of memory transfers and/or memory sizes.

A preliminary version of this work has been presented at the IMPACT'14 workshop [19]. A revised version was published at the International Conference on Compiler Construction (CC'15) [3].

## 7.8. Analysis of X10 Programs

**Participants:** Paul Feautrier, Alain Ketterlin [Inria/CAMUS], Sanjay Rajopadhye [Colorado State University], Vijay Saraswat [IBM Research], Eric Violard [Inria/CAMUS], Tomofumi Yuki.

While, historically, Compsys has applied polyhedral analysis to sequential programs, it was recently realized that it also applies to parallel programs or specifications, with the aim of checking their correctness or improving their performance. The prospect of having to program exascale architectures, with their millions of cores, has led to the development of new programming languages, whose objective is to increase the programmer productivity. Compsys has first applied polyhedral techniques to synchronous languages [24], [25] and pipelined specifications (see Section 7.7), before concentrating on IBM's high-productivity language X10 (see this section as well as Section 7.9) and on the OpenStream language (see Section 7.10).

X10 is based on the creation of independent *activities* (light-weight threads), which can synchronize either by a generalization of the fork/join scheme, or with *clocks*, an improved version of the familiar barriers. X10 is deadlock-free by construction but it is the programmer responsibility to insure determinism by a proper use of synchronizations. Non-determinism bugs may have a very low occurrence probability thus be very difficult to detect by testing, hence the interest for detecting races at compile time. In collaboration with CSU (S. Rajopadhye, T. Yuki) and IBM (V. Saraswat), we first extended array dataflow analysis to polyhedral clock-free X10 programs [34]. We have been working on clocked programs too. Race detection becomes undecidable [35], but realistic problems may still be solved by heuristics.

In cooperation with Eric Violard and Alain Ketterlin (Inria Team Camus, Strasbourg), and in order to obtain a more secure and precise analysis, we are currently attempting to formalize the "happens before" analysis used in these two previous papers [34], [35], using the proof assistant Coq.

## 7.9. Revisiting Loop Transformations with X10 Clocks

**Participant:** Tomofumi Yuki.

Loop transformations are known to be important for performance of compute-intensive programs, and are often used to expose parallelism. However, many transformations involving loops often obfuscate the code, and are cumbersome to apply by hand. In this work, we explored alternative methods for expressing parallelism that are more friendly to the programmer. In particular, we seek to expose parallelism without significantly changing the original loop structure. We illustrated how clocks in X10 can be used to express some of the traditional loop transformations, in the presence of parallelism, in a manner that we believe to be less invasive. Specifically, expressing parallelism corresponding to one-dimensional affine schedules can be achieved without modifying the original loop structure and/or statements.

This work was published at the international workshop on X10 [8].

## 7.10. Static Analysis of OpenStream Programs

**Participants:** Albert Cohen [Inria Parkas team], Alain Darte, Paul Feautrier.

In the context of the ManycoreLabs project (see Section 8.1), we also studied the applicability of polyhedral techniques to the parallel language OpenStream [31]. When applicable, polyhedral techniques are indeed invaluable for compile-time debugging and for generating efficient code well suited to a target architecture. OpenStream is a two-level language in which a control program directs the initialization of parallel task instances that communicate through *streams*, with possibly multiple writers and readers. It has a fairly complex semantics in its most general setting, but we restricted ourselves to the case where the control program is sequential, which is representative of the majority of the OpenStream applications.

In contrast to X10, this restriction offers deterministic concurrency by construction, but deadlocks are still possible. We showed that, if the control program is polyhedral, one may statically compute, for each task instance, the read and write indices to each of its streams, and thus reason statically about the dependences among task instances (the only scheduling constraints in this polyhedral subset). If the control program has nested loops, communications use one-dimensional channels in a form of linearization, and these indices may be polynomials of arbitrary degree, thus requiring to extend to polynomials the standard polyhedral techniques for dependence analysis, scheduling, and deadlock detection. Modern SMT allow to solve polynomial problems, albeit with no guarantee of success; the approach previously developed by P. Feautrier [6] may offer an alternative solution.

The usual way of disproving deadlocks is by exhibiting a schedule for the program operations, a well-known problem for polyhedral programs where dependences can be described by affine constraints. In the case of OpenStream, we established two important results related to deadlocks: 1) a characterization of deadlocks in terms of dependence paths, which implies that streams can be safely bounded as soon as a schedule exists with such sizes, 2) the proof that deadlock detection is undecidable, even for polyhedral OpenStream.

Details of this work are available in a research report [10]. It will be presented at the international workshop IMPACT'16 [2]. Some further developments are in progress for scheduling OpenStream programs using polynomial techniques, see Section 6.4.

## 7.11. Handling Polynomials for Program Analysis and Transformation

**Participant:** Paul Feautrier.

As shown in Section 7.10, many problems in parallel programs analysis and verification can be reduced to proving or disproving properties of polynomials in the variables of the program. For instance, the so-called "linearizations" (replacing a multi-dimensional object by a one-dimensional one) generate polynomial access functions. These polynomials then reappear in dependence testing, scheduling, and invariant construction. It may also happen that polynomials are absent from the source program, but are created either by an enabling analysis, as for OpenStream, or are imposed by complexity consideration. The usual solution is to construct a multi-dimensional function (e.g., a schedule for parallelization or a ranking function for termination [15]), which can then be converted into polynomials by counting. However, a direct approach is preferable, especially when the resulting schedule is to be used for further analysis, e.g., in real-time situations or WCET evaluation.

What is needed here is a replacement for the familiar emptiness tests and for Farkas lemma (deciding whether an affine form is positive inside a polyhedron). Recent mathematical results by Handelman and Schweighofer on the *Positivstellensatz* allow one to devise algorithms that are able to solve these problems. The difference is that one gets only sufficient conditions, and that complexity is much higher than in the affine cases. A paper presenting applications of these ideas to three use cases – dependence testing, scheduling, and transitive closure approximation – was presented at the 5th International Workshop on Polyhedral Compilation Techniques (IMPACT'15) [6] in Amsterdam in January 2015. A tool implementing polyhedral schedules complements this work, see Section 6.6.

## 7.12. Liveness Analysis in Explicitly-Parallel Programs

**Participants:** Alain Darte, Alexandre Isoard, Tomofumi Yuki.

In the light of the parallel specifications encountered in our other works (from Section 7.7 to Section 7.11), we revisited scalar and array element-wise liveness analysis for programs with parallel specifications. In earlier work on memory allocation/contraction (register allocation or intra- and inter-array reuse in the polyhedral model), a notion of "time" or a total order among the iteration points was used to compute the liveness of values. In general, the execution of parallel programs is not a total order, and hence the notion of time is not applicable.

We first revised how conflicts are computed by using ideas from liveness analysis for register allocation, studying the structure of the corresponding conflict/interference graphs. Instead of considering the conflict between two pairs of live ranges, we only consider the conflict between a live range and a write. This simplifies the formulation from having four instances involved in the test down to three, and also improves the precision of the analysis in the general case.

Then we extended the liveness analysis to work with partial orders so that it can be applied to many different parallel languages/specifications with different forms of parallelism. An important result is that the complement of the conflict graph with partial orders is directly connected to memory reuse, even in presence of races. However, programs with conditionals do not even have a partial order, and our next step will be to handle such cases with more accuracy.

Details of this work are available in a research report [13]. It will be presented at the international workshop IMPACT'16 [4].

## 7.13. Extended Lattice-Based Memory Allocation

**Participants:** Alain Darte, Alexandre Isoard, Tomofumi Yuki.

We extended lattice-based memory allocation [20], an earlier work on memory (array) reuse analysis. The main motivation is to handle in a better way the more general forms of specifications we see today, e.g., with loop tiling, pipelining, and other forms of parallelism available in explicitly parallel languages. Our extension has two complementary aspects. We showed how to handle more general specifications where conflicting constraints (those that describe the array indices that cannot share the same location) are specified as a (non-convex) union of polyhedra. Unlike convex specifications, this also requires to be able to choose suitable directions (or basis) of array reuse. For that, we extended two dual approaches, previously proposed for a fixed basis, into optimization schemes to select suitable basis. Our final approach relies on a combination of the two, also revealing their links with, on one hand, the construction of multi-dimensional schedules for parallelism and tiling (but with a fundamental difference that we identify) and, on the other hand, the construction of universal reuse vectors (UOV), which was only used so far in a specific context, for schedule-independent mapping.

This algorithmic work, connected to the parametric tiling of Section 7.7 and the liveness analysis results of Section 7.12, is complemented by a set of prototype scripting tools, see Section 6.3.

Details of this work are available in a research report. It has also been submitted to a conference.

## 7.14. Stencil Accelerators

**Participants:** Steven Derrien [University of Rennes 1, Inria/CAIRN], Xinyu Niu [Imperial College London], Sanjay Rajopadhye [Colorado State University], Tomofumi Yuki.

Stencil computations have been known to be an important class of programs for scientific calculations. Recently, various architectures (mostly targeting FPGAs) for stencils are being proposed as hardware accelerators with high throughput and/or high energy efficiency. There are many different challenges for such design: How to maximize compute-I/O ratio? How to partition the problem so that the data fits on the on-chip memory? How to efficiently pipeline? How to control the area usage? We seek to address these challenges by combining techniques from compilers and high-level synthesis tools.

One project in collaboration with the CAIRN team and Colorado State University targets stencils with regular dependence patterns. Although many architectures have been proposed for this type of stencils, most of them use a large number of small processing elements (PE) to achieve high throughput. We are exploring an alternative design that aims for a single, large, deeply-pipelined PE. The hypothesis is that the pipelined parallelism is more area-efficient compared to replicating small PEs. We have published a work-in-progress paper on this topic at IMPACT'16 [5].

Another type of stencil accelerators that we are working on, in collaboration with Xinyu Niu, targets stencil programs with dynamic dependences (i.e., sparse computations). The collaboration is in the context of the EURECA project [6] where the dynamic reconfigurability of modern FPGAs are used to efficiently handle dynamic access patterns.

## 7.15. PolyApps

**Participant:** Tomofumi Yuki.

Loop transformation frameworks using the polyhedral model have gained increased attention since the rise of the multi-core era. We now have several research tools that have demonstrated their power on important kernels found in scientific computations. However, there remains a large gap between the typical kernels used to evaluate these tools and the actual applications used by the scientists.

PolyApps is an effort to collect applications from other domains of science to better establish the link between the compiler tools and "real" applications. The applications are modified to bypass some of the front-end issues of research tools, while keeping the ability to produce the original output. The goal is to assess how the state-of-the-art automatic parallelizers perform on full applications, and to identify new opportunities that only arise in larger pieces of code.

We showed that, with a few enhancements, the current tools will be able to reach and/or exceed the performance of existing parallelizations of the applications. One of the most critical element missing in current tools is the ability to modify the memory mappings.

# 8. Bilateral Contracts and Grants with Industry

## 8.1. ManycoreLabs Project with Kalray

Compsys was part of 3-years a bilateral contract with Kalray called ManycoreLabs, funded by "Investissements d'avenir pour le développement de l'économie numérique". The goal of this project was to allow the company Kalray, based on a collaboration with several partners, to become the European leader of the market of many-core chips for embedded systems. Industrial partners of this project included Bull, CAPS Entreprise, Digigram, Thales, Renault. Academic partners are CEA, Inria (Parkas, Compsys, and Corse), VERIMAG.

---

[6]http://www.doc.ic.ac.uk/~nx210/2015/09/01/eureca.html

Compsys role was to explore analysis and compilation techniques linked to streaming languages, with the Kalray MPPA platform as long-term target. The research on OpenStream described in Section 7.8 corresponds to extensions of the work package WP 2.5.3. This study showed the need for extending polyhedral techniques to polynomials, which is one of the motivation of the work described in Section 7.11. The work on parametric tiling (Section 7.7), first in the context of FPGA, then of GPUs, was also a first step towards the automatic generation of blocking algorithms for multicores such as the Kalray MPPA.

This project ended in June 2015.

## 8.2. Technological Transfer: XtremLogic Start-Up

The XTREMLOGIC start-up (http://xtremlogic.com/) was initiated, initially with the name Zettice, at the end of 2010 by Alexandru Plesco and Christophe Alias, after the PhD thesis of Alexandru Plesco under the guidance of Christophe Alias, Alain Darte and Tanguy Risset. The goal of XTREMLOGIC is to build on the disruptive technologies emerging from the polyhedral compilation community, and particularly the results obtained in Compsys, to provide the HPC market with efficient and communication-optimal circuit blocks (IP) for FPGA.

The compiler technology transferred to XTREMLOGIC (see Sections 6.2 and 7.5) is the result of a tight collaboration between Christophe Alias and Alexandru Plesco. XTREMLOGIC is one way to spread the polyhedral technology to industry. In 2015, XTREMLOGIC was supported by the Rhône Développement Initiative 2015 (loan).

# 9. Partnerships and Cooperations

## 9.1. Regional Initiatives

### 9.1.1. *In Relation with the LYONCALCUL Initiative*

Compsys follows or participates to the activities of LyonCalcul (http://lyoncalcul.univ-lyon1.fr/), a network to federate activities on high-performance computing in Lyon.

In this context, and with the support of the Labex MILYON (http://milyon.universite-lyon.fr/), Compsys organized in 2013 a thematic quarter on compilation (http://labexcompilation.ens-lyon.fr). A new thematic quarter on high performance computing (HPC) is in preparation for 2016, initiated by Violaine Louvet (Institute Camille Jordan), with the participation of the LIP teams Aric, Avalon, Compsys, and Roma. It will include, in particular, an inter-disciplinary spring school, following the polyhedral school organized in 2013, connecting mathematics (HPC numerical analysis) and computer science (polyhedral optimizations for HPC).

Alain Darte, Alexandre Isoard, and Tomofumi Yuki have also regular exchanges with Violaine Louvet and Thierry Dumont on tiling code optimizations, advising (in an informal way) some of their students during their internships, for implementations on multicore machines and GPUs.

### 9.1.2. *Collaboration with the Verimag lab*

Laure Gonnord, who did her PhD in abstract interpretation at Verimag, re-activated her connection with this group, in particular with N. Halbwachs and D. Monniaux. This led to several joint results, exposed in Sections 7.3 and 7.4. The theme of termination through affine ranking functions was first brought to the attention of Compsys when studying loop transformations for HLS, in the context of the S2S4HLS project with STMicroelectronics. The techniques of Compsys [15] were then extended by Laure Gonnord with D. Monniaux. Conversely, the idea of using Handelman and Schweighofer's theorems to deal with polynomial constraints, as exploited in Section 7.11), was first suggested by D. Monniaux through discussions with Paul Feautrier and some visits at ENS-Lyon.

### 9.1.3. *"PEPS local" with the MMI*

Alain Darte and Laure Gonnord participated to the creation of EMI (Education, Musique et Informatique), an educative inter-disciplinary project ("PEPS de site", coordinated by Natacha Portier, from the MC2 team at LIP, and Yann Orlarey from the Grame laboratory) concerning an experience of musical programming with Faust (a functional audio stream language, with its compiler), in the context of the MMI (Maison des mathématiques et de l'informatique), a place for dissemination.

## 9.2. National Initiatives

### 9.2.1. *French Compiler Community*

In 2010, Laure Gonnord and Fabrice Rastello created the french community of compilation, which had no organized venue in the past. All groups with activities related to compilation were contacted and the first "compilation day" was organized in Lyon.This effort has been quickly a success: the community (http://compilfr.ens-lyon.fr/) is now well identified and 3-days workshops now occur at least once a year (the 10th event has been organized in Sep. 2015). The community is animated by Laure Gonnord and Fabrice Rastello since 2010, and now also by Florian Brandner (ex-Compsys too). Alain Darte, Alexandre Isoard, and Tomofumi Yuki participated to the 10th edition, with talks on "Static Analysis of OpenStream Programs", "Liveness Analysis in the Polyhedral Model", and "PolyApps: Case Study of Polyhedral Compilers using Real Applications" respectively.

Recognized as a sub-group of the CNRS GDR GPL (Software Engineering and Programming), the community is also in charge, since 2014, of organizing one day of the research school "Ecole des jeunes chercheurs en Algorithmique et Programmation" (EJCP). Tomofumi Yuki, in this context, gave a one-day lecture at the 2015 edition.

### 9.2.2. *Collaboration with Parkas group, in Paris*

Alain Darte and Paul Feautrier have regular meetings with Albert Cohen, from the Parkas team at ENS Paris. The current discussions are mostly related to the analysis and compilation of the OpenStream language developed by Parkas, a research topic that started though the ManycoreLabs project (see Section 8.1). The results of Sections 7.10 and 7.11 are related to this collaboration.

### 9.2.3. *Collaboration with Cairn group, in Rennes*

Tomofumi Yuki continues to work with the Cairn group through regular meetings and occasional visits. The topic of the collaboration is in applying compiler techniques for hardware design using high-level synthesis. Section 7.14 presents the results through this collaboration.

### 9.2.4. *Collaboration with Camus group, in Strasbourg*

Paul Feautrier and Tomofumi Yuki have an ongoing cooperation with Alain Ketterlin and Eric Violard (Camus group, Strasbourg) on several subjects connected to the analysis and transformations of X10 programs (see Section 7.8).

## 9.3. European Initiatives

### 9.3.1. *FP7 & H2020 Projects*

Compsys participated to a H2020 proposal (project Verde) on the convergence of compiler tools for hardware accelerators on one side (HLS tools) and programmable accelerators (multicores, GPUs) on the other side. But the project was not selected.

### 9.3.2. *Collaborations with Major European Organizations*

Compsys members participate to the European Network of Excellence on High Performance and Embedded Architecture and Compilation (HiPEAC, http://www.hipeac.net/), either as members or affiliate members. The International Workshop on Polyhedral Compilation Techniques (IMPACT, see Section 9.4.2.2), co-created by Christophe Alias in 2011, is now an annual event of the HIPEAC conference, as an official workshop. The 5th edition, IMPACT'15, was co-chaired by Alain Darte (see http://impact.gforge.inria.fr/impact2015/), while the 6h edition, IMPACT'16, was co-chaired by Tomofumi Yuki (see http://impact.gforge.inria.fr/impact2016/).

## 9.4. International Initiatives

### 9.4.1. *Inria Associate Teams not Involved in an Inria International Labs*

Laure Gonnord and Maroua Maleej are involved in the PROSPIEL Associate Team (Inria/Brazil, https://team.inria.fr/alf/prospiel/), led by Sylvain Collange (Inria Alf), in a collaboration with Fernando Pereira's group in UFMG (Brazil). The PROSPIEL project aims at optimizing parallel applications for high performance on new throughput-oriented architectures: GPUs and many-core processors. Specifically, Laure Gonnord and Maroua Maalej are in charge of designing static analyses for GPUs. Maroua Maleej visited the group of Fernando Pereira in Aug. 2015.

### 9.4.2. *Inria International Partners*

#### 9.4.2.1. Declared Inria International Partners

- Christophe Alias is co-adviser, with Sanjay Rajopadhye from Colorado State University (USA), of the PhD thesis of Guillaume Iooss. The results described in Section 7.6 are part of this collaboration.

- Tomofumi Yuki, who did his PhD with Sanjay Rajopadhye, then a post-doc in the Cairn team in Rennes, continues his collaboration with these two groups, as the results described in Section 7.14 illustrate. He participates regularly, over the net, to the reading group "Melange" of S. Rajodapdhye's group, with CSU students.

- Laure Gonnord and Maroua Maleej have a regular collaboration with Fernando Magno Quintao Pereira from the University of Mineas Gerais (Brazil). The results described in Section 7.2 are part of this collaboration. In Jan.-Feb. 2015, Compsys hosted Fernando Pereira, as a visiting professor.

#### 9.4.2.2. Polyhedral Community

In 2011, as part of the organization of the workshops at CGO'11, Christophe Alias (with C. Bastoul) organized IMPACT'11 (international workshop on polyhedral compilation techniques, http://impact2011.inrialpes.fr/). This workshop in Chamonix was the very first international event on this topic, although it was introduced by Paul Feautrier in the late 80s. Alain Darte gave the introductory keynote talk. After this successful edition (more than 60 people), IMPACT continued as a satellite workshop of the HIPEAC conference, in Paris (2012), Berlin (2013), Vienna (2014). Alain Darte was program co-chair and co-organizer for the past edition, in Amsterdam (2015), while Tomofumi Yuki is program co-chair and co-organizer of the next one, in Prague (2016).

The creation of IMPACT, now the annual event of the polyhedral community, helped to identify this community and to make it more visible. This effort was complemented by the organization of the first (and for the moment unique) school on polyhedral code analysis and optimizations (http://labexcompilation.ens-lyon.fr/polyhedral-school/). A second polyhedral school, more open, because involving themes and researchers from numerical analysis (users of HPC), will be organized in 2016.

Alain Darte also manages two new mailing lists for news (polyhedral-news@listes.ens-lyon.fr) and discussions (polyhedral-discuss@listes.ens-lyon.fr) on polyhedral code analysis and optimizations. Tomofumi Yuki is involved in the development of PolyBench (http://sourceforge.net/projects/polybench), a suite of kernels used for illustrating polyhedral optimizations. He is also developing PolyApps, a set of larger applications to evaluate the gap between kernels and "real" applications, see more details in Section 7.15.

## 9.5. International Research Visitors

### 9.5.1. Visits of International Scientists

*9.5.1.1. Invited Professors*

- Fernando M. Pereira was invited in Jan. 2015 to work with Maroua Maleej and Laure Gonnord on static analyses for pointers.

*9.5.1.2. Internships*

- Tristan Dubois, M1 student from Lyon 1 University, worked for 6 weeks in January-February 2015, on pointer arithmetic in LLVM, supervised by Laure Gonnord.
- Marc Vincenti, M1 student from Lyon 1 University, worked for 6 weeks in January-February 2015, on comparison of termination benchmarks, in the context of the Artefact Evaluation of the PLDI'15 publication [7], whose results are described in Section 7.4.
- Adilla Susungi, a M2 student from Strasbourg University, worked, from March 2015 to July 2015, on the compilation of streaming applications on multi-GPUs, supervised by Christophe Alias. Her internship was funded by Inria.

### 9.5.2. Visits to International Teams

Paul Feautrier has been invited by the University of Passau (Bavaria) in the team of Prof. Christian Lengauer, where he has given a seminar "Toward a Polynomial Model" (September 2015) and held scientific discussions with Armin Groesslinger and other members of the team.

# 10. Dissemination

## 10.1. Promoting Scientific Activities

### 10.1.1. Scientific Events Organization

*10.1.1.1. General Chair, Scientific Chair*

Laure Gonnord is co-chair of the "french compilation community", with Florian Brandner (ENSTA, then Telecom ParisTech) and Fabrice Rastello (Inria Corse).

Alain Darte is general chair of the steering committee of CPC (International Workshop on Compilers for Parallel Computing), which regroups in Europe, every 18 months, a large community of researchers interested in compilers for HPC. Also, Alain Darte and Alexandre Isoard participated to CPC'15 in London (Jan. 15).

*10.1.1.2. Member of Organizing Committees*

Alain Darte was co-organizer of IMPACT'15 (International Workshop on Polyhedral Compilation Techniques) and Tomofumi Yuki is co-organizer of IMPACT'16.

Tomofumi Yuki was publicity co-chair of GPCE 2015 (15th International Conference on Generative Programming: Concepts & Experience).

Alain Darte and Tomofumi Yuki are currently organizing a second polyhedral spring school in 2016, with connections with HPC users from numerical analysis.

### 10.1.2. Scientific Events Selection

*10.1.2.1. Chair of Conference Program Committees*

Alain Darte was program co-chair of IMPACT'15, with Alexandra Jimborean (Uppsala University). Tomofumi Yuki is program co-chair of IMPACT'16, with Michelle Strout (University of Arizona).

Alain Darte was program chair of the topic E2 "Compilers for Embedded Systems" of DATE'15 (International Conference on Design, Automation, and Test in Europe), with Rodric Rabbah (IBM).

*10.1.2.2. Member of Conference Program Committees*

Tomofumi Yuki was a program committee member for the RST Track (Reliable Software Technologies and Communication Middleware) of SAC'16 (31st ACM Symposium on Applied Computing).

Paul Feautrier was a member of the program committees of IMPACT'15, IMPACT'16, and PECCS'15 (5th International Conference on Pervasive and Embedded Computing and Communication Systems).

Christophe Alias was a member of the program committee of IMPACT'16.

Alain Darte was a member of the program committees of PACT'15 (International Conference on Parallel Architectures and Compilation Techniques) and X10'15 (international workshop on X10, part of PLDI'15).

*10.1.2.3. Reviewer*

Paul Feautrier was a reviewer for IMPACT'15, IMPACT'16, PARCO'15, and PACT'15.

Tomofumi Yuki was a reviewer for PACT'15 and X10'15.

Christophe Alias was a reviewer for DATE'15.

Laure Gonnord was a reviewer for VMCAI'15, CGO'15, and PARCO'15.

Alain Darte was a reviewer for DATE'15, IMPACT'15, X10'15, and PACT'15.

### 10.1.3. Journal

*10.1.3.1. Member of Editorial Boards*

Paul Feautrier is a member of the editorial board of IJPP, the "International Journal of Parallel Programming".

*10.1.3.2. Reviewer - Reviewing Activities*

Paul Feautrier was reviewer for "Information and Computation", ACM TODAES, ACM TOPLAS, and IJPP.

Tomofumi Yuki was reviewer for the PARCO journal.

Alain Darte was reviewer for the ACM TACO journal and the "Software Practice and Experience" journal.

Christophe Alias was reviewer for Parallel Computing and IEEE TVLSI.

### 10.1.4. Invited Talks

Paul Feautrier was invited to give a talk on "The Numerical Solution of the Transfer Equation", at a workshop in honor of Roger Cayrel, Paris Observatory, Dec. 2015.

In June 2015, Laure Gonnord was invited at Google, Mountain View and SRI, to give talks about her research about static analyses for compilers.

### 10.1.5. Scientific Expertise

Alain Darte was invited to be part of the scientific committee of the CPU ("cerfication numérique et fiabilité") cluster of excellence (from Bordeaux Idex), and its internal evaluation in Sep. 2015.

In 2015, Maroua Maleej has produced 7 Research Tax Credit documents for Accenture group France as a scientific consultant. The goal is to expertise research done by Accenture project-teams and suggest further ideas by evaluating the state of the art.

## 10.2. Teaching - Supervision - Juries

### 10.2.1. Teaching

Licence:

- Laure Gonnord: Architecture des ordinateurs (TD+TP=40h), L2, Université Lyon 1 Claude Bernard: Spring 2015.
- Maroua Maleej: Algorithmique et Programmation Fonctionnelle et Récursive (TP=28h), L1, Université Lyon 1 Claude Bernard: Fall 2015.

- Christophe Alias: Introduction à la compilation (CM+TD=24h), L3, INSA Centre-Val-de-Loire: Spring 2015.
- Christophe Alias: Concours E3A—épreuve informatique MPSI (correcteur): Spring 2015.

Master:

- Laure Gonnord, Program Analysis and Verification (CM 24h), M1, Ecole Normale Supérieure de Lyon. With David Monniaux.
- Laure Gonnord, Compilation (TP 28h), M1, Ecole Normale Supérieure de Lyon.
- Laure Gonnord, Introduction aux systèmes et réseaux (CM/TP 52h), M2 Pro, Université Lyon 1.
- Laure Gonnord, Compilation (TD/TP 24h), M1, Université Lyon 1 Claude Bernard.
- Laure Gonnord, Complexité (TD 15h), M1, Université Lyon 1 Claude Bernard.
- Laure Gonnord, Temps Réel (TP 24h), M1, Université Lyon 1 Claude Bernard.
- Christophe Alias, Optimisation d'applications embarquées (CM+TD=18h), M1, INSA Centre-Val-de-Loire.
- Christophe Alias, Advanced Compilers: Loop Transformations and High-Level Synthesis (CM 8h), M2, Ecole Normale Supérieure de Lyon.
- Christophe Alias, Compilation (CM 16h), M1, Ecole Normale Supérieure de Lyon.
- Tomofumi Yuki and Christophe Alias, Advanced Compilers: Loop Transformations and High-Level Synthesis, 24h, M2, ENS Lyon.

Master school:

- Laure Gonnord organized in Jan. 2015 a research school entitled "Static analyses in the state-of-the-art compilers" (invited speaker: Fernando Pereira). See http://laure.gonnord. org/pro/research/compil_research_school.html

EJCP:

Tomofumi Yuki has given a one-day lecture at the École Jeunes Chercheurs en Programmation 2015. See http://ejcp2015.inria.fr/.

## *10.2.2. Supervision*

PhD in progress: Guillaume Iooss, "Semantic tiling", started in September 2011, joint PhD ENS-Lyon/Colorado State University, advisors: Christophe Alias and Alain Darte (ENS-Lyon) / Sanjay Rajopadhye (Colorado State University).

PhD in progress: Alexandre Isoard, "Streaming-related code optimizations", started in September 2012, advisor: Alain Darte.

PhD in progress: Maroua Maleej, "Low cost static analyses for compilers", started in October 2014, advisors : Laure Gonnord and Alain Darte, then Laure Gonnord and Frédéric Vivien (Roma team).

## *10.2.3. Juries*

Paul Feautrier was a member of the defense committee of the PhD of Alexis Foulhié (Grenoble, Oct. 2015), entitled "Revisiting the abstract domain of polyhedra: constraints-only representation and formal proof", and was a reviewer for the HDR thesis of Corinne Ancourt (UPMC, May 2015), entitled "Sûreté: de l'analyse à l'instrumentation et à la synthese de code".

Laure Gonnord was a member of the doctoral committee for the evaluation of the first year of the PhD of F. Maurica, Université de la Réunion, in Dec. 2015.

Alain Darte was the reviewer of the PhD of Gergö Barany (Technische Universität Wien, Austria, March 2015), entitled "Integrated Code Motion and Register Allocation".

## 10.3. Popularization

Compsys was involved in the proposal of the inter-disciplinary project EMI, with the maison des mathématiques et de l'informatique (MMI), and the Grame laboratory (computer music). See Section 9.1.3.

Tomofumi Yuki has given a one-day lecture at the École Jeunes Chercheurs en Programmation 2015 (http://ejcp2015.inria.fr/), entitled "Research in Compilers and How it Relates to Software Engineering".

# 11. Bibliography

## Publications of the year

### Articles in International Peer-Reviewed Journals

[1] Q. COLOMBET, F. BRANDNER, A. DARTE. *Studying Optimal Spilling in the Light of SSA*, in "ACM Transactions on Architecture and Code Optimization", January 2015, vol. 11-4, n$^o$ 47, 26 p. [*DOI :* 10.1145/2685392], https://hal.inria.fr/hal-01099016

### International Conferences with Proceedings

[2] A. COHEN, A. DARTE, P. FEAUTRIER. *Static Analysis of OpenStream Programs*, in "6th International Workshop on Polyhedral Compilation Techniques (IMPACT'16), held with HIPEAC'16", Prague, Czech Republic, Proceedings of the IMPACT series, http://impact.gforge.inria.fr/, Michelle Strout and Tomofumi Yuki, January 2016, https://hal.inria.fr/hal-01251845

[3] A. DARTE, A. ISOARD. *Exact and Approximated Data-Reuse Optimizations for Tiling with Parametric Sizes*, in "24th International Conference on Compiler Construction (CC'15), part of ETAPS'15", London, United Kingdom, April 2015, https://hal.inria.fr/hal-01099017

[4] A. DARTE, A. ISOARD, T. YUKI. *Liveness Analysis in Explicitly-Parallel Programs*, in "6th International Workshop on Polyhedral Compilation Techniques (IMPACT'16), held with HIPEAC'16", Prague, Czech Republic, Proceedings of the IMPACT series, Michelle Strout and Tomofumi Yuki, January 2016, http://impact.gforge.inria.fr/ , https://hal.inria.fr/hal-01251843

[5] G. DEEST, N. ESTIBALS, T. YUKI, S. DERRIEN, S. RAJOPADHYE. *Towards Scalable and Efficient FPGA Stencil Accelerators*, in "6th International Workshop on Polyhedral Compilation Techniques (IMPACT'16), held with HIPEAC'16", Prague, Czech Republic, Proceedings of the IMPACT series, January 2016, http://impact.gforge.inria.fr/ , https://hal.inria.fr/hal-01254778

[6] P. FEAUTRIER. *The Power of Polynomials*, in "5th International Workshop on Polyhedral Compilation Techniques (IMPACT'15)", Amsterdam, Netherlands, A. JIMBOREAN, A. DARTE (editors), January 2015, https://hal.inria.fr/hal-01094787

[7] L. GONNORD, D. MONNIAUX, G. RADANNE. *Synthesis of ranking functions using extremal counterexamples*, in "Programming Languages, Design and Implementation", Portland, Oregon, United States, June 2015 [*DOI :* 10.1145/2737924.2737976], https://hal.archives-ouvertes.fr/hal-01144622

[8] T. YUKI. *Revisiting Loop Transformations with X10 Clocks*, in "Proceedings of the ACM SIGPLAN Workshop on X10", Portland, OR, United States, June 2015 [*DOI :* 10.1145/2771774.2771778], https://hal.inria.fr/hal-01253630

### Research Reports

[9]  C. ALIAS, A. PLESCO. *Data-aware Process Networks*, Inria - Research Centre Grenoble – Rhône-Alpes ;
     Inria, June 2015, n⁰ RR-8735, 32 p. , https://hal.inria.fr/hal-01158726

[10] A. COHEN, A. DARTE, P. FEAUTRIER. *Static Analysis of OpenStream Programs*, CNRS ; Inria ; ENS Lyon,
     January 2016, n⁰ RR-8764, 26 p. , Corresponding publication at IMPACT'16 (http://impact.gforge.inria.fr/
     impact2016), https://hal.inria.fr/hal-01184408

[11] A. DARTE, A. ISOARD. *Exact and Approximated Data-Reuse Optimizations for Tiling with Parametric Sizes*,
     LIP - ENS Lyon ; CNRS ; Inria ; UCBL, January 2015, n⁰ RR-8671, 28 p. , https://hal.inria.fr/hal-01103460

[12] A. DARTE, A. ISOARD, T. YUKI. *Extended Lattice-Based Memory Allocation*, CNRS ; ENS Lyon ; Inria,
     November 2015, n⁰ RR-8840, 31 p. , https://hal.inria.fr/hal-01251868

[13] A. DARTE, A. ISOARD, T. YUKI. *Liveness Analysis in Explicitly-Parallel Programs*, CNRS ; Inria ; ENS
     Lyon, January 2016, n⁰ RR-8839, 25 p. , Corresponding publication at IMPACT'16 (http://impact.gforge.
     inria.fr/impact2016), https://hal.inria.fr/hal-01251579

[14] G. IOOSS, S. RAJOPADHYE, C. ALIAS, Y. ZOU. *Mono-parametric Tiling is a Polyhedral Transformation*,
     Inria Grenoble - Rhône-Alpes ; CNRS, October 2015, n⁰ RR-8802, 40 p. , https://hal.inria.fr/hal-01219452

## References in notes

[15] C. ALIAS, A. DARTE, P. FEAUTRIER, L. GONNORD. *Multi-dimensional Rankings, Program Termination,
     and Complexity Bounds of Flowchart Programs*, in "17th International Static Analysis Symposium (SAS'10)",
     Perpignan, France, ACM press, September 2010, pp. 117-133

[16] C. ALIAS, A. DARTE, P. FEAUTRIER, L. GONNORD. *Rank: A Tool to Check Program Termination and
     Computational Complexity*, in "International Workshop on Constraints in Software Testing Verification and
     Analysis (CSTVA'13)", Luxembourg, March 2013, 238 p. , http://hal.inria.fr/hal-00801571

[17] C. ALIAS, A. DARTE, A. PLESCO. *Optimizing Remote Accesses for Offloaded Kernels: Application to
     High-Level Synthesis for FPGA*, in "International Conference on Design, Automation and Test in Europe
     (DATE'13)", Grenoble, France, March 2013, pp. 575-580

[18] Q. COLOMBET, F. BRANDNER, A. DARTE. *Studying Optimal Spilling in the Light of SSA*, in "International
     Conference on Compilers, Architectures, and Synthesis of Embedded Systems (CASES'11)", Taipei, Taiwan,
     ACM, October 2011, pp. 25–34

[19] A. DARTE, A. ISOARD. *Parametric Tiling with Inter-Tile Data Reuse*, in "4th International Workshop on
     Polyhedral Compilation Techniques (IMPACT'14)", Vienna, Austria, S. RAJOPADHYE, S. VERDOOLAEGE
     (editors), January 2014, https://hal.archives-ouvertes.fr/hal-00915831

[20] A. DARTE, R. SCHREIBER, G. VILLARD. *Lattice-Based Memory Allocation*, in "IEEE Transactions on
     Computers", October 2005, vol. 54, n⁰ 10, pp. 1242-1257, Special Issue: Tribute to B. Ramakrishna (Bob)
     Rau

[21] P. FEAUTRIER. *Scalable and Structured Scheduling*, in "International Journal of Parallel Programming", October 2006, vol. 34, n⁰ 5, pp. 459–487

[22] P. FEAUTRIER. *Simplification of Boolean Affine Formulas*, Inria, July 2011, n⁰ RR-7689, http://hal.inria.fr/inria-00609519/PDF/RR-7689.pdf

[23] P. FEAUTRIER. *Dataflow Analysis of Scalar and Array References*, in "International Journal of Parallel Programming", February 1991, vol. 20, n⁰ 1, pp. 23–53

[24] P. FEAUTRIER, A. GAMATIÉ, L. GONNORD. *Enhancing the Compilation of Synchronous Dataflow Programs with a Combined Numerical-Boolean Abstraction*, in "CSI Journal of Computing", 2012, vol. 1, n⁰ 4, 8:86 p.

[25] A. GAMATIÉ, L. GONNORD. *Static Analysis of Synchronous Programs in Signal for Efficient Design of Multi-Clocked Embedded Systems*, in "International Conference on Languages, Compilers, and Tools for Embedded Systems (LCTES'11)", Chicago, USA, April 2011

[26] G. IOOSS, S. RAJOPADHYE, C. ALIAS, Y. ZOU. *CART: Constant Aspect Ratio Tiling*, in "4th International Workshop on Polyhedral Compilation Techniques (IMPACT'14)", Vienna, Austria, S. RAJOPADHYE, S. VERDOOLAEGE (editors), January 2014, https://hal.archives-ouvertes.fr/hal-00915827

[27] M. MAALEJ, L. GONNORD. *Do we still need new Alias Analyses?*, Université Lyon Claude Bernard / Laboratoire d'Informatique du Parallélisme, November 2015, n⁰ RR-8812, https://hal.inria.fr/hal-01228581

[28] D. MONNIAUX, L. GONNORD. *An Encoding of Array Verification Problems into Array-Free Horn Clauses*, July 2015, working paper or preprint, https://hal.archives-ouvertes.fr/hal-01206882

[29] H. NAZARÉ, I. MAFFRA, W. SANTOS, L. OLIVEIRA, F. M. Q. PEREIRA, L. GONNORD. *Validation of Memory Accesses Through Symbolic Analyses*, in "ACM International Conference on Object Oriented Programming Systems Languages & Applications (OOPSLA'14)", Portland, Oregon, United States, October 2014, pp. 791-809, https://hal.inria.fr/hal-01006209

[30] V. PAISANTE, M. MAALEJ, L. BARBOSA, L. GONNORD, F. M. Q. PEREIRA. *Symbolic Range Analysis of Pointers*, in "International Symposium of Code Generation and Optmization (CGO'16)", Barcelone, Spain, March 2016, pp. 791-809 [*DOI : 10.1145/2660193.2660205*], https://hal.inria.fr/hal-01228928

[31] A. POP, A. COHEN. *OpenStream: Expressiveness and data-flow compilation of OpenMP streaming programs*, in "ACM Transactions on Architecture and Code Optimization (TACO)", 2013, vol. 9, n⁰ 4, pp. 1-25

[32] A. TURJAN, B. KIENHUIS, E. DEPRETTERE. *Translating Affine Nested-Loop Programs to Process Networks*, in "International Conference on Compilers, Architecture, and Synthesis for Embedded Systems (CASES'04)", New York, NY, USA, ACM, 2004, pp. 220–229

[33] S. VERDOOLAEGE, H. NIKOLOV, N. TODOR, P. STEFANOV. *Improved Derivation of Process Networks*, in "International Workshop on Optimization for DSP and Embedded Systems (ODES'06)", 2006

[34] T. YUKI, P. FEAUTRIER, S. RAJOPADHYE, V. SARASWAT. *Array Dataflow Analysis for Polyhedral X10 Programs*, in "18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP'13)", Shenzhen, China, ACM, 2013, http://hal.inria.fr/hal-00761537

[35] T. YUKI, P. FEAUTRIER, S. RAJOPADHYE, V. SARASWAT. *Checking Race Freedom of Clocked X10 Programs*, arXiv, 2013, n^o arXiv.1311.4305, http://hal.inria.fr/hal-00907723