

# ACOTES: Advanced Compiler Technologies for Embedded Streaming

*Submission to the Special Issue on European HiPEAC NoE Member's Projects*

Harm Munk, Philippe Dumont,  
Marc Duranton  
NXP Semiconductors, The Netherlands

Razya Ladelsky, Dorit Nuzman,  
Ira Rosen, Uzi Shvadron, Ayal  
Zaks  
IBM Haifa Research Laboratories

Albert Cohen, Cédric Bastoul,  
Cupertino Miranda, Louis-Noël  
Pouchet, Konrad Trifunovic  
Alchemy Group, INRIA Saclay and  
LRI, Paris-Sud 11 University, France

Eduard Ayguadé, Paul Carpenter,  
Roger Ferrer, Xavier Martorell,  
Alex Ramírez, David Ródenas  
Universitat Politècnica de Catalunya,  
Barcelona, Spain

Sebastian Pop  
Compiler Performance Engineering,  
Advanced Micro Devices, Austin, Texas

Menno Lindwer  
Silicon Hive, Eindhoven, The  
Netherlands

Antoni Pop  
Centre de Recherche en Informatique, MINES  
ParisTech, France

Marco Cornero, Erven Rohou, Andrea  
Ornstein  
STMicroelectronics, Cornaredo (MI), Italy

## Abstract

Streaming applications are based on a data-driven approach where compute components consume and produce unbounded data vectors. Streaming oriented systems have become dominant in a wide range of domains, including embedded applications and DSPs. However, programming efficiently for streaming architectures is a very challenging task, having to carefully partition the computation and map it to processes in a way that best matches the underlying multi-core streaming architectures, as well as having to take into account the needed resources (memory, processing, real-time requirements, etc.) and communication overheads (processing and delay) between the processors.

These challenges have led to a number of suggested solutions, whose goal is to improve the programmer's efficiency in developing applications that process massive streams of data on programmable, parallel embedded architectures. StreamIt is one such example. Another more recent approach is that developed by the ACOTES (Advanced Compiler Technologies for Embedded Streaming) project. The ACOTES approach for streaming applications consists of compiler-assisted mapping of streaming tasks to highly parallel systems in order to maximize cost-effectiveness, both in terms of energy and in terms of design effort. The analysis and transformation techniques automate large parts of the partitioning and mapping process, based on the properties of the application domain, on the quantitative information about the target systems, and on programmer directives.

This paper presents the outcomes of the ACOTES project, a 3-year collaborative work of industrial (NXP, ST, IBM, Silicon Hive, NOKIA) and academic (UPC, INRIA, MINES ParisTech) partners, and advocates the use of the Advanced Compiler Technologies that we developed to support Embedded Streaming.

## 1. Introduction

Streaming applications which dominantly process large amounts of data have increasing demands for processing power. This demand stems from several requirements: on the one hand, the amount of processing per data element increases because of higher quality requirements of the result (e.g. video processing). On the other hand, the amount of data per unit of time also increases (e.g. higher communication speeds in wireless networks). This, in fact, calls for higher silicon efficiency, a demand that was met up to a few years ago by designing Application Specific Integrated Circuits (ASICs). The time to design such ASICs is, however, proportional to the complexity of the ASIC; as the complexity of the ASIC grows exponentially, their design becomes economically infeasible. Designers have thus shifted their focus toward programmable platforms, thereby potentially amortizing the design cost across several applications, or even application domains. Programmable platforms however were traditionally unable to meet high throughput requirements, as they were mostly designed for general purpose computation, offering limited opportunities for parallelism.

Several recent architectures do expose parallelism to the application programmer. This, however, shifts the problem of managing

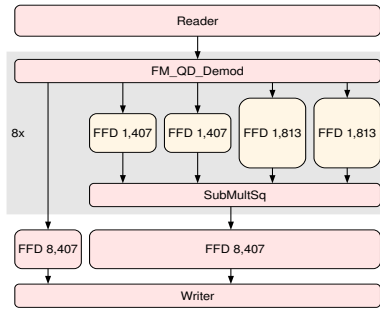


Figure 1. FMradio filters structure

complexity partly from the hardware designer to the software application developer. Exploiting available parallelism optimally requires intimate knowledge of both the application and the target platform. Automating the extraction of parallelism from sequential algorithmic descriptions has proven to be an extremely complex task in general.

In 2006, IBM, Philips (later: NXP Semiconductors), STMicroelectronics, NOKIA, INRIA and UPC initiated the ACOTES project to advance the support of application programmers in parallelising applications on highly parallel architectures. They were later joined by Silicon Hive, and MINES ParisTech. The ACOTES project concentrates on developing tools to assist the application programmer in achieving optimal parallelism. From the outset we decided to use a mainstream language (C), an existing compiler framework (the GNU Compiler Collection - GCC), focus on the data streaming application domain, and target three distinct state-of-the-art multicore architectures (Cell Broadband Engine, xStream processor, and Ne-XVP). This way we were able to concentrate our efforts on support for parallelization across several levels. Data streaming applications typically contain potential for both coarse-grain task-level parallelism across threads, fine-grain data-level parallelism residing inside nested loops of SIMD-like computations, and also memory-level parallelism to optimize data transfers.

### 1.1 Applications

The ACOTES project focuses on the data-streaming application domain. In this paper we present experimental results using three applications from this domain: FM-radio, H264 and Gamma-correction. The project, however, uses several other streaming applications to drive the developments in the project.

*FMradio.* The FM-radio application was extracted and adapted from the GNU Radio project (4). It contains about 500 lines of code. The application receives an input stream, applies a number of filters to it, and finally writes an output stream. Several of the filters apply the same transformation with different configuration parameters. The structure of the filters is shown in Figure 1. The FFD filter is the most time consuming one.

*H.264.* The H.264 application is part of the MPEG-4 standard (28). It consists of a video and audio coder and decoder, achieving high levels of compression, for improved transmission and storage of streaming media files. In this paper we use a subset of the internal algorithms to demonstrate the vectorization capabilities introduced in the GCC compiler.

*Gamma Correction.* The Gamma correction algorithm depicted in Figure 2 appears often as one of the last phases in a typical image signal processing pipeline. The current customer requirement from the Gamma correction algorithm calls for pixel throughput of about 4 cycles-per-pixel for the whole image signal processing

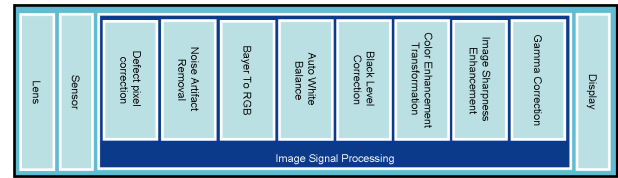


Figure 2. Image Signal Processing Algorithm Pipeline.

pipeline. As we show below, a straight-forward scalar implementation of Gamma correction takes 154 cycles-per-pixel.

The Gamma correction algorithm concentrates on a triple nested loop scanning over the pixels of a 2D image. For each pixel, it searches through an array of thresholds until it finds the threshold interval (plotted along the X-axis of Figure 3) within which the colour value  $x$  lies. The offset, gradient, and threshold of that interval are used to calculate the new pixel value. This application illustrates the potential of outer loop vectorization, developed within ACOTES.

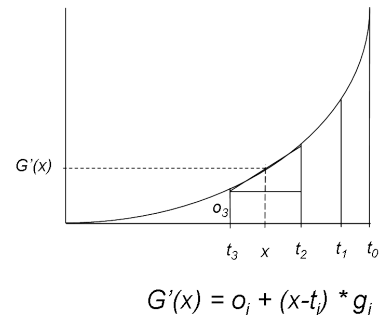


Figure 3. Gamma Correction Filter.

### 1.2 Compiler Structure

GCC, the GNU Compiler Collection is a compiler framework supporting several mainstream programming languages as well as a large array of target processors. The framework is actively maintained by the GCC project of the Free Software Foundation, supervised by the GCC steering committee.

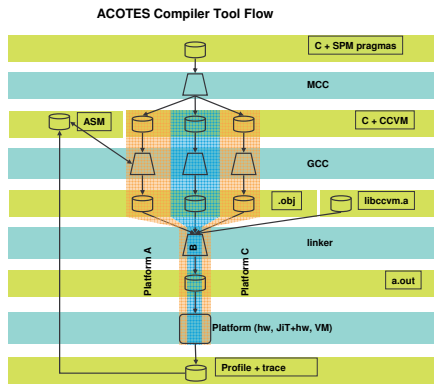
The project consists of a mainstream development part, which results in regular releases, and several development parts often dubbed development branches. The mainstream development is closely supervised by the GCC steering committee delegating responsibilities to leading developers. Results from development branches are incorporated into the mainstream only after scrutiny and approval by designated maintainers. This setup allows for the simultaneous development of a production quality compiler and a set of new and experimental developments.

The structure of GCC is that of a traditional compiler: a front-end for language parsing and semantics checking, a “middle-end” for target-independent optimisations, and a back-end for code generation and target-specific optimisations. It is this combination of support for mainstream languages and targets, industrial mainstream quality, the ability to experiment freely in parallel to the main line of development without interference, a vibrant supporting community and its adoption by a vast number of users that makes GCC an ideal choice for a project such as ACOTES.

### 1.3 ACOTES Project Structure

The ACOTES project is divided into a set of subtopics:

- **Abstract Streaming Machine.** In order to target several distinct architectures, we developed an abstract machine model that captures the streaming characteristics of the platforms in a common model called the Abstract Streaming Machine, or ASM.



**Figure 4.** The overall ACOTES compiler flow: MCC is the Mercurium C-compiler, CCVM is the C Virtual Machine, also called the ACOLib. The other abbreviations are explained in the text.

- **Streaming Programming Model.** In order to use the C-language as the starting point, we developed a number of extensions that express parallelism opportunities in the source, collectively referred to as the Streaming Programming Model, or SPM. As an extra requirement, these language extensions had to be neutral with respect to the core C language. The pragma mechanism is ideal for this purpose. This decision, however, surfaced a shortcoming of GCC, which is not very well suited for experimenting with pragmas in the front-end. This prompted the project to turn to Mercurium (2) as a prototyping tool for implementing the SPM.
- **Middle-end optimisations.** Streaming applications are characterised by nested loops of computations. Often, considerable performance improvements can be gained by applying optimisations to these loops. The project concentrated its effort on the polyhedral model for high-level loop optimisations in the middle-end.
- **Back-end optimisations.** Modern architectures often contain SIMD-like instructions, also called vector instructions. As these instructions are target specific, opportunities for applying them are detected in the back-end of the compiler.
- **Code generation.** Although the code generation is specific for each architecture and outside the scope of the ACOTES project, we decided to concentrate our effort on common back-end algorithms and on common back-end formats.

Apart from these topics, the project has a separate activity to disseminate the project results, part of which is maintaining contact with the GCC-project and HiPEAC network of excellence communities.

This paper describes the implementation of the parallelising support developed in the ACOTES target tool-chain and presents experimental results. Section 2 describes related work, sections 3 and 4 present the Abstract Streaming Machine (compilation target) and the Streaming Programming Model in the compiler front end (extensions to the C-language). Section 5 presents the automatic loop nest optimisations and vectorisation, and the interaction between them. Section 6 describes the target platforms and the code generation phase of the compiler for these architectures. Section 7 presents results, and section 8 concludes.

## 2. Related Work

Since the ACOTES project takes a holistic approach to parallel stream-programming, it spans over the whole flow of compilation down to the runtime system. It has connections with a huge list of related work. This section compares our approach with the most closely related progresses and results in the field of high-performance embedded computing.

StreamIt is a long running project with a source language, publicly available compiler, and benchmark suite. The StreamIt (50) language imposes a hierarchical structure on the program, which must be composed of filters, pipelines, split-joins and feedback loops. It requires the developer to structure the program into separate work functions per filter, in contrast to using pragmas. The StreamIt compiler (26) targets the Raw Architecture Workstation, symmetric multicore architectures, clusters of workstations, and the Cell BE, where aggressive task-level optimizations are performed automatically (31). There is no general purpose machine model similar to the ACOTES ASM, and the ACOTES SPM is much more expressive than the cyclostatic data-flow model of computation underlying StreamIt (12), while still allowing compilation-time task-level optimizations (19; 44) Compared to StreamIt, our approach involves a tight coupling of task- and loop-level optimizations, allowing for more relevant decisions to be taken about adaptation of synchronization, communications, levels of parallelism and locality. The optimization relies on iterative search, we are able to find interesting tradeoffs inaccessible to compilation flows separating those problems into different representations and passes (46). We will illustrate this design on the interplay between task-level and loop-level optimizations, including automatic vectorization.

Sequoia is a well known data parallel language exploiting the structure of data-centric algorithms (18). It severely restricts the model of computation to hierarchical fork-join parallel sections, but allows the programmer to state data affinity to portions of iteration space. For each hardware platform, the application programmer must supply a mapping that takes the “abstract” hierarchy defined in the application, and assigns pieces of it onto specific hardware. This approach requires more effort from the application provider and requires them to learn the memory characteristics of each hardware platform, but it is certainly a pragmatic solution that could be added as an optional feature of the ACOTES programming model.

The new specification of OpenMP (43; 14), version 3.0, supports task parallelism using the new `task` directive. This directive specifies that the serial code within it can be executed by another thread inside the scope of a `parallel` region. In OpenMP, every time the task directive is reached a new task is created to execute its body. In the ACOTES SPM, all the inner tasks are created once when the `taskgroup` directive is reached, and a value is sent on each input stream each time the task directive is reached. This is a form of synchronization that does not exist in the OpenMP 3.0 proposal. However, there are other proposals for OpenMP that add synchronization between threads. Gonzalez et al. (24; 25) propose three new directives: PRED, SUCC and NAME. The NAME directive labels a worksharing, and this label can be used by PRED and SUCC directives, to specify synchronization. Another approach using annotated C is Cell Superscalar (CellSs) (3), which uses a `task` directive to express what are the inputs and outputs at a function level. Each time the function is called, a new task is created and the runtime system takes care of the possible dependencies it may have with other tasks.

StarPU features a stream-oriented model of computation, yet does not involve any language extension (38). Based on multi-versioned kernels, it automates the dynamic balancing and mapping of tasks and data over heterogeneous, accelerator-centric par-

allel architectures. It would be an interesting target to map code generated by GCC and to port the ACOTES runtime system.

StreamRoller (30) is a stream compiler for the Cell Broadband Engine, which uses the SGMS algorithm to split stateless kernels, partition the graph, and statically schedule. Task fission and fusion is translated into an Integer Linear Programming (ILP) problem, which is solved using the commercial CPLEX solver (27).

Gedae is a proprietary graphical programming environment for streaming signal processing applications in the defense industry. Unlike ACOTES, the developer specifies the mapping of the program onto the target, and the compiler generates the executable according to this mapping (34).

We selected GNU Radio as our motivating example (4). It is a framework developed in C++ and Python. GNU Radio allows to express graphs of filters and connections described using Python. Filters are usually constructed as C++ classes. GNU Radio comes with its own task scheduler and the system can be deployed on multiple architectures, even including FPGAs. GNU Radio provides more than 100 different basic blocks that can be combined to achieve the goal of the application. New blocks may be added to add new functionality. Both StreamIt and GNU Radio are designed for signal processing applications, and require the program to be written specifically in terms of streaming blocks.

### 3. ASM

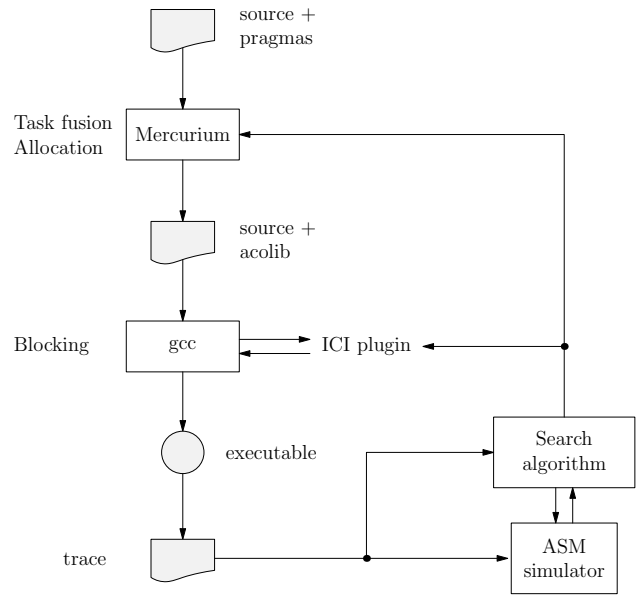
The Abstract Streaming Machine (8; 9) is the compiler’s description of the target multiprocessor system. The ASM defines the search space for the compiler’s partitioning and allocation algorithm by specifying the system topology and performance characteristics, and providing constraints on the mapping; e.g. memory sizes.

We have implemented the ASM in a coarse-grain simulator, which estimates the performance of a candidate mapping on the given target without requiring compilation or execution. In addition to improving the speed of the compiler, the simulator allowed us to start work on the search algorithm before the compiler’s transformation infrastructure was complete, and is repeatable because there is no experimental error.

Figure 5 shows the structure of the ACOTES compiler, including the ASM simulator. The compilation flow is iterative, with a heuristic search algorithm determining a candidate mapping, which is compiled using Mercurium (2) and GCC. The mapping is provided to GCC through a plugin using the Interactive Compilation Interface (ICI) (20).

#### 3.1 The ASM machine model

The topology of the target platform is given by an unstructured bipartite graph  $H = (V, E)$  where  $V = P \cup M \cup I$  is the set of vertices, a disjoint union of processors,  $P$ , and memories,  $M$ , in one partition and interconnects,  $I$ , in the other. The edges,  $E$ , serve only to define the topology. Figure 6 shows two example targets: (a) a Cell processor, and (b) a four-core shared memory machine. Each processor, interconnect, and memory is defined using the parameters summarized in Figure 7, and described in detail below. Figure 7(a) and 7(b) give the machine descriptions for the example targets, measured on Cell and estimated for a 4-processor SMP. The ASM defines the machine characteristics that are visible to software, including the ACOlib runtime system, so it may not exactly match the physical hardware. For example, the Operating System in a Playstation 3 makes only six of eight SPEs available to software, and the mapping from virtual to physical core is not known. We assume that any processors available to the stream program will not be time-shared with other applications while the stream program is running.



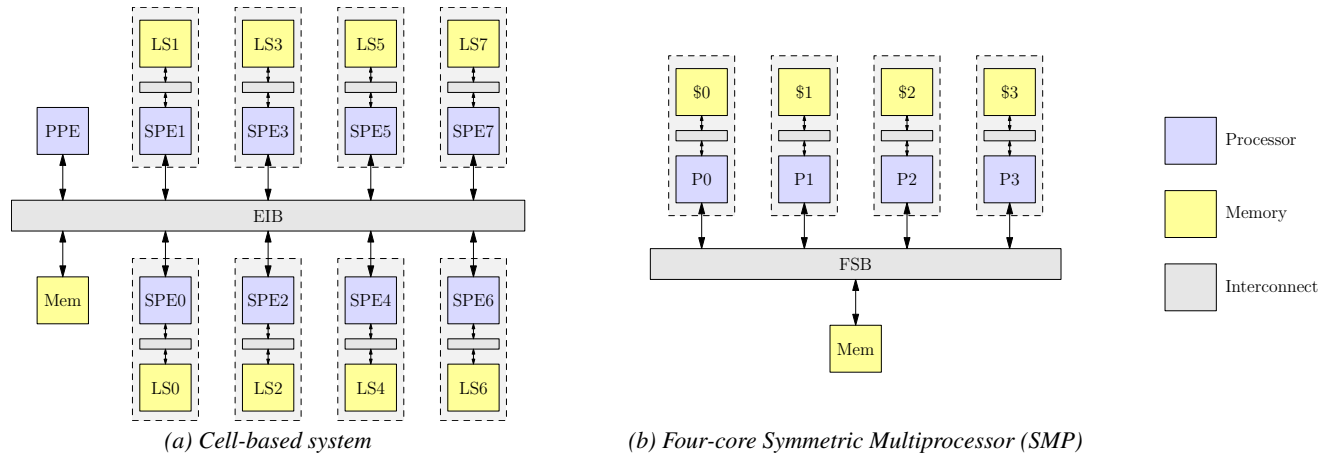
**Figure 5.** The feedback loop of the ACOTES compiler: a heuristic search algorithm controls Mercurium and GCC

Each processor is defined using the parameters in Figure 7(a). The ASM supplements the compiler’s back-end machine description, so the details of the processor’s ISA and micro-architecture are not duplicated in the ASM. The `addressSpace` and `hasIO` parameters provide constraints on the mapping, but are not required to estimate the performance of a valid mapping. The former defines the local address space of the processor; i.e. which memories are directly accessible through ordinary load-store instructions, and where they appear in virtual memory, and is used to place stream buffers. The latter defines which processors can perform system IO, and is a simple way to ensure that tasks that need system IO are mapped to an appropriate processor.

Each memory is defined using the parameters in Figure 7(b). The latency and bandwidth figures may be used by the compiler to refine the estimate of the execution time of each task, but are not used directly by the coarse-grain simulator. The memory sizes are used to determine where to place communications buffers, and provide constraints on the blocking factors.

Each interconnect is defined using the parameters shown in Figure 7(c). The graph topology is given by the `elements` parameter, which lists the processors and memories that are adjacent to the bus. Each interconnect is modelled as a bus with multiple channels, which has been shown to be a good approximation of the performance observed in practice when the processors and memories on a single link are equidistant (23). Each bus has a single unbounded queue to hold the messages ready to be transmitted, and one or more channels on which to transmit them. Streams are statically allocated onto buses, but the choice of channel is made at runtime. The `interfaceDuplex` parameter defines for each processor or memory whether it can simultaneously read and write on different channels. The bandwidth and latency of each channel is controlled using four parameters: the start latency ( $L$ ), start cost ( $S$ ), bandwidth ( $B$ ) and finish cost ( $F$ ). In transferring a message of size  $n$  bytes, the latency of the link is given by  $L + S + \lceil \frac{n}{B} \rceil$  and the cost incurred on the link is  $S + \lceil \frac{n}{B} \rceil + F$ . This is a natural model for distributed memory machines, and is equivalent to the assumption of cache-to-cache transfers on shared memory machines.

The ASM simulator assumes that the only significant traffic on an interconnect is the transfer of messages related to streams.



**Figure 6.** Topology of two example targets for the ASM

Hence each processor should have some private memory — either a local store or a cache. If it is a local store, the compiler must allocate the stream buffers in this local store. If it is a cache, the ASM assumes that it is sufficiently effective so that the cache miss traffic on the interconnect is low.

Hardware routing is controlled using the `interfaceRouting` parameter, which defines for each processor whether it can route messages from this interconnect onto another interconnect that it is adjacent to. Each entry can take the value `storeAndForward`, `cutThrough`, or `None`.

### 3.2 The ASM program model in the simulator

The coarse-grain simulator models the stream program as a directed graph  $G = (T, S)$  where  $T$  is the set of vertices representing tasks and  $S$  is the set of edges representing streams. The graph does not have to be acyclic, but it must be connected.

Note that a task may have irregular data-dependent behaviour. We therefore divide tasks into *subtasks*, which are the basic unit of sequencing. A subtask pops a fixed number of elements from each input stream and pushes a fixed number of elements into each output stream. In detail, the work function for a subtask is divided into three consecutive phases: first, the *acquire phase* calls `Iport_acquire` or `Oport_acquire` for each stream to obtain the next set of full input buffers and empty output buffers. Second, the *processing phase* works locally on these buffers, and is modelled using a fixed or normally-distributed processing time. Finally, the *release phase* calls `Iport_pop` to discard the input buffers, and calls `Oport_push` to send the output buffers, releasing the buffers in the same order they were acquired. This three-stage model is not a fundamental requirement of the ASM, and was introduced as a convenience in the implementation of the simulator, since the ACOTES compiler will generate subtasks of this form.

Streams are defined by the size of each element, the location and the length of the distinct producer and consumer buffers (distributed memory) or the single shared buffer (shared memory). These buffers do not have to be the same length. Streams are point-to-point, so each stream has exactly one producer task and one consumer task, but those tasks may access the same stream from more than one subtask.

### 3.3 Definition and sequencing of irregular tasks

The coarse-grain simulator uses the sequential semantics of the SPM program to control the sequencing of subtasks in the stream program. A task is controlled by its *subtask tree*, which is built up

from subtasks, *If* nodes and *While* nodes. Each *If* or *While* node is associated with a similar statement in the sequential program.

When the simulator is to be executed in the trace-driven mode, the executable is instrumented to record the outcome each time a control statement is executed. A control statement is an *if* or *while* statement in the original SPM program that controls one or more subtask trees. The resulting sequence of outcomes is known as a control variable, and takes the values 1 or 0 for an *if* statement, or the non-negative iteration count for a *while* statement. When the simulator is used in the trace-driven mode, the program model is driven by the set of control variables taken from the trace.

The set of control variables may be reused with a different partition or allocation. It usually cannot be reused with a different blocking factor, or after compiler transformations such as loop interchange or distribution.

## 4. SPM and the Front-End Compiler

The Streaming Programming Model (SPM) designed in the context of the ACOTES project is implemented using extensions to the C language. It consists of a set of *pragmas* extending the serial code semantics. The main requirements of the SPM are: to be easy to learn, easy to use and reuse, and to support task and data parallelism. We think that OpenMP (43) can therefore serve as a good basis to develop our SPM: OpenMP can be learned, applied, and tested incrementally, which is convenient for new programmers in the streaming field.

### 4.1 Elements in the SPM

The SPM adds three basic elements to OpenMP: *streaming tasks*, *streams*, and *ports*. These elements are grouped together in *taskgroups*, the only compound abstraction. Applications are represented in the SPM as multiple tasks connected via point-to-point data streams. Each task may be viewed as an independent process with all its data being private. Communication and synchronization among tasks happens only via streams. A stream is directed, and we refer to its two end points (*ports* from now on) from the point of view of the task, so that the producer has an *output port* to generate data into the stream, and the consumer has an *input port* to read data from the stream. The two ends are permanently connected together. The consumer task blocks when it tries to read from an empty input stream, and the producer blocks when it tries to write to a full output stream.

Using the scheme described above, task parallelism is supported in a traditional way with the addition of having communication

Parameter	Type	Description	Cell	SMP
name	String	Unique name in platform namespace	'SPE $n$ '	'CPU $n$ '
clockRate	Fixed pt.	Clock rate, in GHz	3.2	2.4
hasIO	Bool	True if the processor can perform IO	False	True
addressSpace	List	List of the physical memories in the system that are addressable by this processor and their virtual address	[('LS $n$ ', 0)]	[('Mem', 0)]
pushAcqCost	Int	Cost, in cycles, of acquiring a producer buffer (before waiting)	448	20
pushSendFixedCost	Int	Fixed cost, in cycles, of pushing a block (before waiting)	1104	50
pushSendUnit	Int	Number of bytes per push transfer unit	16384	0
pushSendUnitCost	Int	Incremental cost, in cycles, to push pushUnit bytes	352	0
popAcqFixedCost	Int	Fixed cost, in cycles, of popping a block (before waiting)	317	50
popAcqUnit	Int	Number of bytes per pop transfer unit	16384	0
popAcqUnitCost	Int	Incremental cost, in cycles, to pop popUnit bytes	0	0
popDiscCost	Int	Cost, in cycles, of discarding a consumer buffer (before waiting)	189	20

(a) Definition of a processor

Parameter	Type	Description	Cell	SMP
name	String	Unique name in platform namespace	'LS $n$ '	'Mem'
size	Int	Size, in bytes	262144	2147483648
clockRate	Fixed pt.	Clock rate, in GHz	3.2	0.4
latency	Int	Access latency, in cycles	2	4
bandwidth	Int	Bandwidth, in bytes per cycle	128	8

(b) Definition of a memory

Parameter	Type	Description	Cell	SMP
name	String	Unique name in platform namespace	'EIB'	'FSB'
clockRate	Fixed pt.	Clock rate, in GHz	1.6	0.4
elements	[String]	List of the names of the elements (processors and memories) on the bus	['PPE', 'SPE0', ..., 'SPE7']	['CPU0', ..., 'CPU3']
interfaceDuplex	[Bool]	If the bus has more than one channel, then define for each processor whether it can transmit and receive simultaneously on different channels	[True, ..., True]	[False, ..., False]
interfaceRouting	[Enum]	Define for each processor the type of routing from this bus: <i>storeAndForward</i> , <i>cutThrough</i> , or <i>None</i>	[None, ..., None]	[None, ..., None]
startLatency	Int	Start latency, $L$ , in cycles	80	0
startCost	Int	Start cost on the channel, $S$ , in cycles	0	0
bandwidthPerCh	Int	Bandwidth per channel, $B$ in bytes per cycle	16	16
finishCost	Int	Finish cost, $F$ , in cycles	0	0
numChannels	Int	Number of channels on the bus	3	1
multiplexable	Int	False for a hardware FIFO that can only support one stream	True	True

(c) Definition of an interconnect

**Figure 7.** Parameters of the Abstract Streaming Machine and values for two example targets (measured on Cell and estimated for a four-core SMP)

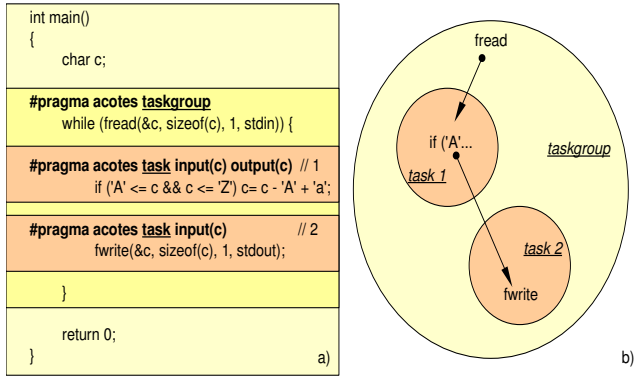


Figure 8. Example program (a), and representation of tasks (b)

channels between tasks. Data parallelism is supported through the ability of replicating a task into a number of instances, allowing to run each instance in a different core on disjoint input and/or output data.

There are several distinctions between the execution model of the SPM and that of OpenMP, namely:

- In the SPM, streaming tasks are created all at once when a *taskgroup* is entered. This is contrary to OpenMP where a thread creates each task in a parallel region upon encountering it dynamically.
- Tasks in SPM are permanent, meaning that they are alive while there is input data for them to process. This implies the automatic generation of an implicit loop: *while (there-is-input-data) { ... }*, enclosing the code in the body of the task, contrary to what is done in OpenMP. Only when the input streams have no more data, can the task finish. A *taskgroup* ends when all its tasks have finished.
- Contrary to OpenMP which supports access to shared data, any data accessed by a task must be either private or acquired through the input streams of the task. The SPM defines specific situations where global data can be accessed through a well-defined interface.

#### 4.2 Streaming execution model

An SPM program start executing serially as a single process. Upon entering a *taskgroup*, tasks are created and the program starts processing data in streams. Figure 8a shows a simple example, which converts an input stream read from *stdin* to lower case, and then writes the resulting stream to *stdout*. Figure 8b shows the scheme of tasks built using the SPM pragmas. As can be observed in the drawing, the *taskgroup* is used as a container for the annotated tasks. Arrows represent streams along the direction in which data circulates.

Observe that extra clauses are attached to the SPM task pragmas to express data transfers through the streams. The clauses *input* and *output* receive one variable for each input and output stream that should be established, respectively. A task can receive input data either from the taskgroup or from a previous task having the same variable designated as output. A task can send output data either to a subsequent task having the same variable designated as input, or to the taskgroup. This way, the taskgroup itself is always used when there is no previous or subsequent task having the corresponding variable as output or input.

The SPM also supports the exploitation of data parallelism. In the streaming context, this implies processing chunks of data taken from streams in parallel. This is accomplished by replicating an

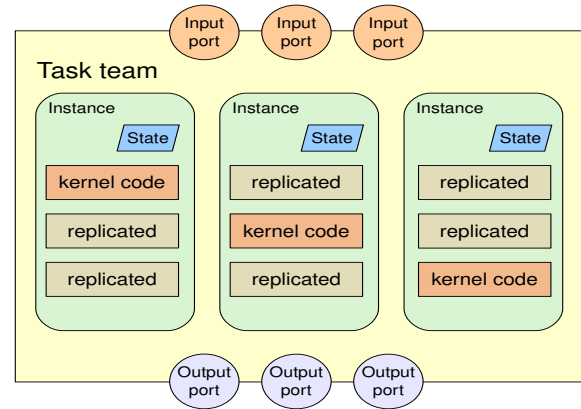


Figure 9. Representation of a team of tasks

SPM task into multiple instances (called a *team*), where each instance processes a different set of input data chunks. Figure 9 represents the execution of a team consisting of three instances. Tasks can have many inputs and many outputs. When splitting a task into a team of instances, each input stream can be split so that data chunks are distributed among the instances, or it can be replicated so that all instances will receive all data chunks. The team contribute to the output streams as a single task, where each output chunk is generated by a single instance at every task iteration. The runtime system implements this combined contribution by keeping track of which instance is the leader at every iteration. See in Figure 9 how the complete kernel code (updating state and contributing output) is executed by a single instance at a time, while the code needed only to update the state of each instance is executed in a replicated way.

#### 4.3 Front-end compiler and runtime library (ACOLib)

The transformations outlined above are implemented in the *Mercurium C Compiler* (2; 9), which is a source-to-source compiler, generating C code. We use Mercurium to translate the full set of SPM pragmas and clauses into calls to a runtime library which supports task and stream management. The resulting C output is then compiled with GCC, and linked with the runtime library to obtain the final binary file.

The initial approach for a runtime system in the ACOTES project was to implement a simple library (ACOLib) supporting the functionality needed for tasking and streams. We developed an implementation of ACOLib for generic SMP environments, but designed it in such a way that the library could also work with possibly distributed memory architectures and local stores.

ACOLib supports streaming execution in SMP environments, as shown in Figure 10a which depicts a snippet of the code generated for the taskgroup example presented in Figure 8a. Observe how two tasks are initialized, their ports created and connected, and then they are started. The two tasks are alive while the taskgroup reads characters from the file and sends them to *Task 1*. Figure 10b shows the code resulting from the transformation of *Task 1* of the same example. This task reads characters from its input port, processes them, and writes the result onto its output port. The task will remain alive as long as there are characters available in its input stream.

<pre> task_t task1; task_init(&amp;task1, task1_outline, 0); task_t task2; task_init(&amp;task2, task2_outline, 0); task_oport(task0, 0, sizeof(char), 1, 1); task_iport(task1, 0, sizeof(char), 1, 0 + 1, (void *) 0, 0); task_oport(task1, 1, sizeof(char), 1, 1); task_iport(task2, 0, sizeof(char), 1, 0 + 1, (void *) 0, 0); port_connect(task1, 1, task2, 0); port_connect(task0, 0, task1, 0); task_start(task1); task_start(task2); while (fread(&amp;c, sizeof(c), 1, stdin)) {     oport_acquire(0, 1);     memcpy(oport_peek(0), &amp;c, sizeof(char));     oport_push(0, 1); } task_wait(task2); task_wait(task1); </pre> <p>a) Taskgroup</p>	<pre> void task1_outline(task_t __task) {     trace_instance_begin();     char c;     iport_acquire(0, 1);     oport_acquire(1, 1);     while (task_allopen())     {         memcpy(&amp;c, iport_peek(0), sizeof(char));         if ('A' &lt;= c &amp;&amp; c &lt;= 'Z')             c = c - 'A' + 'a';         memcpy(oport_peek(1), &amp;c, sizeof(char));     }     iport_pop(0, 1);     oport_push(1, 1);     iport_acquire(0, 1);     oport_acquire(1, 1); } task_close(); </pre> <p>b) Task 1</p>
---	--

Figure 10. Code generated from the taskgroup example

#### 4.4 Integration in OpenMP and GCC

In parallel to developing the full set of SPM pragmas and their support by Mercurium and ACOLib, we developed an alternative proposal suggesting a minimal extension (44) of the OpenMP3.0 standard (43). This approach leverages the knowledge of OpenMP thereby improving the learning curve while preserving, to a certain extent, the semantic of the SPM pragmas. The additional compiler support required by this extension is implemented directly in GCC and its runtime library libGOMP.

In order to provide stream programming support in OpenMP, the minimal necessary extension is to allow the use of `lastprivate` clauses on task constructs. This extension does not change the existing semantic of the `lastprivate` clause. The semantic of `firstprivate` and `lastprivate` clauses is very close to SPM’s input and output clauses. The `firstprivate` clause corresponds to data that is consumed by the task (flows in), while the `lastprivate` clause corresponds to data that is produced by the task (flows out). The explicit knowledge about data flow between tasks helps the compiler’s static analysis and facilitates generating stream communication between tasks.

To illustrate the OpenMP extension as well as the similarity of this approach with the SPM, we propose in Figure 11 an implementation of the example code of Figure 8a using OpenMP3.0 extended annotations. We implemented this support for streams in GCC’s libGOMP runtime library (45), which also simplifies the toolchain by removing the dependence on the source to source *Mercurium C Compiler*.

#### 4.5 MSF: Multicore Streaming Framework

As described in subsection 4.3 above, the ACOLib runtime library developed to support tasks and streams was implemented for generic SMP environments, i.e. with shared memory. In order to apply our streaming framework to distributed memory architectures such as that of the Cell/B.E., we are making use of an underlying Multicore Streaming Framework (MSF). On one hand, MSF provides efficient support for code management, as the code that the tasks execute is preloaded in advance of its execution. And on the other hand, it provides efficient support for data transfers between tasks, which may be running on the same or on different processors, as well as with shared or distributed memory. The abstraction layer implementing such data transfers resembles that of ACOLib, offering stream-based communication.

```

int main()
{
    char c;
#pragma omp parallel
    {
#pragma omp single
        {
            while (fread (&c, sizeof (c), 1, stdin))
            {
#pragma omp task firstprivate (c) lastprivate (c)
                if ('A' <= c && c <= 'Z')
                    c = c - 'A' + 'a';
            }

#pragma omp task firstprivate (c)
                fwrite (&c, sizeof (c), 1, stdout);
        }
    }
    return 0;
}

```

Figure 11. Implementation of the example of Figure 8a using OpenMP extended annotations.

The use of MSF in ACOTES augments the general tool-chain of ACOTES which starts with Mercurium translating SPM pragmas into ACOLib calls. Instead, for the Cell/B.E., it makes direct use of MSF facilities for setting up tasks in remote processors and establishing stream communications among them. Furthermore, MSF also provides the runtime management and scheduling of task executions according to available data in the heterogeneous and distributed environment of the Cell/B.E.

MSF provides a generic programming model for parallel platforms and an API for programming directly using its facilities. Although its first implementation is on the Cell/B.E., MSF can be implemented on any programmable platform. Tasks can be directed to various processors based on their processing characteristics. On the Cell/B.E the target processors are either PPEs or SPEs. The tasks are then compiled by a back-end compiler that deals with standard sequential code, and vectorization can be applied to achieve better performance on the target processor. Once available on a specific platform the same front-end compilation techniques can be performed independently of the underlined architecture. Using the information provided by the ASM, applications can be optimized for each platform by changing runtime parameters that are used by MSF.

#### 4.6 Silicon Hive’s ANSI-C Language Extensions, underpinning SPM

The ACOTES Streaming Programming Model defines tasks and blocking ports for communication between these tasks.

Silicon Hive, one of the ACOTES partners, already developed and used an in-house C-compiler (HiveCC) prior to the project. In the course of the project, Silicon Hive adapted HiveCC to the developing insights of the project, providing feedback to the other partners. What was found during use of the blocking communication ports of ACOTES, was that synchronization (i.e. the blocking behaviour) and data communication should take place at different granularities. This allows the compiler to expose parallelism between communication and computation. The ANSI-C language extensions for synchronisation, as discussed below, were the result. We show how these synchronisation mechanisms are being used in the context of the gamma correction algorithm that was introduced in Section 1.1. In this context, we will also make use of other language extensions, such as: attributes, built-in types, function inlining, operator overloading, and register structures, as illustrated in Figure 12 (for reference, the unoptimized version, which is used as an input to the ACOTES auto-vectorizer is shown in Fig-



```

HIVE_TERNARY_OPERATOR( ?, OP_vec_mux
                      , tvec, tflags, tvec, tvec )
HIVE_BINARY_OPERATOR ( >>, OP_vec_asrrnd
                      , tvec, tvec, tvec )
HIVE_BINARY_CALL     ( >, static inline, isp_vec_gt_c
                      , tflags, tvec, int )
#define LINESZV LINESZ/VECSZ
SYNC_WITH(0) tvec MEM(VMEM) inBuf[2*LINESZV];
SYNC_WITH(1) tvec MEM(VMEM) outBuf[2*LINESZV];
...
// Initialize correction coefficients:
int ta[SN]={...}, oa[SN]={...}, ga[SN]={...};
// double buffering; pre-fetch & sync with buf 1:
signal_inputFormatter() SYNC(0);
for( line=0; line<LINES_PER_FRAME; line++ ) {
    int buf=line&1; // determine buffer
    // double buffering; pre-fetch next line:
    signal_inputFormatter() SYNC(0);
    // double buffering; wait for current line:
    wait_inputFormatter() SYNC(0);
    for( c=0; c<LINESZV; c++ ) {
        tvec x = inBuf[buf*LINESZV+c];
        int t=thrh[0], o=offs[0], g=grad[0];
        for( i=0; i<SN-1; i++ ) {
            tflags flag = x>ta[i]; // overloaded '>'
            // overloaded vector '?' operator:
            o = flag ? o : OP_vec_clone(oa[i+1]);
            g = flag ? g : OP_vec_clone(ga[i+1]);
            t = flag ? t : OP_vec_clone(ta[i+1]);
        }
        # pragma hivecc unroll
    }
    // overloaded vector '>>' operator:
    outBuf[buf*LINESZV+c] = o+(x-t)*g>>SCALE;
    # pragma hivecc unroll=6, pipelining=0
}
// Signal output buffer full, sync with output buf.
signal_DMA() SYNC(1);
// Wait until DMA finished previous line.
if(line) wait_DMA() SYNC(1);
}

```

**Figure 12.** Gamma-correction algorithm optimized with SiliconHive’s SPM

ure 18). HiveCC supports several pragmas by which the programmer can steer the compiler to reveal parallelism. HiveCC also provides built-in preprocessor definitions which enable the programmer to develop code independent of the underlying system architecture.

Taken together, these language extensions are needed to efficiently implement SPM on SiliconHive processors.

*Type attributes* specify aspects of variables such as location, access routes, and synchronization with volatile operations. *Expression attributes* allow the user to specify the relative timing of associated statements, aliasing relations between variables, and the mapping of code on specific function units. MEM and SYNC\_WITH in Figure 12 are example for such attributes. These define input/output as mapped on a specific memory and with which side-effect statements their accesses need to be synchronized.

*Built-in Types* allow HiveCC to generate non-ANSI-C operations on non-ANSI-C types. These types result from processors being defined with datapaths of arbitrary width (i.e. not restricted to the widths of ANSI-C datatypes). The additional types are available as signed or unsigned integers of processor-defined width and vectors with processor-defined numbers of elements and element widths. If the associated operators have been overloaded (see below), these types may be used in arithmetic expressions. In the optimized gamma correction algorithm in Figure 12, pixel colour components are mapped onto the vector type tvec. In the case of ISP2300 this is a vector of 32 elements, each element being 16 bits wide.

*Custom Operations* can be called as intrinsic functions, which is needed when the overloading mechanism cannot distinguish operations based on signature alone. In fact, all operations defined on SiliconHive processors may be called as intrinsic functions. Besides intrinsics for regular and custom operations, HiveCC also provides a number of *Pseudo Operations*. These operations are interpreted by the simulator, e.g. to print cycle counts on the simulation console.

*Operator Overloading* allows the programmer to assign normal operator symbols to functions, including intrinsic functions for specific processor operators. HiveCC supports overloading of unary, binary, ternary, and assignment signatures. The optimized Gamma correction code provides three examples of overloading (see Figure 12). By applying overloading to inline functions, the programmer can develop machine-independent code without loss of performance. This is illustrated in the following code example:

```

inline static tvec fmux(tvec c,tvec t,tvec e) {
# if HAS_mux
    return( OP_mux(c,t,e) );
# else
    return c?t:e;
}
HIVE_TERNARY_CALL(?,static,fmux,tvec,tvec,tvec,tvec)

```

*Register Structures* are an extension to the ANSI-C register qualifier. By adding this qualifier to a newly defined type, the programmer indicates that all elements of the type are to be mapped onto registers.

In conclusion, the new synchronisation mechanism is needed to allow data to be communicated at e.g. pixel granularity, while synchronisation takes place at line granularity. This way, rather than keeping all volatile loads and stores together, the programmer is free to communicate data when it becomes available (streaming) and the compiler is free to schedule the resulting loads and stores within the constraints imposed by the synchronisation mechanism, exposing maximal parallelism.

The other type attributes and built-in types are needed in this context, because they allow the compiler to reason about the many different non-ANSI-C datatypes that embedded processors may support. Lastly, the overloading mechanism supports the programmer in reasoning about his code, which is particularly important when communication and computation need to be alternated.

## 5. Compiler Middle-End and Loop-Level Optimizations

Through the SPM, the programmer exposes much of the pipeline and data parallelism in stream computations. The goal of the ACOTES project is to minimize the burden of manually adapting a stream computing program to a new architecture. This adaptation is partly managed by the runtime system for the SPM, and partly by the compiler middle-end. Program transformations are indeed necessary to adjust the synchronization grain, and to tune the exploitation of the architecture’s memory hierarchy w.r.t. the temporal and spatial locality of the streaming program. One typically splits such transformations into task-level and loop-level optimizations. Both address parallelism, locality and specialization, but generally at different levels (the levels of the memory hierarchy and the levels of parallelism in the target). In the following, we will focus on the loop-level optimizations, although our design is extensible to task-level optimizations such as static task pipelining, fusion and blocking (26). Such an extension would require an adequate intermediate representation of the task-level data flow; this is still work in progress, based on Feautrier’s proposal to extend the polyhedral representation to the full network of tasks (19). We will conclude

this section with the study of some tradeoffs between the exploitation of thread-level parallelism, fine-grain SIMD parallelism (e.g., vectors) and memory locality. Those tradeoffs are particularly important when optimizing a streaming application for multiple levels of parallelism, and considering power efficiency and compute density metrics.

### 5.1 Loop-Nest Optimizations

Loop nest optimizations are important compiler transformations to fully exploit the features of a given architecture. Current compiler techniques using syntax-based representations and hard-wired pattern matching techniques are not able to provide peak performance on complex architectures and/or with multiple on-chip processing units.

In order to enable long and/or complex sequences of program transformations, we rely on a powerful algebraic program representation called the *polyhedral model*, where a sequence of transformations is represented as a single affine mapping function.

In scientific and engineering applications, most of the execution time is spent in nested loops. The *polyhedral model* views a dynamic instance (iteration) of each program statement as an integer point in a well defined subspace called *polyhedron*. Basing on this representation a dependence graph is built, which represents dependences among pairs of statement execution instances (iterations in loops).

The polyhedral model is classically applicable to *static control parts* (SCoP), that is loop nests in which the data access functions and loop bounds are affine combinations of surrounding loop iteration variables and *global parameters*.

*Polyhedral representation of programs.* Each dynamic instance of a statement  $S$  is denoted by a pair  $(S, \mathbf{i})$ , where  $\mathbf{i}$  is iteration vector which contains values for the loop indices of the surrounding loops, from outermost to innermost. If loop bounds are affine expressions of outer loop indices and global parameters (usually, symbolic constants representing problem size) then the set of all iteration vectors around statement  $S$  might be represented by polytope  $\mathcal{D}_S$  which is called *iteration domain* of the statement  $S$ . Let  $\mathbf{g}$  be the vector of *global parameters*. The iteration domain set  $\mathcal{D}_S$  is defined by

$$\mathcal{D}_S = \{ \mathbf{i} \mid \mathcal{D}_S \times (\mathbf{i}|\mathbf{g}|1)^t \geq \mathbf{0} \}$$

*Polyhedral dependences.* Dependences in a SCoP are represented as a Dependence Graph (DG). DG is directed multigraph  $\text{DG} = (V, E)$  where each vertex represent a statement and each edge  $e^{S_i \rightarrow S_j} \in E$  from  $S_i$  to  $S_j$  represents a dependence polyhedron from dynamic instance of  $S_i$  to dynamic instance of  $S_j$ . The dependence polyhedron is a subset of cartesian product of iteration domains  $\mathcal{D}_{S_i}$  and  $\mathcal{D}_{S_j}$ . Dependence polyhedron for edge  $e$  is denoted as  $\mathcal{P}_e$ .

*Access functions.* For each statement  $S$ , we define two sets  $\mathcal{W}_S$  and  $\mathcal{R}_S$  of  $(M, f)$  pairs, each pair representing a reference to variable  $M$  being written or read in the statement  $S$ ;  $f$  is the *access function* mapping iterations in  $\mathcal{D}_S$  to memory locations in  $M$ .  $f$  is a function of loop iterators and global parameters. The access function  $f$  is defined by a matrix  $F$  such that

$$f(\mathbf{i}) = F \times (\mathbf{i}|\mathbf{g}|1)^t.$$

Subscript function returns a vector whose dimensionality is equal to the dimensionality of an array  $M$ .

*Scheduling function.* Iteration domains define exactly the set of executed dynamic instances for each statement. However, this algebraic structure does not describe the order in which each statement instance has to be executed with respect to other statement instances (46). A convenient way to express the execution order for each statement instance is to give each instance an execution date. It is obviously impractical to define all of them one by one since

the number of instances may be either very large or unknown at compile time. An appropriate solution is to define, for each statement  $S$ , a *scheduling function*  $\theta_S$  mapping instances of  $S$  to multidimensional timestamps (vectors). For tractability reasons, we restrict these functions to be affine, and we will use matrix operations on homogeneous coordinates (additional dimension equal to the constant 1) to represent affine functions. For the sake of transformation composition or search space exploration (22),  $\theta_S$  is often broken into dedicated blocks: a matrix  $A^S$  operating on iteration vectors, a vector  $\beta^S$  for static (multidimensional) statement ordering, and a matrix  $\Gamma^S$  to parameterize the schedule and to model pipelining:

$$\Theta_S = \left[ \begin{array}{ccc|ccc|c} 0 & \cdots & 0 & 0 & \cdots & 0 & \beta_0^S \\ A_{1,1}^S & \cdots & A_{1,d^S}^S & \Gamma_{1,1}^S & \cdots & \Gamma_{1,d_g}^S & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 & \beta_1^S \\ A_{2,1}^S & \cdots & A_{2,d^S}^S & \Gamma_{2,1}^S & \cdots & \Gamma_{2,d_g}^S & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{d^S,1}^S & \cdots & A_{d^S,d^S}^S & \Gamma_{d^S,1}^S & \cdots & \Gamma_{d^S,d_g}^S & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 & \beta_{d^S}^S \end{array} \right]$$

$$\theta_S(\mathbf{i}_S) = \Theta_S \times \begin{pmatrix} \mathbf{i}_S \\ \mathbf{g} \\ 1 \end{pmatrix}$$

As an example we will consider the pseudocode in Figure 17a. The loop kernel is composed of three statements:  $S_1$ ,  $S_2$  and  $S_3$ . The iteration domains for statements  $S_1$  and  $S_2$  are the following:

$$D^{S_1} = \left[ \begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & -1 & 0 & 0 \end{array} \right] \quad \begin{array}{l} 0 \leq i \\ i \leq M - 1 \end{array}$$

$$D^{S_2} = \left[ \begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 1 & -1 & 0 \end{array} \right] \quad \begin{array}{l} 0 \leq i \\ i \leq M - 1 \\ 0 \leq j \\ j \leq K - 1 \end{array}$$

(domain for  $S_3$  is the same as domain for  $S_2$ ). Domain for statement  $S_1$  has single iterator dimension (corresponding to iteration variable  $i$ ), and two parameter dimensions (corresponding to  $M$  and  $K$ ). Domain for statement  $S_2$  has two iterator dimensions (corresponding to iteration variables  $i$  and  $j$ ). There are no array data access function for statement  $S_1$  because it does not access any array. Data access functions for statement  $S_2$  are the following:

$$\mathcal{W}^{S_2} = \left\{ \right\}$$

$$\mathcal{R}^{S_2} = \left\{ \left( \|x\|, [ 1 \ 1 \ 0 \ 0 \ 0 \ 0 ] \right), \left( \|c\|, [ 0 \ 1 \ 0 \ 0 \ 0 \ 0 ] \right) \right\} \quad \begin{array}{l} x[i+j] \\ c[j] \end{array}$$

There are no write accesses in statement  $S_2$ , while there are two read data accesses: one from array  $x$  and other from array  $c$ . Original (corresponding to original input code) scheduling functions for statement  $S_1$  and  $S_2$  are given:

$$\begin{aligned}
A^{S_1} &= \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} & A^{S_2} &= \begin{bmatrix} 1 & 0 & \\ 0 & 1 & \\ & & 1 \end{bmatrix} \\
\beta^{S_1} &= \begin{bmatrix} 0 & 0 & \\ & 0 & 0 \\ & & 0 \end{bmatrix}^t & \beta^{S_2} &= \begin{bmatrix} 0 & 1 & 0 \\ & 0 & 0 \\ & & 0 \end{bmatrix}^t \\
\Gamma^{S_1} &= \begin{bmatrix} 0 & 0 & \\ & 0 & 0 \\ & & 0 \end{bmatrix} & \Gamma^{S_2} &= \begin{bmatrix} 0 & 0 & \\ & 0 & 0 \\ & & 0 \end{bmatrix} \\
\text{i.e. } \Theta^{S_1} &= \begin{bmatrix} 0 & 0 & | & 0 \\ 1 & 0 & | & 0 \\ 0 & 0 & | & 0 \end{bmatrix} & \text{i.e. } \Theta^{S_2} &= \begin{bmatrix} 0 & 0 & 0 & | & 0 \\ 1 & 0 & 0 & | & 0 \\ 0 & 0 & 1 & | & 1 \\ 0 & 1 & 0 & | & 0 \\ 0 & 0 & 0 & | & 0 \end{bmatrix}
\end{aligned}$$

note that  $\Gamma^{S_1}$  and  $\Gamma^{S_2}$  are all zeros, since the schedule does not depend on global parameters.  $A^{S_1}$  and  $A^{S_2}$  are identity matrices.

Program optimization in polyhedral model is usually done in three steps: (1) static analysis of input program resulting in algebraic representation of static control loop nests and construction of dependence graph, (2) transformations of polyhedral abstraction (based on linear algebra and integer linear programming machinery) without touching syntax representation of original program (3) code generation (going back into syntax representation). Note that step (3) is done only once: all the transformations (sequences) are operating on polyhedral (algebraic) representation.

The following table lists the main loop transformations that we can freely compose and parameterize in our framework; see Allen and Kennedy’s book for reference on those transformations (1) and Girbal et al. for details on the encoding and composition invariants (22):

Transformation name	Matrices involved
Interchange, skewing (unimodular)	A
Strip-mining, tiling	D, $\Theta$
Pipelining (multidimensional)	$\beta$
Parametric pipelining (multidimensional)	$\Gamma$
Reversal (unimodular)	A, $\Gamma$
Motion, fusion, fission (distribution)	$\beta$
Privatization, contraction	F

Considerable advances in dependence analysis, optimization and parallelization heuristics, and code generation proved that polyhedral model is scalable enough to be used in industrial tools. Yet, the problem of devising the optimal transformation sequences for optimizing locality and enabling parallelism is still a topic of considerable research efforts.

There are two approaches to optimizing programs using polyhedral model: static analytical modelling and iterative optimization. The former builds an analytical model and tries to statically predict the best possible loop transformation. Iterative optimization takes a feedback-directed approach, building different versions of the input program by applying different optimizations and choosing the one that gives the best performance gains. Those two approaches are complementary: analytical modelling can miss the best optimization opportunity but it takes just a one pass to complete. On contrary, iterative optimization might search for the best optimization but the search space might be huge, taking many iterations to complete, a challenge for its adoption in production compilers.

Our approach combines an analytical model and an iterative, feedback-directed approach. We rely on the loop tiling framework of Bondughula et al. (6), and on the associated tool Pluto, to extract blocked/tiled loop nests that exhibit one or more parallel loop levels. This framework includes a heuristic to select the shapes of the tiles, to maximize coarse grain parallelism while minimizing communication. This heuristic happens to behave consistently well on our benchmarks. However, another heuristic is proposed to deal with the combinatorics of loop fusion/distribution, and this one is not robust enough to achieve good performance on a wide variety of

benchmarks and on multiple architectures. We thus replaced the fusion/distribution heuristic with an iterative search approach, building a search space of legal, unique transformations before looking for a proper tiling scheme with Bondughula’s heuristic. The iterative search is adapted from the more general technique of Pouchet et al. (46), with new invariants enabling to focus the search space to more relevant transformations (for the practicality of the iterative search). Our results on 2 UTDSP and 2 BLAS2 benchmark are reported in Figures 13, 14, 15, 16, considering two different quad-core general-purpose processors. Those results demonstrate (1) the potential of our loop-nest optimizer compared to the state-of-the-art (Intel’s compiler), and (2) the wide performance differences among the different transformation sequences explored by the iterative search method.

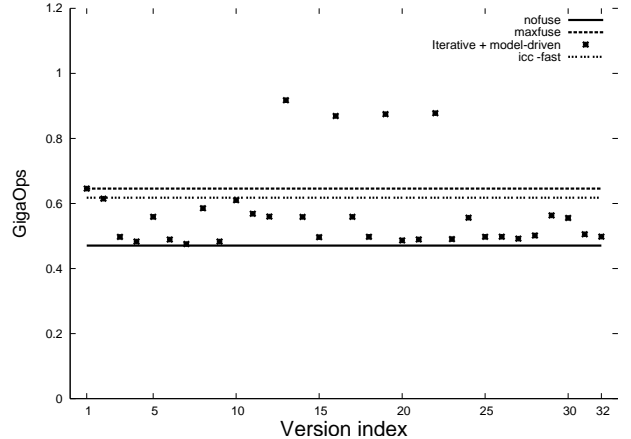


Figure 13. RegDetect on Intel Q6600

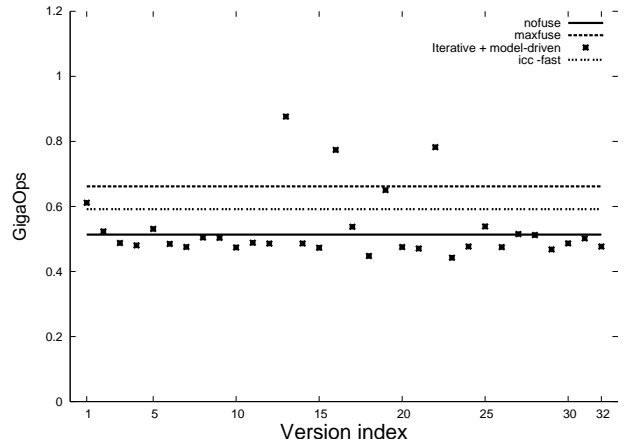


Figure 14. RegDetect on AMD Phenom 9850

## 5.2 Vectorization

Vectorization is the process of converting scalar source code into code that operates on vectors of data elements, making use of SIMD/vector instructions. Automatic vectorization by a compiler typically focuses on loops, where occurrences of an instruction across different loop iterations operate on different data elements. It can be seen as a downstream stage of the loop nest optimizer, where the selection of target-dependent instructions and access patterns come into play. Vectorization is known to be one of the most effective ways to exploit fine-grain data-level parallelism, and is

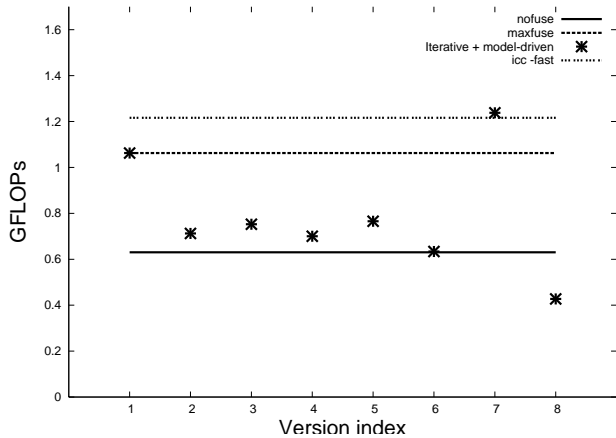


Figure 15. GEMVER:  $N = 8000$  on Intel Q6600

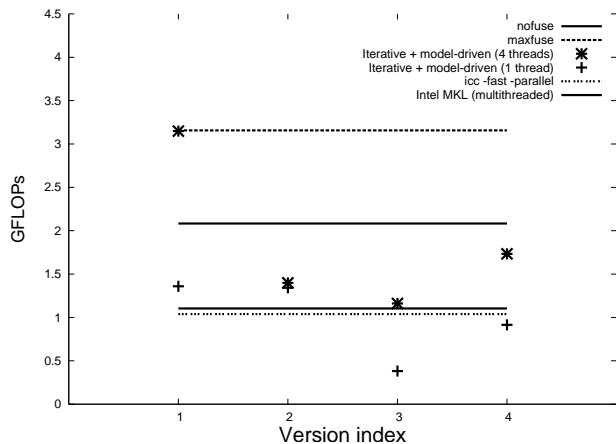


Figure 16. Doitgen on Intel Q6600

especially important for streaming architectures because their processing units typically contain vector (SIMD) units (see Section 7) to take advantage of the abundant data-level parallelism available in streaming applications (see Section 1.1). Vectorization has therefore been identified as one of the key compiler optimizations to be addressed by the ACOTES project.

The auto-vectorizer available since GCC 4.0 (36) is capable of vectorizing inner-most loops with memory references of unit or power-of-two strides, that may or may not be aligned, and that may include multiple data-types and type conversions, reductions, and other special idioms. The two main enhancements that were missing and were identified as important for the streaming domain are a cost-model, and the capability to vectorize outer-loops. One of the most important goals of the cost model, beyond facilitating informed decision making by the vectorizer, was to be part of an interface to exchange data and guidance with the high-level loop optimizer. This is described in detail in Section 5.3. The rest of this section focuses on the in-place outer-loop vectorization capability we developed.

Outer loop vectorization refers to vectorizing a level of a loop nest other than the inner-most, which can be beneficial if the outer loop exhibits greater data-level parallelism and locality than the inner-most loop. Figure 17c shows the result of vectorizing the outer  $i$  loop of the loop nest in Figure 17a, assuming Vector Length  $VL=4$  and  $M$  divisible by 4. Notice that the innermost  $j$  loop contin-

ues to advance in steps of 1 (compared to 4 in the case of innermost loop vectorization depicted in Figure 17b), computing 4 results for 4 successive  $i$  iterations simultaneously.

Outer-loops may have longer iteration counts, smaller strides, more efficient computation constructs, lower initialization and finalization overheads than those in innermost loops, as well as greater potential for promoting vector register reuse, thereby allowing us to leverage the many vector-registers often available in streaming architectures.

As mentioned above, data-streaming applications are characterized by nested loops of SIMD-like computations. The high-level data-reuse carried by the outer-loops in these loop nests can be detected and exploited only if operating at the level of the outer-loop. For this reason we have implemented an in-place vectorization approach that directly vectorizes the outer-loop (48; 53; 37; 35; 40; 42), instead of the traditional approach of interchanging an outer-loop with the inner-most loop, followed by vectorizing it at the inner-most position (1). The cost model we developed is capable of guiding the compiler which of these two alternatives is expected to be more profitable (as explained in the following Section).

Operating directly on the level of the outer-loop allows detecting high-level data reuse opportunities that are carried by the outer-loop, as well as fine grained data reuse opportunities related to the handling of alignment. We developed an optimization tapping such opportunities, incorporated within the outer-loop vectorizer. This optimization detects situations in which the misalignment of a vector load in a nested inner-loop is not invariant in the inner-loop (which happens e.g. if the inner-loop stride  $S$  is smaller than the Vector Length ( $VL$ )), yet the different misalignments across consecutive inner-loop iterations repeat themselves to form a cyclic group of  $VL/S$  distinct misalignments (if  $S$  evenly divides  $VL$ ). This is the case in the example in Figure 17c where  $S=1$  and  $VL=4$ . In this case we can achieve fixed misalignment by unrolling the inner-loop by  $VL/S$ . Fixed misalignment across iterations can be vectorized much more efficiently, as the misalignment (and in turn, the permutation masks to extract the desired data) can be computed once before the loop (instead of in each iteration). Moreover, each such cyclic group of loads exhibits a high rate of overlap in the data that is being fetched, and can be optimized by removing redundant loads from the unrolled iteration.

Note, however, that such unrolling may result in high register pressure, which on some architectures may result in register spilling, incurring high overhead that masks away the above benefits. For this reason, depending on the required unrolling factor, this optimization may not be suitable for architectures with too few vector registers, but is especially appropriate for streaming architectures, that often include a relatively large number of vector registers (e.g. 128 on the Cell SPE).

We evaluated these techniques on a set of multimedia benchmarks. Our implementation of in-place outer-loop vectorization, achieves speedup factors of  $2.92x$  on average across this set of benchmarks, compared to  $1.21x$  achieved by innermost loop vectorization. Outer-loop vectorization provides superior speedups for most benchmarks due to smaller strides and/or larger outer-loop-counts than those in the inner-loop, and/or by avoiding a reduction-epilog penalty. The optimization for fixed misalignment and data reuse using unrolling is capable of further boosting the performance obtained by outer-loop vectorization, to achieve an average speedup factor of  $4.98x$  (for detailed results see (42)).

These techniques are also applicable to the main computation kernels in the streaming applications described in Section 1, namely H.264 and Gamma-correction (in FMradio the inner-most loops are the best choice for vectorization).

*Gamma correction.* The complete gamma-correction algorithm we used for this evaluation is shown in Figure 18. As the reader

<pre> for (i=0; i&lt;M; i++){ S1: s = 0 for (j=0; j&lt;K; j++){ S2: s += x[i+j] * c[j] } S3: y[i] = s } </pre> <p style="text-align: center;"><b>(a) Scalar</b></p>	<pre> for (i=0; i&lt;M; i++){ vs[0:3] = {0,0,0,0} for (vj=0; vj&lt;K; vj+=4){ vc = c[vj:vj+3] vs[0:3] += x[i+vj:i+vj+3] * vc } y[i] = sum(vs[0:3]) } </pre> <p style="text-align: center;"><b>(b) Inner-loop vectorized</b></p>	<pre> for (vi=0; vi&lt;M; vi+=4){ vs[0:3] = {0,0,0,0} for (j=0; j&lt;K; j++){ vc = {c[j],c[j],c[j],c[j]} vs[0:3] += x[vi+j:vi+3+j] * vc } vy[vi:vi+3] = vs[0:3] } </pre> <p style="text-align: center;"><b>(c) Outer-loop vectorized</b></p>
---	---	--

**Figure 17.** FIR-filter Vectorization

```

volatile int inBuf[LINESZ]; // global arrays
volatile int outBuf[LINESZ];
...
// Initialize correction coefficients:
int ta[SN]={...}, oa[SN]={...}, ga[SN]={...};
for( line=0; line<LINES_PER_FRAME; line++ ) {
signal_inputFormatter(); // Line buffer empty
wait_inputFormatter(); // Processor stalled,
// until line buffer is full.
for( c=0; c<LINESZ; c++ ) {
int x = inBuf[c];
int t=thrh[0], o=offs[0], g=grad[0];
// Search for correction interval:
for( i=0; i<SN-1; i++ ) {
int flag = x>ta[i];
// Found interval; set pixel coeffs:
o = flag ? o : oa[i+1];
g = flag ? g : ga[i+1];
t = flag ? t : ta[i+1]; }
// Calculate output:
outBuf[c] = o + asrrnd((x-t)*g,SCALE );}
// Line output buffer full; Signal&wait for DMA.
signal_DMA();
wait_DMA(); } // Stalled until DMA is ready.

```

**Figure 18.** gamma-correction algorithm

will have noticed, the algorithm continues searching after the interval has been found, but will no longer update the three coefficients. Thus, the algorithm is predictable and more amenable to vectorization. Vectorizing the outer-loop that scans over pixels in a row (the  $c$  loop), is most effectively done when the inner-most loop (that searches through the array of thresholds,  $i$  loop in the Figure) is completely unrolled (this means that the pixel-loop becomes the inner-most loop in the nest, and so regular inner-loop vectorization can be applied), achieving an improvement factor of 1.81x/10.15x over the sequential version on PowerPC970/Cell-SPU respectively. The super-linear speedup on the Cell-SPU (the Vectorization Factor is 4) is due to the heavy penalty for branches on the SPU (which the vectorizer converts into conditional vector operations) and the data-rotate overhead for operating on scalar data in vector registers (which is avoided when the code is vectorized). Outer-loop vectorization can be used if the inner-most loop is not completely unrolled, e.g. if its loop-count is unknown at compile time or too large to completely unroll. In the particular case at hand however the inner-most loop-count ( $SN - 1$ ) is a compile-time known constant 3.

As a point of comparison, we note that these optimizations can be expressed using Silicon Hive’s SPM (as described in Section 4.6). A manually optimized implementation of gamma correction for the ISP2300 is shown in Figure 12. It applies these optimizations (vectorization, loop unrolling) and others (software pipelining with double buffering and explicit synchronization) achieves a 200-fold acceleration, bringing performance up to 0.76

cycles/pixel. Vectorization is expressed using overloaded comparison and conditional operators and by making the pixel loop iterate over vectors of pixels rather than a single pixel, achieving a  $30\times$  improvement factor - close to the theoretical speedup factor  $VF = 32$ . While the data-level parallelism exploited by the compiler using GCC auto-vectorization is comparable to the one achieved manually using Silicon Hive’s SPM, the other optimizations that the latter applies are not yet supported by the ACOTES toolchain. In gamma-correction these further increase data-level parallelism by a factor of 3 on the ISP2300.

*H.264*. The main computation kernel in H.264 consists of two modes: vertical and horizontal. The first mode, has a consecutive access-pattern in the inner-most loop, and an outer-loop that scans through different rows. Here only inner-loop vectorization is applicable, and it obtains a  $5.2\times$  to  $7.6\times$  performance improvement on PowerPC970/Cell-SPU respectively. The second mode has a consecutive access-pattern in the outer-loop and a row-size stride in the inner-loop. In-place outer-loop vectorization is applied here and achieves a  $10\times$  to  $11\times$  performance improvement on PowerPC970/Cell-SPU respectively. The alternative of first interchanging the two loops in the nest and then applying inner-loop vectorization (rather than vectorizing the outer-loop in-place) achieves a speedup of  $6.8\times$  to  $11\times$  on PowerPC970/Cell-SPU. In-place outer-loop vectorization is better on PowerPC970 than the interchange based approach due to improved locality (which the Cell-SPU is less sensitive to as it does not have a cache). Reasoning about such tradeoffs and selecting between these alternatives is the role of the cost-model, which is discussed in more detail in the next section.

### 5.3 Interaction Between Loop-Nest Optimizations and Vectorization

Vectorization involves low-level, target-specific considerations and transformations, which currently exclude it from being part of the polyhedral framework. In this section, we make a first step in this direction, building a performance model for automatic vectorization that integrates seamlessly within the polyhedral representation. Figure 19 summarizes this integration step in the context of the GCC compilation flow. We address a key adaptation problem when porting a streaming application to a new target architecture; it facilitates educated decision making on how to best apply loop transformations while considering the subsequent effects of vectorization.

To profitably apply vectorization to a loop-nest, an intelligent decision needs to be made as there are often several alternatives to choose from, each with its associated positive or negative performance gain. One such *default* alternative is to keep the original, unvectorized code. Several key costs impact the expected performance of vectorized code, including: strides of accesses to memory, loop trip counts, reduction operations across loops and more. These factors depend on the loop form being vectorized, however, and must be adapted if certain loop transformations are applied to

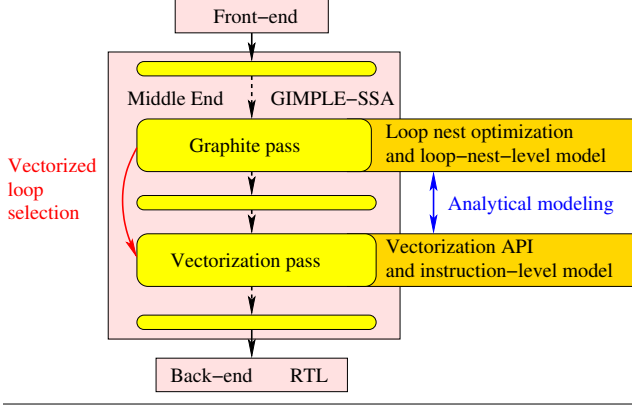


Figure 19. GCC compilation flow

the loop nest. Here we describe how to integrate these factors into the polyhedral model, thereby adapting them seamlessly as different loop transformations are being considered (without explicitly transforming the code), facilitating efficient application of both loop transformations and vectorization.

The underlying assumption of vectorization is that the kernel of a loop usually executes faster if vectorized than if not, but that associated overheads may hinder the vectorized version, diminishing its speedup compared to the original scalar version, and more so for loops that iterate a small number of times. Indeed, if the number  $N$  of iterations of a loop is smaller than its *vectorization factor*  $VF$ , there is no potential for speeding it up using vectorization; on the contrary, vectorizing such a loop may only slow down its execution due to additional preparatory actions and checks. Furthermore, even if  $N$  is larger than  $VF$ , the number of iterations of the vectorized loop, although positive, may not suffice to out-weigh the overheads incurred by vectorization.

Alongside each transformation, the cost for individual statement is updated. We can iteratively perform different transformations in polyhedral model by changing schedule  $\theta^S$  for each statement  $S$  and computing the cost function of the new schedule. At the end we pick the best possible schedule, based on the minimal cost.

In a classical polyhedral framework, access functions for array references are represented as affine expressions. One may compose this subscript function with static knowledge about the data layout of an array. For each array reference, one may form a *linearized memory access function*  $\ell$ , capturing the stream of memory access locations as function of the iteration vector:

$$\ell(\mathbf{i}) = b + (\mathbf{L}^i | \mathbf{L}^g | \omega) \times (\mathbf{i} | \mathbf{g} | 1)^t = b + \mathbf{L}^i \mathbf{i} + \mathbf{L}^g \mathbf{g} + \omega$$

where  $b$  is the base address of the array and  $(\mathbf{L}^i | \mathbf{L}^g | \omega)$  is the vector of coefficients that encodes the layout information for data array (assuming row-major data layout).  $b$  is typically not known at compilation time; nevertheless, the most relevant property of  $b$  (for our purpose) is its alignment modulo the vector size, which is generally available from language and `malloc` alignment guarantees. It is crucial for computing the cost of vectorized data load/store instructions, which constitutes the majority of the vectorizer overhead.

After applying the schedule transformation, the new time-stamp vector is expressed as follows (we ignore  $\beta$  vector component of the schedule for the purpose of data access modelling):

$$\mathbf{s} = (\mathbf{A} | \Gamma) \times (\mathbf{i} | \mathbf{g})^t = \mathbf{A} \mathbf{i} + \Gamma \mathbf{g}$$

thus, original iteration vector  $\mathbf{i}$  is given by:

$$\mathbf{i} = \mathbf{A}^{-1}(\mathbf{s} - \Gamma \mathbf{g})$$

which gives us the new, transformed linearized access function:

$$\ell'(\mathbf{s}) = b + \mathbf{L}^i \mathbf{A}^{-1} \mathbf{s} + (\mathbf{L}^g - \mathbf{L}^i \mathbf{A}^{-1} \Gamma) \mathbf{g} + \omega$$

with a new vector of coefficients:

$$\mathbf{L}' = (\mathbf{L}^i \mathbf{A}^{-1} | \mathbf{L}^g - \mathbf{L}^i \mathbf{A}^{-1} \Gamma | \omega)$$

This shows that linearized access functions, on which the total vectorization cost depends, is transformed automatically with the scheduling transformations. Thus, we do not need to generate code in order to compute the vectorization cost after applying a set of loop transformations – the vectorization cost is the direct function of scheduling matrix. For the rest of presentation we focus on  $\mathbf{L}^i$  part of the linearized access function coefficient vector.

Cost model is based on modelling the vectorization cost per statement, given its modified iteration domain  $\mathcal{D}^S$  and  $\theta_S$ . The cost function for statement  $S$  is the following:

$$\begin{aligned} c_{vec}(S, d_v) &= \frac{|\mathcal{D}^S|}{VF} \left( \sum c_{vect.instr} \right) + \\ &\quad \sum_{m \in (\mathcal{W}_S)} \left( f_a + \frac{|\mathcal{D}^S|}{VF} (c_{vect.store} + f_m) \right) + \\ &\quad \sum_{m \in (\mathcal{R}_S)} \left( f_a + \frac{|\mathcal{D}^S|}{VF} (c_{vect.load} + f_s + f_m) \right) \end{aligned}$$

Where  $|\mathcal{D}^S|$  is the total number of iterations for statement  $S$  (taking into the account all the nested loops around the statement  $S$ ).  $VF$  is the vectorization factor, whereas  $d_v$  represents the depth of the loop meant for vectorization. It is implicitly assumed, that function depends on modified iteration domain  $\mathcal{D}^S$  and  $\theta_S$  (through the change of linearized access function).

Given the desired loop level  $d$  we can extract the stride  $\delta$  of memory access with respect to that loop by simply looking at element  $d$  of  $\mathbf{L}^i$  vector:

$$\delta_d = \mathbf{L}_d^i$$

stride at the vectorized loop level is obtained by taking the relevant element from linearized access function and multiplying it by  $VF$ :

$$\delta_{d_v} = \mathbf{L}_{d_v}^i \cdot VF$$

Factor  $f_s$  considers the penalty of load<sup>1</sup> instructions accessing memory addresses with a *stride* across the loop being vectorized. Typically, unit-strided (i.e. consecutive) accesses to memory are supported most efficiently by vector loads and stores, incurring minimal if any overhead. However, accesses to non-unit strided addresses may require additional data unpack or pack operations, following or preceding vector load or store instructions, respectively (41). For example,  $VF$  scalar accesses to memory addresses with stride  $\delta_{d_v}$  across the loop being vectorized may require  $\delta_{d_v}$  vector loads (each with cost  $c_1$ ), followed by  $\delta_{d_v} - 1$  vector extract odd or extract even instructions (each with cost  $c_2$ ), to produce one vector holding the desired  $VF$  elements. On the other hand, if several accesses to the same address are vectorized together (i.e.  $\delta_{d_v} = 0$ ), a vector “splat” instruction is often required to propagate the loaded value across all elements of a vector (with cost  $c_0$ ). Equation 1 shows how factor  $f_s$  is computed as a function of the stride  $\delta_{d_v}$ .

$$f_s = \left\{ \begin{array}{ll} \delta_{d_v} = 0 : & c_0 \\ \delta_{d_v} = 1 : & 0 \\ \delta_{d_v} > 1 : & \delta_{d_v} \cdot c_1 + (\delta_{d_v} - 1) \cdot c_2 \end{array} \right\} \quad (1)$$

<sup>1</sup> Storing vectors with strided access is not yet implemented in current GCC vectorizer.

Factor  $f_a$  considers the *alignment* of loads and stores. Typically, accesses to memory addresses that are aligned on  $VF$ -element-boundaries are supported very efficiently, whereas other accesses may require loading two aligned vectors from which the desired unaligned  $VF$  elements are extracted (for loading) or inserted (for storing).

This alignment overhead may be reduced considerably if the stride  $\delta$  of memory addresses accessed across loop levels  $d_v+1..d^S$  is a multiple of  $VF$ , since the misalignment remains constant inside the vectorized loop. If this is the case there is the opportunity to reuse loaded vectors and use invariant extraction masks. By having the transformed linearized access function:

$$\ell(\mathbf{i}) = b + L_1^i i_1 + \dots + L_{d_v}^i i_{d_v} + \dots + L_{d^S}^i i_{d^S} + L^g \mathbf{g} + \omega$$

it is easy to check if misalignment inside the vectorized loop remains constant – the coefficients from  $L_{d_v+1}^i$  to  $L_{d^S}^i$  (corresponding to strides of all inner loops of the vectorized loop) have to be a multiple of  $VF$ .

If the misalignment is constant inside the vectorized loop we also check if the base address which is accessed on each first iteration of the vectorized loop ( $d_v$ ) is known to be aligned on  $VF$ -element-boundary; if so then there is no need for re-aligning any data:  $f_a = 0$ . This is done by considering strides across outer-loops (enclosing the vectorized loop, if exist), and initial alignment properties such as array alignment. In order to check the alignment in outer loops, we need to check if coefficients from  $L_1^i$  to  $L_{d_v-1}^i$  are multiple of  $VF$ .

By putting together all considerations for alignment, cost can be modelled as:

$$f_a = \left\{ \begin{array}{ll} \text{aligned} & : 0 \\ \text{var. misalign.} & : |\mathcal{D}_{1..d_v-1}^S|(c_1 + c_3 + c_4) \\ \text{fixed misalign.} & : |\mathcal{D}_{1..d_v-1}^S|(c_1 + c_3) + |\mathcal{D}_{1..d_v-1}^S|(c_1 + c_4) \end{array} \right\} \quad (2)$$

where  $c_3$  represents the cost of building a mask based on the misalignment amount,  $c_4$  is representing the cost of extraction or insertion and  $c_1$  is the vector load cost.  $|\mathcal{D}_{1..d_v-1}^S|$  denotes the number of iterations around the vectorizer loop level.

The vectorization factor  $VF$  of a loop is determined according to the size of the underlying vector registers and the smallest datatype size of variables appearing inside the loop. Each individual vector register will thus be able to hold  $VF$  variables of this small size. However, if there are variables in the loop of larger size, storing  $VF$  copies of them will require multiple vector registers, which in turn implies that the associated instructions need to be replicated. Factor  $f_m$  records the extra overhead that is associated with this replication.

Additional factors that depend on the specific machine resources available may also impact the performance of vectorization, such as the size of register files, available ILP, and complex vector instructions.

By taking an example kernel in Figure 17 the linearized access function for arrays  $\mathbf{x}$  and  $\mathbf{c}$  are as follows:

$$\ell_x(\mathbf{i}) = b + ( 1 \ 1 \ 0 \ 0 \mid \omega ) \times (\mathbf{i}|\mathbf{g}|1)^t$$

$$\ell_c(\mathbf{i}) = b + ( 0 \ 1 \ 0 \ 0 \mid \omega ) \times (\mathbf{i}|\mathbf{g}|1)^t$$

After applying loop interchange transformation by having the interchange matrix as follows:

$$A^{S_2} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

the new access functions would become as follows:

$$\ell'_x(\mathbf{i}) = b + ( 1 \ 1 \ 0 \ 0 \mid \omega ) \times (\mathbf{i}|\mathbf{g}|1)^t$$

$$\ell'_c(\mathbf{i}) = b + ( 1 \ 0 \ 0 \ 0 \mid \omega ) \times (\mathbf{i}|\mathbf{g}|1)^t$$

Notice that strides has changed. If we choose to vectorize the innermost loop, before the transformation, the access stride in matrix  $\mathbf{c}$  with respect to the vectorized loop was 1, while after the loop interchange, the stride with respect to the vectorized loop is 0. The cost function is updated accordingly

By applying different loop interchange transformations and choosing different loops to vectorize, the performance of the resulting vectorization varies considerably. Our model was able to predict the best possible combination of loop interchange and outer/inner vectorization strategy.

We evaluated our approach by introducing our model into the polyhedral framework of GCC<sup>2</sup>, and comparing its performance estimates for different loop interchanges and vectorization alternatives against actual execution runs of a set of benchmarks. The set of benchmarks includes a rate 2 interpolation (*interp*), block finite impulse response filter (*bkfir*), an  $8 \times 8$  discrete cosine transform for image compression (*dct* (33)), 2D-convolution by  $3 \times 3$  filters for edge detection (*conv*), a kernel from H.264 (*H264*), video image dissolve (*dissolve*), weight-update for neural-nets training (*alvinn*) and a  $16 \times 16$  matrix-matrix multiply (*MMM*) (including a transposed version *MMM\_trans*).

On the SPU, in all but one case (*alvinn*) the model correctly predicted the best vectorization technique. Using the cost-model driven approach, we obtain an average speedup of 3.5x over the scalar version, which is an improvement of 36% over the optimized in-place outer-loop vectorization technique, and 2.3x times better than the innermost vectorization approach, on average.

On the PPC970, the cost model mispredicts in 3 cases (*interp*, *bkfir* and *alvinn*). The Overall speedup obtained by the cost-model driven approach is 2.9x over the scalar version, an improvement of 50% over outer-opt, and 2.3x times better than innermost loop vectorization, on average.

Incorporating the vectorizer into the polyhedral framework is still work in progress, and not fully automated. Until auto vectorization is fully integrated into the polyhedral interface, the vectorizer provides information to the high level loop optimizer via an interface which exports its analysis utilities, and allows the high level loop optimizer to query which loops could/would be vectorized, thereby improving the effectiveness of the high level loop optimization heuristics, and has the potential of speeding up the exploration of loop transformations. The following utilities were identified as useful to assist the high level optimizer to pick the right loops to optimize for vectorization:

1. given a loop, return whether the loop is vectorizable or not;
2. given a loop, return iteration count threshold for profitable vectorization;
3. given a loop and an associated number of iterations, return performance impact from vectorizing the loop.

This API is meant to assist the high level loop optimizer when strip-mining/coarsening loops for the auto parallelizer vs. the auto vectorizer. It is implemented in GCC, and is available to be used by passes outside vectorization.

## 6. Compiler Back-End

This section describes the last stage of the ACOTES tool-chain, involving compiler back-ends and code generation for the relevant streaming platforms considered, including the runtime support of these platforms. Previous sections described uniform, largely

<sup>2</sup> Graphite, <http://gcc.gnu.org/wiki/Graphite>

target-independent stages of the ACOTES tool-chain, from programming model and programmer annotations, through high-level transformations and down to low-level transformations using GCC. From this point several alternatives were considered, to support code generation for the different streaming targets provided by the partners.

### 6.1 Common Compiler Backend Format

In order to increase the amount of sharing among the partners, we studied the suitability of adopting common compiler back-end formats. Several candidates existed for this purpose. A natural one was the GCC RTL format, since GCC4 is already used as the basis for the higher level parts of the project.

Indeed for one platform, namely the Cell Broadband Engine, the native GCC RTL format was employed. The GNU tool-chain served the development of the Cell processor from its early architecture exploration and programming environment creation (17; 52), and so GCC already provided the ability to generate code for both the PPE (PowerPC Engine) and the SPE's (Synergistic Processing Engines) of Cell. Additional enhancements were developed to improve auto-vectorization and other transformations for Cell, that are needed in order to program and compile streaming applications effectively (47).

Alternatively, we considered the suitability of higher level formats, such as the Java bytecode or the Low Level Virtual Machine (LLVM) (32). Finally, the ECMA-335 Common Language Infrastructure (CLI) standard (13), at the basis of the well known Microsoft .NET technology, was seriously considered. The motivations were multiple: first of all it is an open standard, providing a clear separation between multiple compiler front-ends producing CLI, and multiple CLI consumers, providing therefore a natural support for multiple ISA mapping; secondly, it supports the semantics of multiple languages, including C, which is used for our project. And finally, thanks to Microsoft, it is a widely adopted technology, with growing interest also in the open source community, as indicated by the Mono (39), DotGNU Portable.NET (49) or ILDJIT (7) projects.

The outcome of the investigation (10) is that GCC RTL is indeed confirmed as a natural candidate as compiler back-end format. For what concerns the high-level processor-independent formats, the characteristics of CLI bytecode and data representation, make it an excellent choice for ACOTES: the framework has been standardized by two different entities, ECMA and ISO, and is stable. It can be used as a processor independent format or can be specialized for a particular target. Additional information can be added to the code, variables and types to drive further optimizations. Finally it supports both managed and unmanaged environments.

In order to exploit the GCC4-based developments of other partners and to share visibility on this format within the project, STMicroelectronics developed a GCC4 to CLI translator (11; 51). We leveraged from existing efforts in this direction by the open source community that we consolidated and complemented by adding the streaming information. This translator is free software.

### 6.2 Split Compilation

Split (or multi staged) compilation refers to the fact that the compilation process may be decomposed into several steps. In our context, the first step occurs on the workstation of the developer. It takes C code as input and generates a program representation in CLI format. The second step occurs on the final device where the application is run. It reads the CLI format and generates the native binary format. This latter step can be designed to run at *install-time* or at *run-time*, depending on the system.

We proposed to take advantage of this two-step situation to transfer the complexity of compiler optimizations as much as pos-

sible to the first stage. When an optimization cannot be applied — either because it is target-dependent, or because it may increase code size too much, or because it is too costly to be applied at runtime — it might still be considered: an analysis can be run and its results encoded into annotations embedded in the intermediate format. The second stage can rely on the annotations and skip expensive analysis to implement straightforward code transformations. Annotations may also express the hardware requirements or characteristics of a piece of code (I/O required, benefits from hardware floating point support, etc.)

This multi stage process makes it possible to apply to embedded systems the most recent aggressive techniques like iterative compilation (5) or transformation in the polyhedral model (22).

## 7. Target Streaming Platforms

The ACOTES project targets four streaming platforms. Of these four, two (the Cell B.E. and the HiveFlex ISP2300) are in production, while the other two (the exStream and the Ne-XVP) are in the design stage.

### 7.1 Cell B.E.

The Cell/B.E. (16) is a heterogeneous multicore processor that consists of an IBM 64-bit Power Architecture core, called the IBM PowerPC processor element (PPE), augmented by eight specialized single-instruction, multiple-data (SIMD) coprocessors (See Figure 6(a)). These coprocessors, called synergistic processor elements (SPEs) (15), provide data-intensive processing. The SPE accelerators operate from a local storage that contains instructions and data for a single SPE; this is the only memory directly accessible from the SPE. Memory access is performed via a DMA-based interface using copy-in/copy-out semantics.

The PPE is fully compliant with the 64-bit PowerPC Architecture and can run 32-bit and 64-bit operating systems and applications. The SPEs are independent processors, each running its own individual application programs that are loaded by the application that runs on the PPE. The SPEs depend on the PPE to run the operating system, and, in many cases, the top-level control thread of an application. The PPE depends on the SPEs to provide the bulk of the application performance. The SPEs are designed to be programmed in high-level languages and support a rich instruction set that includes extensive single-instruction, multiple-data (SIMD) functionality. However, just like conventional processors with SIMD extensions, use of SIMD data types is preferred, not mandatory.

The PPE and the SPEs are compiled separately by different back-end compilers and then linked together to compose an application. CELL/B.E. introduces multilevel parallelism, which users must exploit to gain the best performance. The lower level is the SPE level in which vectorization of inner and then outer loops should be explored. The higher level is functional parallelism level that should enable to distribute the processing between several SPEs.

A first step towards enabling auto-vectorization for the Cell is to model the SPE registers and instructions as vector registers and instructions in the compiler. In GCC this is done in special machine-description files in the GCC port for the Cell SPE. Once this model is in place the existing auto-vectorization capabilities of the compiler are transparently enabled.

### 7.2 Ne-XVP architecture from NXP semiconductors

The Ne-XVP (Nexperia eXtreme Video - or Versatile- Processor) architecture from NXP Semiconductors Research allows to get high silicon efficiency and low power for streaming applications like video processing (video encoding and decoding, frame rate con-



version, image improvement), while keeping a high level of adaptation to algorithmic change and application diversity. It is therefore a scalable architecture based on duplication of programmable or not - cores sharing a common memory structure. The right selection of the characteristics of the cores allows getting a very high efficiency tuned for each application domain while keeping a low cost of ownership by reducing the verification and validation costs, and reusing common module elements.

*multi-core* approach is well represented now in the industry: the clock race has reached its limits and the era of multi-core is coming. But contrary to the typical approach for multicores (heterogeneous), the Ne-XVP architecture uses in a better way the available silicon, starting from the observation that existing cores are optimized for single core environment, and not for multicore. An optimum is better reached when the optimization is done globally than by adding separately optimized elements (at least for non linear systems, which is typically the case of the Design Space of architectures). But, of course, the elementary cores still should be very efficient for the application domain. It is why the architecture uses a VLIW architecture, well suited for embedded applications, and mainly the architecture which was optimized for years for video processing in Philips then NXP: the TriMedia Instruction Set Architecture (ISA).

At the core level, Ne-XVP is not making a new architecture from the TriMedia ISA, it only plays with the standard architecture parameters such as the number of issue slots, the number of registers, the composition and number of Functional Units and the size of data and instruction cache. The main change is the support for multi-threading in cores, by adding extra registers banks and minimum logics, allowing a fast context switch in a core. It enables a single physical core to look like several cores from a software or application view. The reason of adding multi-threading to cores is two fold:

- It eased the software and the port of applications: a multi-threaded application can work either in different cores, or on a single core, so the number of threads is virtualized. A multi-threaded application can work on different instances of the Ne-XVP architecture that physically differs by their number of cores.
- Multithreading makes each core more insensitive to latency: if a thread is waiting due to a cache miss or latency due to a functional unit, then another task might kick in, increasing the use of the core. Even if multithreading adds complexity to the cores, it allows also to decrease the requirements of some part of the cores (latency, bypass, etc), globally adding a very small area.

To reach the high level of Mops/W, the Ne-XVP will then use the maximum parallelism, but not all parts of applications can be parallelized. This is the well-known *Amdahl's law*: it is the irremovable sequential part of the application that will limit the ultimate performance of the architecture on a particular application. Therefore, the architecture has a part that runs fast to execute the sequential part of the application.

This fact leads to a heterogeneous approach, where a core is optimized for sequential tasks while others can work as *pool* for parallel operations. It should be noted that increasing the different type of cores beyond two did not really bring further improvement. More details can be found in (29).

The specialization and efficiency of the types of core can also be increased by adding specific instructions (*custom-ops* in TriMedia terminology). The instructions can of course be implemented in a specific functional unit in the cores (if they are simple), or be shared by different cores, leading to embedding specific coprocessors in the array of cores, depending on the complexity and granularity of

the function. This approach is taken in the Ne-XVP approach for the CABAC coprocessor and for coprocessors helping for the task managements and memory coherency. From a point of view of cost of development, rather complex functions are shared (-co-joining-) and therefore integrated as specific coprocessor, because it will not impact the design and verification of each cores, which is always tricky and expensive.

The resources are also tuned to the characteristics of the application domains. For example, the coherence is only required at certain particular points (communication between tasks, end of tasks, etc), which allows to implement a coherence mechanism that has a very light footprint.

To ease the scalability to various size of video picture, data partitioning is often used to parallelize applications. The hardware task scheduler accelerator takes benefit of this and dispatches the tasks to the available cores, allowing a same binary code to be executed without modification on different instances of the Ne-XVP architecture with various numbers of cores. Data partitioning also permit the *co-joining* of instruction caches, reducing the silicon footprint of the system.

All those considerations led to the Ne-XVP architecture template in Figure 20. The architecture is composed of several cores of two types – core1 and core2. The first type optimized for sequential code execution is instantiated only once to address the Amdahl's law bottleneck, whereas the second type is replicated many times to carry out the parallel workload. Each core, multithreaded, has a data and instruction cache, where the latter can be shared among several cores. Hardware Coherence Coprocessors maintain data cache coherence in hardware, which also governs the cache2cache tunnels. The Synchronization unit features various hardware accelerators (Hardware Synchronization Unit, Hardware Task Scheduler, Task Scheduling Unit) for fast inter-task synchronization and communication, which may be deployed independently of each other.

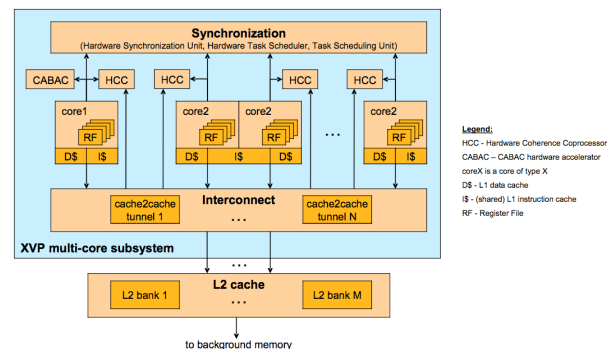


Figure 20. Ne-XVP architecture.

Compared to a *simple* multicore system based on standard “off-the-shelf” TriMedia TM3270 with 128KB data cache, 64KB instruction cache, 5 issue slots, 128 registers, without multithreading and coprocessors, the Ne-XVP architecture is 16 times more efficient in silicon area (and power) for applications such as H.264 decoding in CMOS 45nm technology. Special Architecture Features Targeting Selected Applications The Ne-XVP architecture has a large versatility concerning its programming model. The software can see the Ne-XVP hardware like:

- A common address space architecture. The Hardware Coherence Coprocessor ensures the coherence of the various versions of the data in different caches,
- A distributed memory architecture. Physically, the L1 caches are distributed, and cache-to-cache communications are done by the configurable tunnels.

- A data driven model: The Hardware Task Scheduler allows activating the various threads running on the cores only when data are ready and available.
- A multi-threaded architecture: the HTS and the cores allow implementing a model similar to the "Cell superscalar" (CellSs).
- A streaming architecture: the typical *Push, Pop, Peak, Sneak* functions can be efficiently emulated with the standard Ne-XVP instructions.

Several *streaming* applications were used in the definition of the Ne-XVP architecture; more particularly image processing functions, codec functions (H.264 at standard and Super-HD resolution), image enhancement and 3D Graphics.

For streaming applications mainly composed of a succession of kernels, Ne-XVP can be nearly as efficient as dedicated hardware: each kernel can be mapped on a core, and each core can be self-sufficient, running only with its internal local memory and registers. For the more complex applications, like H.264 super-HD, Ne-XVP allows to easily implement various thread level parallelisms. It could be semi-static, where the compiler does not have to guess the dependencies, which are checked at run-time by the Hardware Task Scheduler. The experiments show that performance is increasing quasi proportionally with the number of cores. The programming model is facilitated by the shared memory architecture, which makes it easier for programmers familiar to *C* or *C++* programming language.

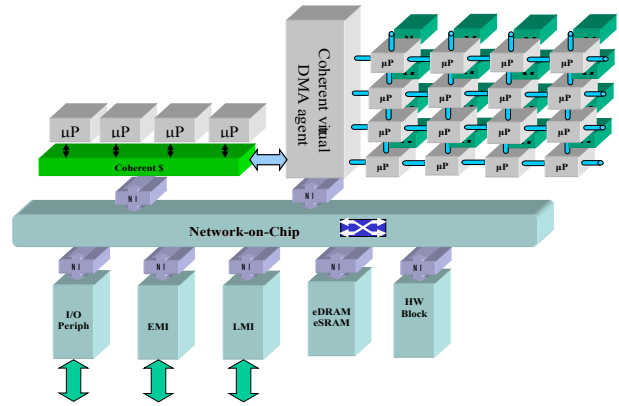
For the ACOTES compiler chain, its various elements are exercised to efficiently use the Ne-XVP architecture: the SMP allows dispatching tasks on different cores (or hardware threads in a core). The communication scheme of ACOTES allows to totally use the efficiency of the Hardware Coherence Coprocessor, requesting coherence only when it is required. Correct placement (linked to the ASM) allows minimizing the inter-core communication and keeps busy the tunnel units. Finally, the cores, based on the TriMedia VLIW Instruction Set Architecture, benefit from the loop and vector processing improvement of the ACOTES GCC compiler chain. The ACOLib library can also use the hardware mechanisms implemented in the various coprocessors to efficiently support the application.

### 7.3 xStream architecture from STMicroelectronics

The STMicroelectronics xStream architecture is based on the convergence of communication and computing as a way to solve scalability and programmability of high-performance embedded functionalities, such as graphics, multimedia and radio subsystems. In addition it addresses some of the increasingly challenging design and silicon fabrication issues at the architecture, micro-architecture and design levels through the use of advanced techniques such as voltage and frequency scaling (via a globally asynchronous locally synchronous model, GALS), local clock generators adapted to local silicon process variations (to increase yield and fault tolerance), skew insensitive design (via mesochronous and delay insensitive network on chip links) and use of regular cell design flows (for silicon manufacturability).

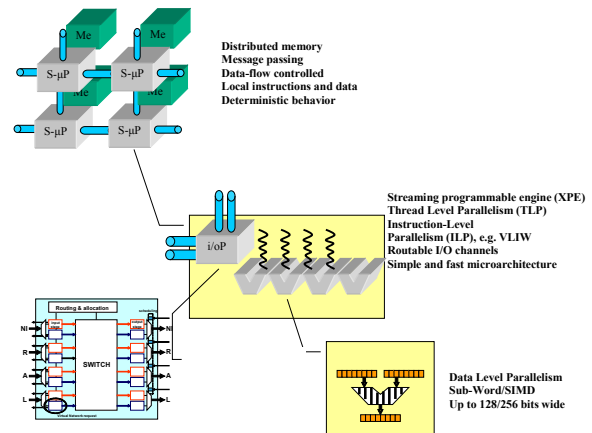
Figure 21 illustrates a high-level view of a complete system embedding an instance of the xStream processor fabric.

The system is composed of the traditional "host processing" part on the left side, which we have depicted as a Symmetric Multiprocessing (SMP) subsystem for future scalability, while the entity on the top-right end of the picture above is the 'streaming engine' of the xStream architecture. It is meant to be addressing the needs of data-flow dominated, highly computational intensive semi-regular tasks, typical of many embedded products. The streaming nature of the kernels mapped onto it makes it possible to design a semi-regular fabric of programmable engines interconnected via a rela-



**Figure 21.** High-level view of an embedded system containing an xStream computing fabric.

tively simple network of point to point channels. The structure of the tasks that is mapped onto the computational fabric is more similar to a pipeline of 'filters' rather than a set of tasks explicitly communicating between them, while the latter model is not ruled-out by the xStream template, including up to more traditional lock-based synchronization distributed parallel tasks. It mostly focuses on providing a range of facilities, both HW and SW, to fully exploit data-flow centric applications and programming models. The fabric supports a number of simultaneous software pipelines running on the processing elements as to accommodate complex applications and also provide load balancing and latency hiding capability. A property of the streaming fabric is to support very high internal data bandwidth, throughput and computationally intensive tasks.



**Figure 22.** The xStream fabric and node structure.

The processing elements of the streaming fabric are relatively simple programmable processors or engines with a general purpose but simple basic ISA that can be extended with SIMD or Vector mode instructions. The engines include a set of features for improving performance and efficiency, such as sub-word parallelism, wide data-paths, simple pipelines, etc. At the same time they execute instructions fetched from local memories instead of caches, a great simplification at the pipeline forefront. Local memory is also used for wide data accesses. The engines are connected between them, and the interconnect functionality plays one of the critical

roles in this picture. In fact it is quite the essence of the system to be able to provide a self synchronizing support for software pipelines. This is achieved with a set of lightweight routers, very similar to the ones being defined for network-on-chip replacements of standard bus infrastructures; but with more freedom for simplification, due to the constrained nature of the communication patterns versus a generic system back-bone NoC. The fabric is not limited to exploitation of programmable engines, in fact it is entirely possible to use hybrid approaches where some elements of the array can be fixed functions, implemented in a classic ASIC design flow, if such function are of a critical, very well know and fixed nature. Likewise the pipelines can be attached on the periphery of it to I/O channels that go to dedicated peripherals or specific I/O functions such as cryptographic engines or similar.

The xStream architecture provides native support for streaming applications, in particular for what concerns

- communication links,
- communication primitives, which are mapped directly into native instructions for maximum efficiency,
- memory operations for both internal and external memory access and
- a processing element well suited for data intensive computations (the xPE).

The interconnection network of the xStream is composed of a set routers connected by physical links (for example a 2D mesh topology). Each router has one (or more) end-point connected to it, which can be a consumer and/or producer of the packets flowing through the network. The end-points can be one of the following: An xStream accelerator processor or xPE, High-bandwidth IO links to the outside world or DMA channels.

Router-to-router connection as well as router-to-end-point connection is implemented over a single physical channel, whose parameters can be tuned for width, speed, arbitrations, local buffering, etc.

To improve performance and simplify low-level deadlock avoidance, the Network On chip supports virtual channels that are used to multiplex physical connections. The architecture also supports a communication feature at a higher level of abstraction: Virtual channels are implemented by using multiple "virtual" queues managed by a Flow Controller (xFC) at each end-point node of the fabric. A packet can be written into a specific "virtual" queue only if the credit based end-to-end flow-control allows it, this guards the application from running into so called high-level application deadlocks that might arise when the data-flow graphs do not satisfy certain conditions. Additionally, virtual queues can be of virtually any length, limited only by the available amount of local memory per each node, this feature greatly extend the freedom of software mapping tools to explore a larger legal solution space for the mapping problems and objective functions (such as local bottlenecks, load-balancing, etc.), and increases the potential of finding better optimized configurations. Virtual queues at the inputs and outputs of the xStream xPE processor are exposed at the instruction set level so that they can be effectively treated similarly to fully disambiguated sequential memory references by optimizing compilers.

The following ISA operations are defined for the queues of each xPE and can be described functionally as:

- pop destination\_reg, queue\_identifier
- push source\_reg, queue\_identifier
- peek destination\_reg, index, queue\_identifier
- qsync queue\_identifier

The queue\_identifier is a global queue identifier and specifies the source or destination nodes to which the operation is directed as well as the specific virtual queue on which the operation is to be performed.

The pop operation retrieves data from the head of the queue and fills the destination\_reg.

The push operation pushes the content of the source\_reg into the tail of the queue.

In practice queues are defined as variable length in the implementation through a mechanism that uses local memory to the xPE as a back-log temporary storage implemented by the xFC.

This same mechanism allows also to support a third primitive, which is the peek operation. Peek behaves like a pop issued after a successive number of pops for index-1 elements but does not actually remove the elements from the queue. Effectively it provides partial random access support for the queue itself. A peek operation is only allowed for a queue with back-log storage in local memory that has been declared to be large enough to contain a number of elements greater or equal to the maximum peek index used for a given input queue. In practice, peek blocks if less than index elements are present in the virtual queue. It may also be used for synchronization purposes and block based processing instead of pure streaming.

Finally the qsync primitive guarantees that an output virtual queue is drained of all of the data pushed through it and locally stored at the destination end, this instructions is required in a number of different scenarios, both for pipeline set up and shut downs, but also to allow advanced management of virtual queues and context switches.

The xStream architecture template can accommodate various kinds of computing nodes from programmable ones to hardwired functions. To complete the template with a suitable programmable element we have designed a highly parametric programmable engine that we call xStream processing element (xPE).

The xPE is optimized for stream oriented, data-flow dominated, performance demanding computation. The target application range is especially focused on embedded multimedia, telecom, and signal processing applications, Specific requirements that were used for the definition of the engine architecture and microarchitecture were:

Low silicon area or more exactly, high computing density in terms of MIPS per physical gate,

High computational power with outstanding power figures in term of MIPS/mW, MIPS/MHz, associated with relatively high (for the embedded world) operating frequencies, support for high level programming languages and compiler friendliness, especially for ISA orthogonality when it comes to vector/SIMD execution semantics. Finally the xPE is highly design time configurable with a wide range of tuning knobs and optional features that enable extensive trade offs for area, power, and compute density for a single xStream node to adapt the fabric granularity to the specific application domain an instance of it is design for.

The xPE microarchitecture is highly streamlined and "minimalist" core architecture to tune system frequency and limit core size by shaving off most of the complexity required for more general purpose microprocessors and media processors. The xPE execution semantics is VLIW coupled with a modular and scalable design based on configurable VLIW slices.

Local memories and interfacing are optimized for managing and accessing data streams, as well as wide vector instructions operating on packed data words to exploit available DLP. The xPE supports a fine-grained multithreading to exploit task level parallelism and, once again, to ease the data-flow application mapping tasks and to achieve the more conventional latency hiding benefits of multithreading.

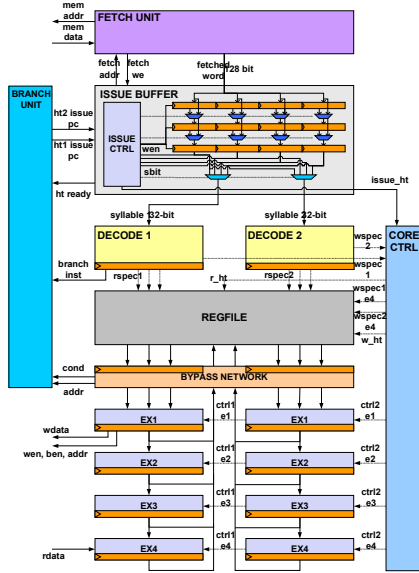


Figure 23. The xPE slice data-path

Each xPE slice includes two vector integer/floating point general-purpose ALUs, one vector integer/floating point general-purpose multiplier unit, one vector load/store unit with special stream access operations (push/pop to/from queues), independent fetch/issue unit with supporting features to synchronize multiple slices, shared registers for cluster communication and some simple branch prediction features. The xPE data-path supports either 32-bit operands or 128-bit packed vector operations ( $4 \times 32$ -bit or  $8 \times 16$ -bit), has an extensible register file with multiple views, as 64 general-purpose, 32-bit wide registers and one with 32 general-purpose, 128-bit wide vector registers. The xPE pipeline is 9 stages with up to 4 stages of execution depending from specific functional units' latencies. Full bypassing and forwarding is supported.

The vector extensions ISA include three source operand operations for SIMD most of which capable of a sub-word permutation operand for greater flexibility to support automated compiler vectorization.

The xPE ISA is designed for maximum flexibility for further extensions, but with code density very much in mind in terms of efficiency as most of the code will be fetched from relatively small local memory. A short summary of the ISA features is:

- Integer, fixed point, and floating point operations.
- Special instructions for thread synchronization.
- Instructions are encoded with 24-bit syllables.
- Up to two syllables are packed in a bundle and executed in parallel.
- Wide immediate values and various bundle extension formats.
- Variable bundle size to optimize code size.

#### 7.4 Silicon Hive HiveFlex ISP2300 Subsystem Architecture

Silicon Hive's basic system-level processing template is an IP-based structure containing an arbitrary number of processors, interconnected through buses, point-to-point connections, and/or streaming connections. Figure 24 depicts the template for a single processor. The processors within a system may all be different. Each may use various amounts of different types of parallelism: sub-operation-

level parallelism (pipelining within operations), operation-level parallelism (domain specific operations with many inputs and outputs), data-level parallelism (vector, SIMD), and instruction-level parallelism (ILP, VLIW). Using this template, Silicon Hive develops domain-specific multi-processor subsystems for streaming applications.

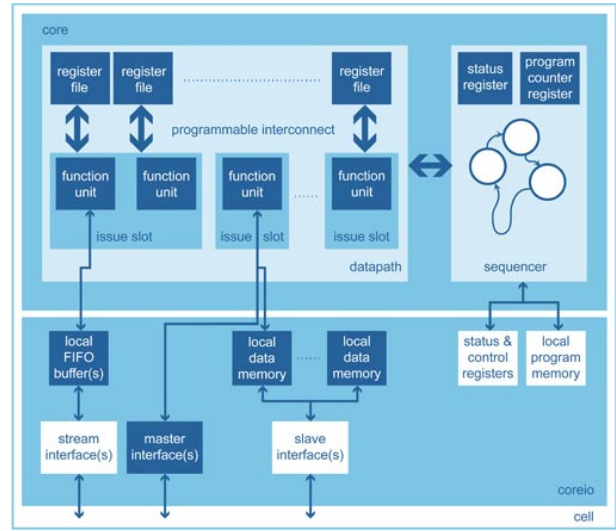


Figure 24. Silicon Hive Processor Template.

The ACOTES project targets streaming applications and has chosen three applications from three different streaming domains: H264 video processing, FM radio modulation (communication), and gamma correction (image processing). It was found that each of these domains exhibit different kinds of parallelism and thus require different processor architectures. A communications processor needs to rely more on ILP. On the image signal processing is more regular and thus can benefit most from extensive use of vectorisation. Lastly, video processing is a mix of control and image processing. Thus, a video processor would consist of multiple smaller vector processors, combined with a scalar control processor.

Next to applying the template features, as described above, the processors themselves need to be scalable, in terms of the above architectural parameters. Before committing a processor design to silicon, the architectural parameters are fixed. For example, the required performance points for a communications processor are obtained by scaling the number of issue slot clusters (ranging from 5 to 20 issue slots for complex arithmetic). The image signal processor (ISP2300) has a fixed set of 8 vector issue slots, but its vector and element sizes need to be fixed (typically, they scale from 4 to 128 and from 8 to 16, respectively).

This section discusses HiveGo CSS 31xx camera subsystem in more detail. It contains amongst others a HiveFlex ISP2300 processor, configured with 32-element vectors, each element being 16 bits wide. However, the concepts apply equally well to the typical communications or video processors.

Processor subsystems are based on a flexible template. And, within families, the different processor instantiations also have wide variations in architectural configurations. Thus, as an additional requirement to being able to deal with the different types of parallelism, mentioned above, the software development environment must be able to target a very wide range of different system and processor architectures.

Because of the above wide range of different architectures that need to be supported simultaneously, the requirement is that tools are not changed or generated to fit the target architecture. The number of different tools to be supplied would explode. Thus, each implied tool (e.g. compiler front-end, scheduler, linker, browser, simulator, etc.) needs to reads system description and needs to target its operation to the system.

Subsystems consist of processors, system-level IP blocks, and stream processing code.

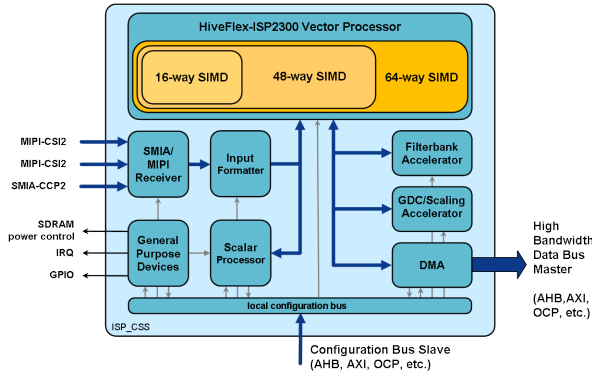


Figure 25. Block Diagram of HiveGo CSS 31xx Subsystem.

Figure 25 depicts the HiveGo CSS 31xx camera subsystem. In this system, HiveFlex ISP2300 and Scalar Processor are ANSI-C programmable processors. InputFormatter, FilterbankAccelerator, GDC/ScalingAccelerator, and DMA are weakly programmable autonomous modules. For the purposes of this paper, they may be considered parallel processors, streaming data into and out of memories inside ISP2300. Their operation needs to be synchronized with the application running on ISP2300.

InputFormatter reads raw data coming from the sensor through the SMIA/MIPI interface. The data is packed to fit the vector elements of ISP2300. ISP2300 itself performs a chain of operations on each pixel, whereby the raw sensor data is converted from a Bayer format into weighed RGB pixels. In addition to that, typical camera operations, such as white balance correction and gamma correction are performed. For typical convolutional filtering (motion blur, etc.) and scaling, the software architect may program the operation of the associated parallel processors to fit between certain software loops. See Figure 2, where co-processor-implemented functions, such as GDC/Upscaling phase, may be inserted between the color enhancement and sharpness enhancement phases, and where certain image enhancement phases would typically be executed on the Filterbank block.

## 8. Tool-Chain Evaluation

In this section we demonstrate how the concepts, tools and optimizations presented thus far can be used together in one vertical toolchain. We use the FMradio application (see section 1) as an example, and show how each step in the overall ACOTES framework can be applied to it.

The machine used for evaluation is a 4-core Power6 with 2-way SMT in each core, 64KB L1 cache (cache line size 128B) and 2MB L2 cache per core, running under Linux. Experiments were done using all 8 hardware threads (i.e. with affinity set to “all”).

### 8.1 Using the SPM

The parallelism in FMradio can be exposed using the Acotes directives and Front-End tools described in Sections 4.1 to 4.3 (version

```

taskgroup
{
  Reader (&pair)

  for (i = 0; i < 8; i++) {
    task input (pair) output (fm_qd_bp)
    FM_QD_Demod (pair, &fm_qd_bp)
    task input (fm_qd_bp) output (band_11)
    FFD (1.813, fm_qd_bp, &band_11)
    task input (fm_qd_bp) output (band_12)
    FFD (1.813, fm_qd_bp, &band_12)
    task input (fm_qd_bp) output (band_21)
    FFD (1.407, fm_qd_bp, &band_21)
    task input (fm_qd_bp) output (band_22)
    FFD (1.407, fm_qd_bp, &band_22)
    task input (band_11, band_12, band_21, band_22) output (ffd_bp)
    subMultSq (band_11, band_12, band_21, band_22, &ffd_bp)
  }
  task input (fm_qd_bp) output (band_2)
  FFD (8.407, fm_qd_bp, &band_2)
  task input (ffd_bp) output (band_3)
  FFD (8.407, ffd_bp, &band_3)
  task input (band_2, band_3)
  {
    stereo_sum (band_2, band_3, &output)
    Writer (output)
  }
}

```

Figure 26. Annotated FMradio

from now on denoted “SPM”) or using standard OpenMP pragmas, with the minor additional extension described in Section 4.4 (further denoted “GCC-SPM”). The GCC-SPM’s implementation, in the “streamOMP” branch of GCC, is still ongoing, so we will provide results based on manual streamization for this version. We demonstrate and evaluate both approaches.

After analysing FMradio, we observed that up to 13 tasks can be created and linked using streams. Of such tasks, 5 of them are very light, and we have joined them in two sets of 3 and 2 tasks. This leaves a total of 10 exploitable tasks. We have annotated the code with the Acotes directives, compiled it with the Acotes compiler, and linked using ACOLib. Figure 26 shows the structure of the resulting code. Observe how it describes the same structure as presented in Figure 1.

The alternative OpenMP-based implementation of the FMradio code presented on figure Figure 26 only requires minute modifications, like replacing the input and output clauses by their OpenMP counterparts `firstprivate` and `lastprivate` and adding the appropriate `parallel` and `single` directives. Figure 27 shows the resulting OpenMP annotated code. Here we only rely on GCC, with OpenMP enabled through the `-fopenmp` compiler option.

```

#pragma omp parallel
#pragma omp single
{
  Reader (&pair)

  for (i = 0; i < 8; i++) {
    #pragma omp task firstprivate (pair) lastprivate (fm_qd_bp)
    FM_QD_Demod (pair, &fm_qd_bp)
    #pragma omp task firstprivate (fm_qd_bp) lastprivate (band_11)
    FFD (1.813, fm_qd_bp, &band_11)
    #pragma omp task firstprivate (fm_qd_bp) lastprivate (band_12)
    FFD (1.813, fm_qd_bp, &band_12)
    #pragma omp task firstprivate (fm_qd_bp) lastprivate (band_21)
    FFD (1.407, fm_qd_bp, &band_21)
    #pragma omp task firstprivate (fm_qd_bp) lastprivate (band_22)
    FFD (1.407, fm_qd_bp, &band_22)
    #pragma omp task firstprivate (band_11, band_12, band_21, band_22) lastprivate (ffd_bp)
    subMultSq (band_11, band_12, band_21, band_22, &ffd_bp)
  }
  #pragma omp task firstprivate (fm_qd_bp) lastprivate (band_2)
  FFD (8.407, fm_qd_bp, &band_2)
  #pragma omp task firstprivate (ffd_bp) lastprivate (band_3)
  FFD (8.407, ffd_bp, &band_3)
  #pragma omp task firstprivate (band_2, band_3)
  {
    stereo_sum (band_2, band_3, &output)
    Writer (output)
  }
}

```

Figure 27. OpenMP annotated FMradio

**Table 1.** Speedup results obtained from FMradio

Version	Only Stream.	Only Vect.	Stream. + Vect.
SPM	1.1	2	1.4
GCC-SPM	2.6	2	3.6

We evaluate the impact of streamization using the above two approaches on Power6. The sequential code is identical in both cases. The SPM streamized version produced only 1.1x speedup factor over the sequential code because of degraded cache behavior and synchronization overhead. Refer to Table 1 for a comparison of the performance achieved by each technique explained in this section.

However, in the case of the GCC-SPM version, these problems were fixed by an optimization (45) implemented in GCC’s OpenMP runtime library, `libGOMP`, whereby the allowed patterns of accesses to streams preclude false sharing of cache lines between producers and consumers. More specifically, this optimization consists of increasing the granularity of accesses to streams by aggregating the reading (resp. writing) of multiple elements in read (resp. write) windows. The size of such windows is a multiple of the size of a cache line, which ensures that producers and consumers never access simultaneously the same cache lines. Thanks to this optimization, this version achieves a 2.6x speedup factor over the sequential code.

## 8.2 Loop-level optimization and vectorization

FMradio does not exhibit enough nested loops to illustrate the need for complex loop transformations. Loop fusion is the only relevant one, but happens to be always compatible with vectorization in this example; it is systematically applied together with task-level fusion. The compiler (GCC) then proceeds to apply auto-vectorization, as one of the final Middle-End optimization passes. The main computation kernel in FMradio is an inner-loop that scans through the input buffer and computes the sum of products with the array of coefficients. The impact of vectorization is 1.3x/1.4x improvement factor over the streamized versions (as described in Section 8.1) using “SPM”/“GCC-SPM” respectively). Alignment handling and reduction overhead (to finalize the summation) are the main factors that explain the gap between the theoretical speedup factor from vectorization ( $4\times$ ) and the speedups we observed.

Relative to the sequential (non-streamized) version, the impact of vectorization is higher, achieving a  $2\times$  speedup. These speedups are summarized in Table 1.

## 8.3 Code generation

The streamized and vectorized code proceeds through the compilation flow, reaching the back-end target-dependent compilation passes, all the way through final code generation, using the machine description files (of Power6 in this example), where the Acolib/OpenMP constructs of the “SPM”/“GCC-SPM” are translated to pthreads library calls. The overall impact of streamizing and vectorizing FMradio is 1.4x/3.6x respectively.

## 9. Conclusions

In this paper we presented and demonstrated the framework developed by the ACOTES project. ACOTES includes partners from both industry and academia, whose goal is to improve programmer’s productivity using: (1) automatic simulation and compilation techniques to abstract the underlying multi-core hardware from the programmer, and (2) programmer hints (pragmas) that define the inputs, outputs and control variables of the computation, hinting to

the underlying compilation system where the borders of the components are. The actual components are then built based on an abstract representation of the platform called the Abstract Streaming Machine (ASM). The ASM expresses the processing thread-level and data-level parallelism capabilities available, and in addition communication overhead (processing and delay) between the processors. The automatic compiler transformations then base their parallelism related optimization decisions on the pragmas and the resources needed by each constructed component mapped to each processor. These techniques and tools were demonstrated using the FMradio streaming program, starting from its programming using the ACOTES pragmas, through its multiple levels of compilation, all the way to actual execution on a real streaming architecture.

## References

- [1] R. Allen and K. Kennedy. *Optimizing Compilers for Modern Architectures*. Morgan Kaufmann Publishers, 2001.
- [2] J. Balart, A. Duran, M. González, X. Martorell, E. Ayguadé, and J. Labarta. Nanos Mercurium: a Research Compiler for OpenMP. In *Proceedings of the European Workshop on OpenMP 2004*, October 2004.
- [3] P. Bellens, J.M. Perez, R.M. Badia, J. Labarta, B.S. Center, U.P.C.B.N. II, and J. Girona. CellSs: a Programming Model for the Cell BE Architecture. *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing*, 2006.
- [4] E. Blossom. GNU Radio: tools for exploring the radio frequency spectrum. *Linux Journal*, 2004(122), 2004.
- [5] François Bodin, Toru Kisuki, Peter M. W. Knijnenburg, Mike F. P. O’Boyle, and Erven Rohou. Iterative compilation in a non-linear optimisation space. In *Workshop on Profile and Feedback-Directed Compilation (FDO-1)*, in conjunction with PACT ’98, October 1998.
- [6] Uday Bondhugula, Albert Hartono, J. Ramanujam, and P. Sadayappan. A practical automatic polyhedral parallelization and locality optimization system. In *ACM Conf. on Programming Languages Design and Implementation (PLDI’08)*, Tucson, AZ, USA, June 2008.
- [7] Simone Campanoni, Giovanni Agosta, and Stefano Crespi Reghizzi. A parallel dynamic compiler for CIL bytecode. *SIGPLAN Notices*, 43(4):11–20, 2008.
- [8] Paul Carpenter, Alex Ramirez, Xavier Martorell, David Rodenas, and Roger Ferrer. Report on Streaming Programming Model and Abstract Streaming Machine Final Version. Deliverable D2.2, IST ACOTES Project, September 2008.
- [9] Paul Carpenter, David Rodenas, Xavier Martorell, Alex Ramirez, and Eduard Ayguade. A Streaming Machine Description and Programming Model. In *Lecture Notes in Computer Science, Embedded Computer Systems: Architectures, Modeling and Simulation, Springer Berlin/Heidelberg*, pages 107–116, August 2007.
- [10] Marco Cornero, Erven Rohou, Andrea Ornstein, and Razya Ladelsky. Report on Back-end Formats. Deliverable D5.3, IST ACOTES Project, December 2007.
- [11] Roberto Costa, Andrea C. Ornstein, and Erven Rohou. CLI back-end in GCC. In *GCC Developers’ Summit*, pages 111–116, July 2007.
- [12] D. G. Messerschmitt E. A. Lee. Static scheduling of synchronous data flow programs for digital signal processing. *IEEE Trans. Computers*, 36(1):24–25, 1987.
- [13] ECMA International, Rue du Rhône 114, 1204 Geneva, Switzerland. *Common Language Infrastructure (CLI) Partitions I to IV*, 4th edition, June 2006.
- [14] Ayguade Eduard, Nawal Copty, Alejandro Duran, Jay Hoeflinger, Yuan Lin, Massaioli Federico, Ernesto Su, Priya Unnikrishnan, and Zhang Guansong. A proposal for task parallelism in OpenMP. *Submitted to IWOMP2007*, 2007.
- [15] B. Flachs et al. A Streaming Processor Unit for a CELL Processor. *ISSCC Dig. Tech. Papers*, Paper 7.4:134–135, Feb 2005.

- [16] D. Pham et al. The Design and Implementation of a First-Generation CELL Processor. *ISSCC Dig. Tech. Papers*, Paper 10.2:184–185, Feb 2005.
- [17] M. Gschwind et al. An Open Source Environment for Cell Broadband Engine System Software. *IEEE Computer*, 40(6):37–47, June 2007.
- [18] Kayvon Fatahalian, Daniel Reiter Horn, Timothy J. Knightld, Larkhoon Leem, Mike Houston, Ji Young Park, Mattan Erez, Manman Ren, Alex Aiken, William J. Dally, and Pat Hanrahan. Sequoia: programming the memory hierarchy. In *ACM/IEEE conference on Supercomputing (SC'06)*, 2006.
- [19] P. Feautrier. Scalable and modular scheduling. In A. D. Pimentel and S. Vassiliadis, editors, *Computer Systems: Architectures, Modeling and Simulation (SAMOS'04)*, number 3133 in LNCS, pages 433–442. Springer-Verlag, 2004. An updated version will appear in *IJPP*.
- [20] G. Fursin and A. Cohen. Building a Practical Iterative Interactive Compiler. In *1st Workshop on Statistical and Machine Learning Approaches Applied to Architectures and Compilation (SMART'07), colocated with HiPEAC 2007 conference*, 2007.
- [21] A.S. Terechko G. Al-Kadi. A hardware task scheduler for embedded video processing. In *Proceedings of the 4th International Conference on High Performance and Embedded Architectures and Compilers, Paphos, Cyprus*, number 5409 in LNCS, pages 140–152, 2009.
- [22] S. Girbal, N. Vasilache, C. Bastoul, A. Cohen, D. Parelo, M. Sigler, and O. Temam. Semi-automatic composition of loop transformations for deep parallelism and memory hierarchies. *Intl. J. of Parallel Programming*, 34(3):261–317, June 2006. Special issue on Microgrids.
- [23] S. Girona, J. Labarta, and R.M. Badia. Validation of Dimemas communication model for MPI collective operations. *Proc. EuroPVM/MPI*, 2000.
- [24] M. Gonzalez, E. Ayguade, X. Martorell, and J. Labarta. Complex Pipelined Executions in OpenMP Parallel Applications. *International Conference on Parallel Processing (ICPP R2001)*, September 2001.
- [25] M. Gonzalez, E. Ayguade, X. Martorell, and J. Labarta. Exploiting pipelined executions in OpenMP. *Parallel Processing, 2003. Proceedings. 2003 International Conference on*, pages 153–160, 2003.
- [26] M.I. Gordon, W. Thies, and S. Amarasinghe. Exploiting coarse-grained task, data, and pipeline parallelism in stream programs. *Proc. ASPLOS 2006*, pages 151–162, 2006.
- [27] ILOG. <http://www.ilog.com/products/cplex/>. CPLEX Math Programming Engine.
- [28] Coding of Moving Pictures International Organization for Standardisation, ISO/IEC JTC1/SC29/WG11 and Audio. Overview of the MPEG-4 Standard, <http://www.chiairiglione.org/mpeg>.
- [29] A. Terechko J. Hoogerbrugge. A multithreaded multicore system for embedded media processing. *Transactions on High-Performance Embedded Architectures and Compilers*, 4(2), 2008.
- [30] M. Kudlur and S. Mahlke. Orchestrating the execution of stream programs on multicore platforms. *Proceedings of the 2008 ACM SIGPLAN conference on Programming language design and implementation*, pages 114–124, 2008.
- [31] Manjunath Kudlur and Scott Mahlke. Orchestrating the execution of stream programs on multicore platforms. In *ACM Conf. on Programming Languages Design and Implementation (PLDI'08)*, pages 114–124, June 2008.
- [32] Chris Lattner and Vikram Adve. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *Proceedings of the 2004 International Symposium on Code Generation and Optimization (CGO'04)*, Palo Alto, California, Mar 2004.
- [33] C. G. Lee. Utdsp benchmarks. <http://www.eecg.toronto.edu/corinna/DSP/infrastructure/UTDSP.html>, 1998.
- [34] W.I. Lundgren, K.B. Barnes, and J.W. Steed. Gedae: Auto Coding to a Virtual Machine.
- [35] D. Naishlos, M. Biberstein, S. Ben-David, and A. Zaks. Vectorizing for a simdd dsp architecture. In *CASES*, 2003.
- [36] Dorit Naishlos. Autovectorization in gcc. In *the GCC Developer's summit*, pages 105–118, June 2004.
- [37] V. Ngo. Parallel loop transformation techniques for vector-based multiprocessor systems. *Ph.D. thesis, U. of Minn.*, 1994.
- [38] Maik Nijhuis, Herbert Bos, Henri Bal, and Cédric Augonnet. Mapping and synchronizing streaming applications on cell processors. In *International Conference on High Performance Embedded Architectures and Compilers (HiPEAC'09)*, Paphos, Cyprus, January 2009.
- [39] Novell. The Mono Project. <http://www.mono-project.com>.
- [40] D. Nuzman, M. Namolaru, A. Zaks, and J. H. Derby. Compiling for an indirect vector register architecture. In *CF*, 2008.
- [41] D. Nuzman, I. Rosen, and A. Zaks. Auto-vectorization of interleaved data for SIMD. In *PLDI*, 2006.
- [42] Dorit Nuzman and Ayal Zaks. Outer-loop vectorization - revisited for short SIMD architectures. In *PACT*, October 2008.
- [43] OpenMP Organization. OpenMP Application Program Interface, v. 3.0, May 2008.
- [44] Antoniu Pop and Sebastian Pop. A Proposal for lastprivate Clause on OpenMP task Pragma. Technical report, MINES ParisTech, CRI - Centre de Recherche en Informatique, Mathématiques et Systèmes, 35 rue St Honoré 77305 Fontainebleau-Cedex, France, January 2009.
- [45] Antoniu Pop, Sebastian Pop, Harsha Jagasia, Jan Sjödin, and Paul H J Kelly. Improving GNU compiler collection infrastructure for streamization. In *Proceedings of the 2008 GCC Developers' Summit*, pages 77–86, 2008.
- [46] Louis-Noel Pouchet, Cedric Bastoul, Albert Cohen, and John Cavazos. Iterative optimization in the polyhedral model: Part II, multidimensional time. In *ACM Conf. on Programming Language Design and Implementation (PLDI'08)*, Tucson, Arizona, June 2008.
- [47] Ira Rosen, Ben Elliston, Revital Eres, Alan Modra, Dorit Nuzman, Ulrich Weigand, Ayal Zaks, and David Edelsohn. Compiling Effectively for Cell B.E. with GCC. In *14th Workshop on Compilers for Parallel Computing (CPC)*, January 2009.
- [48] R. G. Scarborough and H. G. Kolsky. A vectorizing fortran compiler. *IBM J. of R & D*, 30(2):163–171, March 1986.
- [49] Southern Storm Software, Pty Ltd. DotGNU project. <http://dotgnu.org>.
- [50] StreamIt Language Specification Version 2.1, 2006.
- [51] Gabriele Svelto, Andrea Ornstein, and Erven Rohou. A stack-based internal representation for GCC. In *First International Workshop on GCC Research Opportunities (GROW09), in conjunction with HiPEAC 2009*, pages 37–48, January 2009.
- [52] U. Weigand. Porting the GNU Tool Chain to the Cell Architecture. *Proceedings of the 2005 GCC Summit, Ottawa, Canada*, June 2005.
- [53] Michael Wolfe. *High Performance Compilers for Parallel Computing*. Addison Wesley, 1996.